
dpEmu

Release 0.1.0

dpEmu Team

Feb 07, 2020

CONTENTS

1	Introduction	1
2	Getting Started	3
3	User Manual	21
4	Case studies	33
5	Module documentation	207

INTRODUCTION

dpEmu is a Python library for emulating data problems in the use and training of machine learning systems. dpEmu can help you to:

- Generate errors in your training and/or testing data in a controlled and documentable manner.
- Run one or more machine learning models on your data using different values for the error parameters.
- Visualize the results.

GETTING STARTED

2.1 Tutorial I: Error Generation Basics

Let's start by importing some of the packages, classes and functions that we will use in the tutorial.

```
[1]: import matplotlib.pyplot as plt
import numpy as np

from dpemu.dataset_utils import load_mnist
from dpemu.filters.common import GaussianNoise, Missing
from dpemu.nodes import Array, Series
```

In this tutorial we will be using the famous MNIST dataset of handwritten digits. dpEmu provides a convenience function for downloading the dataset. The dataset is split into training and testing data, but we can ignore the testing data for now and only keep the training data:

```
[2]: x, y, _, _ = load_mnist()
```

It's a good idea to start by exploring the shape and the data type as well as the minimum and maximum values of the input data:

```
[3]: print(f"shape: {x.shape}")
print(f"dtype: {x.dtype}")
print(f"min: {x.min()}, max: {x.max()}")

shape: (60000, 784)
dtype: float64
min: 0.0, max: 255.0
```

Now let's do the same for the output data:

```
[4]: print(f"shape: {y.shape}")
print(f"dtype: {y.dtype}")
print(f"min: {y.min()}, max: {y.max()}")

shape: (60000,)
dtype: float64
min: 0.0, max: 9.0
```

It looks like our input consists of 60000 rows where each row is a 784 pixel (i.e. 28×28) black and white image of a handwritten digit. The output corresponding to each row is its correct label.

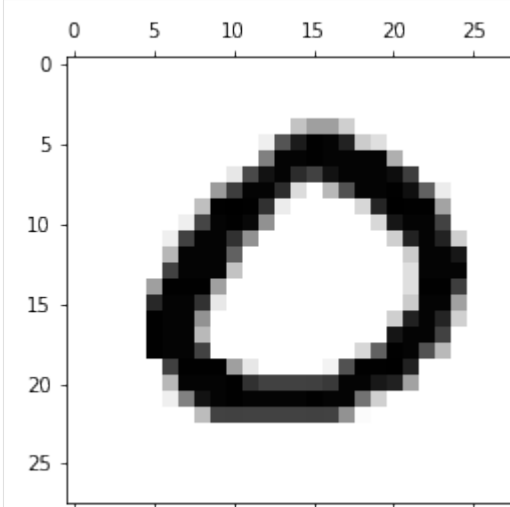
We could work with the whole dataset, but for the purposes of this tutorial a small subset will suffice:

```
[5]: n = 1000  
xs = x[:n]  
ys = y[:n]
```

Now let's pick a data point at random and display the image and its label.

```
[6]: ind = np.random.randint(n)  
plt.matshow(xs[ind].reshape((28, 28)), cmap='gray_r')  
print(f"The label of the image at index {ind} is {ys[ind]}")
```

The label of the image at index 435 is 0.0.



Now that we know our data (superficially at least) we can start adding errors. First we must model the shape of the data as a tree. If that sounds complicated, don't worry – it's ridiculously easy!

Since the inputs are an indexed collection of images, it's natural to represent them as a series of arrays, each array corresponding to a single image. Let's do just that:

```
[7]: image_node = Array()  
series_node = Series(image_node)
```

The Series node is the root of the tree, and the Array node is its only child.

We can now add one or more error sources. Error sources are known as Filters in dpEmu parlance, and they can be attached to Array nodes (and indeed some other kinds of nodes which we will not discuss in this tutorial).

```
[8]: gaussian_noise_source = GaussianNoise("mean", "std")  
image_node.addfilter(gaussian_noise_source)
```

The GaussianNoise Filter does exactly what it sounds like: it adds noise drawn from a Normal distribution. The constructor takes two String arguments – namely, identifiers for the parameters (the mean and the standard deviation) of the distribution. We will provide the values of these parameters when we want to generate the errors.

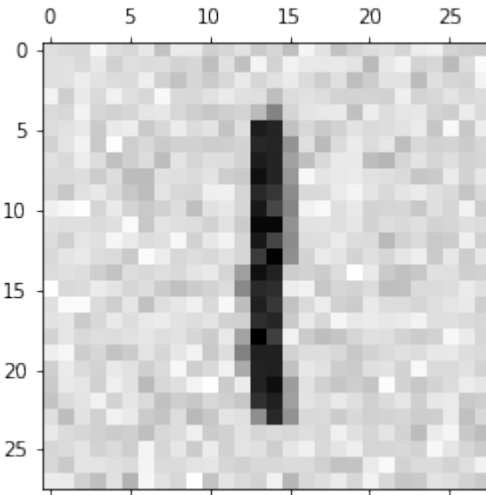
Now let's try applying our error generating tree!

```
[9]: params = {"mean": 0.0, "std": 20.0}  
errorified = series_node.generate_error(xs, params)
```

It really is that easy! Now let's plot a random image from our data subset into which errors have been introduced:


```
[10]: ind = np.random.randint(n)
      plt.matshow(errorified[ind].reshape((28, 28)), cmap='gray_r')
```

```
[10]: <matplotlib.image.AxesImage at 0x7fd6ce330978>
```



We are not limited to one error source (i.e. Filter) per node. Let's add another one:

```
[11]: image_node.addfilter(Missing("probability", "missing_value"))
```

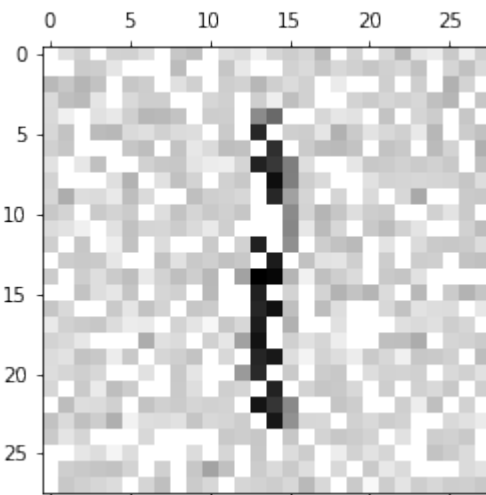
The Missing Filter takes each value in the array and changes it to NaN with the probability of our choice.

Now let's apply the modified error generating tree to the same subset of data:

```
[12]: params = {"mean": 0.0, "std": 20.0, "probability": .3, "missing_value": np.nan}
      errorified = series_node.generate_error(xs, params)
```

```
      plt.matshow(errorified[ind].reshape((28, 28)), cmap='gray_r')
```

```
[12]: <matplotlib.image.AxesImage at 0x7fd6ce29d6a0>
```



Congratulations! This concludes the first tutorial. There is much more to explore, but you now know enough to get started. We hope you enjoy using dpEmu!

The notebook for this tutorial can be found [here](#).

2.2 Tutorial II: Custom Filters

The dpEmu package provides many filters for common use cases, but it might occur that no built-in filter exists for what you want. To solve this problem, you can create your own filters by inheriting the `Filter`-template. In this tutorial, we will create a custom filter for converting input images from RGB to grayscale.

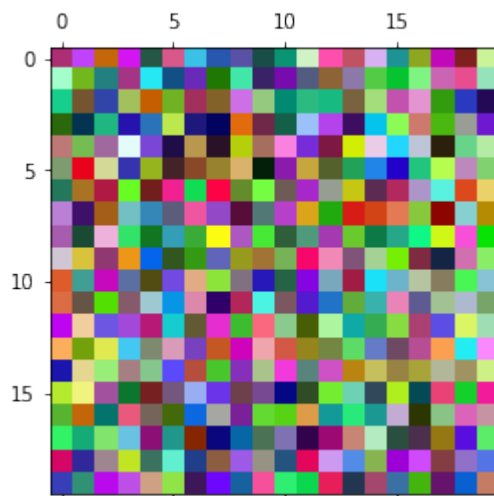
First, let's generate some input data. We'll use numpy's `randomstate` with a set seed to ensure repeatability. The `randomstate` object simply replaces the `np.random`-part of any function call you want to make to numpy's random module.

```
[1]: import matplotlib.pyplot as plt
import numpy as np

rand = np.random.RandomState(seed=0)
img = rand.randint(low=0, high=255, size=(20, 20, 3))
data = np.array([img for _ in range(9)])
```

Let's see what these random images look like:

```
[2]: plt.matshow(data[0])
[2]: <matplotlib.image.AxesImage at 0x7fa3a2a4b0d0>
```



Before we write our own filter, let's rehearse what we learned in the previous tutorial and try modifying the image with an existing filter. We'll use the Resolution filter here.

```
[3]: from dpemu.nodes import Array, Series
from dpemu.filters.image import Resolution

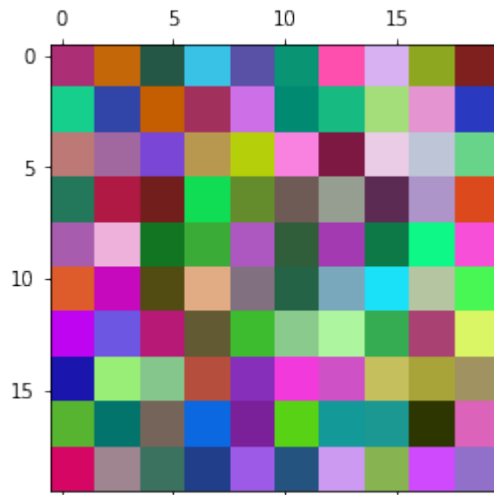
image_node = Array()
series_node = Series(image_node)
resolution_filter = Resolution("scale")
image_node.addfilter(resolution_filter)

params = {"scale" : 2}
errorified = series_node.generate_error(data, params)
```

Let's quickly check that the output is what we'd expect:

```
[4]: plt.matshow(errorified[0])
```

```
[4]: <matplotlib.image.AxesImage at 0x7fa3d6ff9590>
```



Now we're ready to write our own filter to replace the built-in resolution filter. To do this, we want to inherit the `Filter`-class in `dpemu.filters.filter`. When we inherit the class, we'll want to define a constructor and override the `apply`-function which applies the filter to the input data. The first parameter of the function is the input data. In this case we do not need the other two parameters, but we'll learn how to use them later on in this tutorial.

```
[6]: from dpemu.filters.filter import Filter
```

```
class Grayscale(Filter):
    def __init__(self):
        super().__init__()
    def apply(self, node_data, random_state, named_dims):
        avg = np.sum(node_data, axis=-1) // 3
        for ci in range(3):
            node_data[:, :, ci] = avg
```

The code here is very simple. In the `apply`-function, we just take the mean of the color channels, and then assign that to each of them. Note that Filters must always maintain the dimensions of the input data. That is why the processed image will still contain 3 color channels, despite being in grayscale.

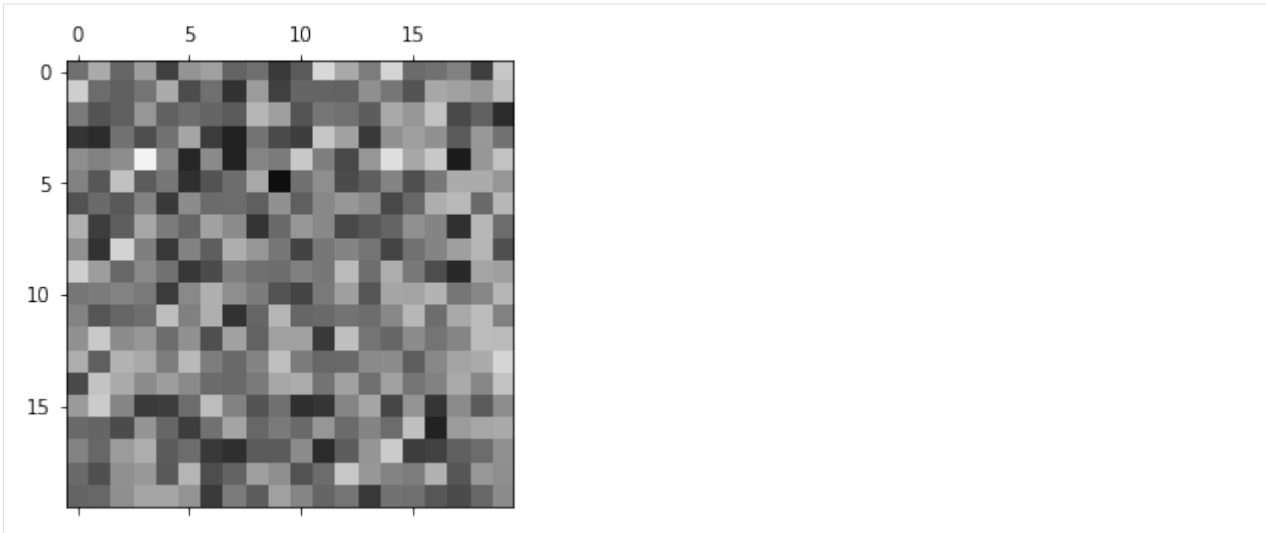
Let's try our filter. Unlike before, we do not pass any parameters when calling `generate_error`, since our filter doesn't take any parameters.

```
[8]: image_node = Array()
series_node = Series(image_node)
grayscale_filter = Grayscale()
image_node.addfilter(grayscale_filter)

params = {}
errorified = series_node.generate_error(data, params)

plt.matshow(errorified[0])
```

```
[8]: <matplotlib.image.AxesImage at 0x7fa397cff7d0>
```



Now that we understand how a simple filter works, let's incorporate randomness and parameters. We'll change the filter so that every pixel gets replaced with its grayscale equivalent with some probability given as a parameter:

```
[7]: class RandomGrayscale(Filter):
    def __init__(self, probability_id):
        super().__init__()
        self.probability_id = probability_id
    def apply(self, node_data, random_state, named_dims):
        inds = random_state.rand(node_data.shape[0], node_data.shape[1]) < self.
        ↪probability
        avg = np.sum(node_data[inds], axis=-1) // 3
        for ci in range(3):
            node_data[inds, ci] = avg
```

If you read the code carefully, you might notice that `self.probability` is not defined anywhere. For convenience, for every variable ending in `_id`, a value will be assigned to the variable without the `_id`-suffix from the `params-list` passed to `generate_error`. For example, if we have variable `probability_id`, which is set to "probability" in the initializer, and we call `generate_error` with a dictionary containing the key-value pair "probability" : 0.5, the value of `self.probability` will be 0.5.

In function `__init__`, we take the identifier for our probability parameter. This will be used as the key to find the value of `self.probability` from the `params-dictionary`.

In `apply`, we randomize the positions where the pixel will be replaced by its grayscale equivalent, then replace them as we did previously. For repeatability, we use the numpy `RandomState` passed to the function as the second parameter.

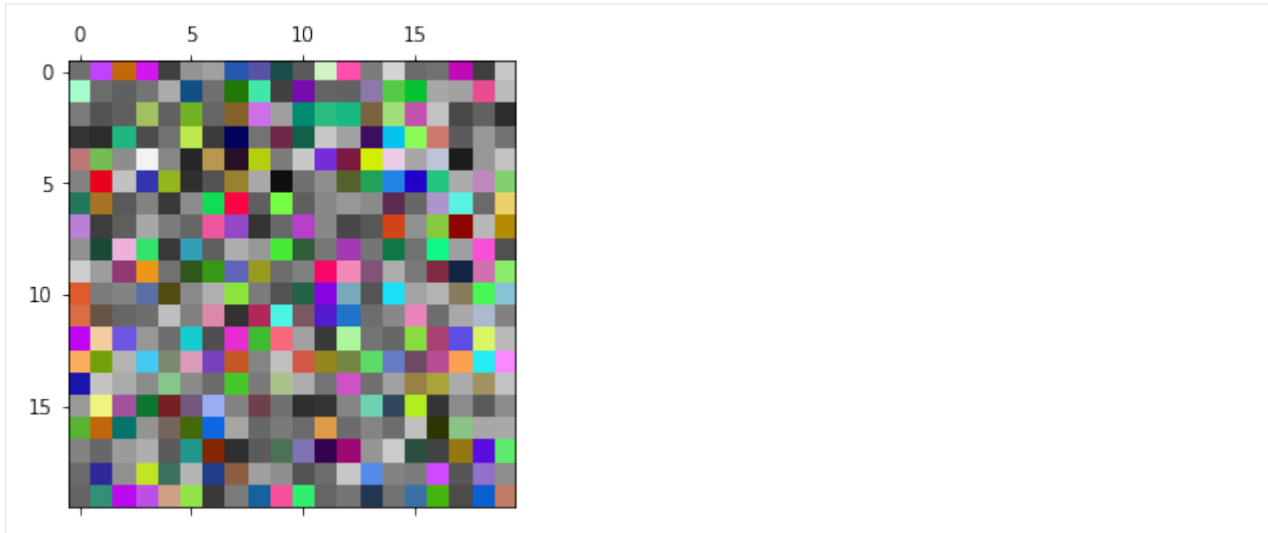
Let's inspect the results:

```
[8]: image_node = Array()
    series_node = Series(image_node)
    grayscale_filter = RandomGrayscale("probability")
    image_node.addfilter(grayscale_filter)

    params = {"probability" : 0.5}
    errorified = series_node.generate_error(data, params)

    plt.matshow(errorified[0])

[8]: <matplotlib.image.AxesImage at 0x7f9c692bb160>
```



Now, the only thing we are yet to understand is the `named_dims` parameter to apply. For this, imagine our tuple of ten images is a video. We'll create a filter that makes the pixels in the video decay into grayscale over time. Once a pixel turns grayscale, it won't turn back, and by expectation half of the pixels will decay in the amount of frames given as a parameter.

Most of the new code is just randomizing the times at which individual pixels decay. Note that we make a copy of the `random_state` to ensure that we generate the same times for every image in the series. This is to ensure pixels don't gain back their colors after decaying.

```
[10]: from math import log
      from copy import deepcopy

      class DecayGrayscale(Filter):
          def __init__(self, half_time_id):
              super().__init__()
              self.half_time_id = half_time_id

          def apply(self, node_data, random_state, named_dims):
              shape = (node_data.shape[0], node_data.shape[1])
              times = deepcopy(random_state).exponential(scale=(self.half_time / log(2)),
              ↪size=shape)
              inds = times <= named_dims["time"]
              avg = np.sum(node_data[inds], axis=-1) // 3
              for ci in range(3):
                  node_data[inds, ci] = avg
```

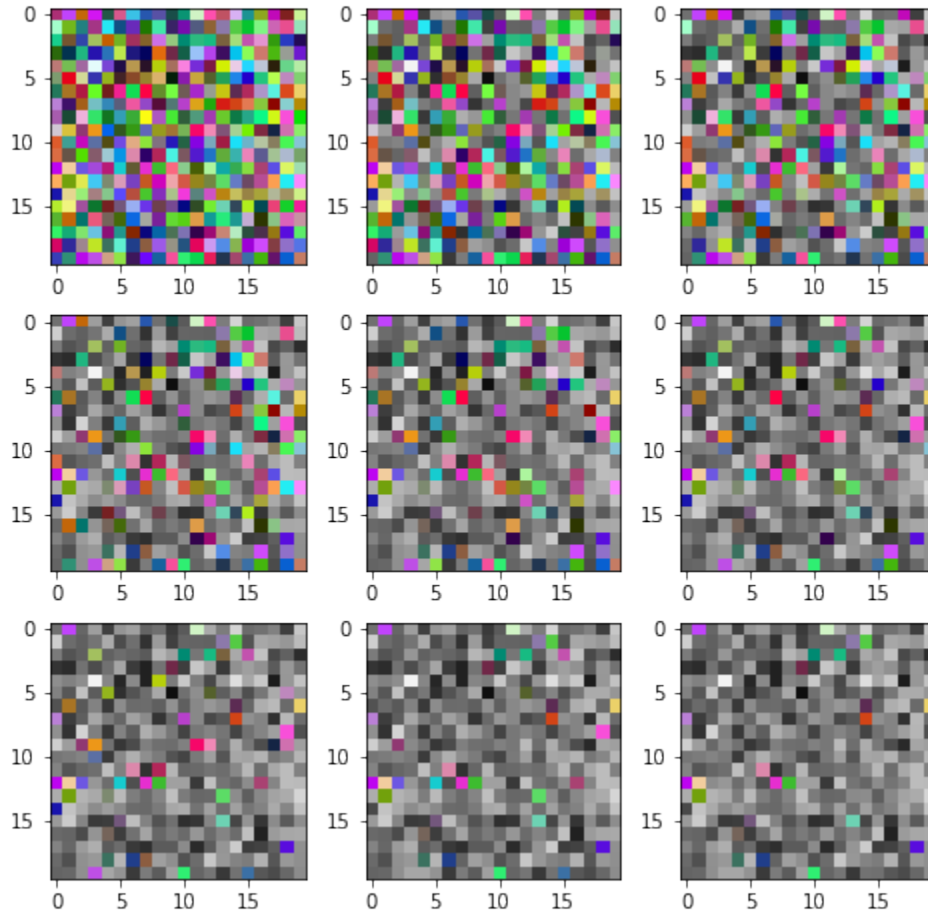
We use `named_dims` on the line `inds = times <= named_dims["time"]`. This line creates a mask of all pixels that decay before the current time, which is given by `named_dims["time"]`. To use this filter, we'll have to tell the series-node that the dimension it is iterating over is called "time":

```
[13]: image_node = Array()
      series_node = Series(image_node, "time")
      grayscale_filter = DecayGrayscale("half_time")
      image_node.addfilter(grayscale_filter)

      params = {"half_time" : 2}
      errorified = series_node.generate_error(data, params)
```

Finally, to show multiple images in the same plot, we'll have to do some more work:

```
[14]: fig = plt.figure(figsize=(8, 8))
      for i in range(1, 9 + 1):
          fig.add_subplot(3, 3, i)
          plt.imshow(errorified[i-1])
      plt.show()
```



And we have achieved the desired effect! This concludes the second tutorial.

The notebook for this tutorial can be found [here](#).

2.3 Tutorial III: Using Runner Advanced I

The purpose of this tutorial is to learn more about combining error generation with our Runner and using it to compare different ML models' performance with errorified data. In this tutorial, instead of MNIST dataset, we'll be using smaller 8x8 images of handwritten digits for performance reasons. The usecase is comparing some clustering algorithms using data with different amount of error.

First we have to import the required modules. Let's also disable some annoying warnings.

```
[1]: import warnings
      from abc import ABC, abstractmethod

      import matplotlib.pyplot as plt
```

(continues on next page)

(continued from previous page)

```

import numpy as np
from hdbscan import HDBSCAN
from numba.errors import NumbaDeprecationWarning, NumbaWarning
from numpy.random import RandomState
from sklearn.cluster import KMeans
from sklearn.metrics import adjusted_rand_score, adjusted_mutual_info_score

from dpemu import runner
from dpemu.dataset_utils import load_digits_
from dpemu.filters.common import Missing
from dpemu.ml_utils import reduce_dimensions
from dpemu.nodes import Array
from dpemu.plotting_utils import visualize_scores, visualize_classes, \
    print_results_by_model

warnings.simplefilter("ignore", category=NumbaDeprecationWarning)
warnings.simplefilter("ignore", category=NumbaWarning)

/home/thalvari/PycharmProjects/dpEmu/venv/lib/python3.7/site-packages/sklearn/
↳externals/six.py:31: DeprecationWarning: The module is deprecated in version 0.21
↳and will be removed in version 0.23 since we've dropped support for Python 2.7.
↳Please rely on the official version of six (https://pypi.org/project/six/).
    "(https://pypi.org/project/six/).", DeprecationWarning)
/home/thalvari/PycharmProjects/dpEmu/venv/lib/python3.7/site-packages/sklearn/
↳externals/joblib/__init__.py:15: DeprecationWarning: sklearn.externals.joblib is
↳deprecated in 0.21 and will be removed in 0.23. Please import this functionality
↳directly from joblib, which can be installed with: pip install joblib. If this
↳warning is raised when loading pickled models, you may need to re-serialize those
↳models with scikit-learn 0.21+.
    warnings.warn(msg, category=DeprecationWarning)

```

The first phase of our pipeline is loading the dataset. The smaller Digits dataset contains 1797 images.

```

[2]: def get_data():
    return load_digits_()

```

Next we have to define our root node for error generation. Now the data is in one huge Numpy array and each row contains all pixel values for one image. Thus the root node's type needs to be Array. As for filters, we'll be using the Missing filter from the previous tutorial.

```

[3]: def get_err_root_node():
    err_root_node = Array()
    err_root_node.addfilter(Missing("probability", "missing_value_id"))
    return err_root_node

```

In this step we are defining the different error parameters for our filter we want to test the ML models with.

```

[4]: def get_err_params_list():
    p_steps = np.linspace(0, .5, num=6)
    err_params_list = [{"probability": p, "missing_value_id": 0} for p in p_steps]
    return err_params_list

```

The preprocessor code is run only once after the errorified data is generated, but before any of the ML models. The purpose of a preprocessor is to modify the errorified data to a usable format for the ML models. In this example we are using our premade pipeline for dimensionality reduction. Note that every dictionary element you add to the last return parameter will create a new column in the resulting Dataframe returned by the Runner. In this tutorial we will be using the reduced errorified data in our visualizations.

```
[5]: class Preprocessor:
    def __init__(self):
        self.random_state = RandomState(42)

    def run(self, _, data, params):
        reduced_data = reduce_dimensions(data, self.random_state)
        return None, reduced_data, {"reduced_data": reduced_data}
```

Next we have to define the ML models. Every model has to have a public run method, which gets the preprocessed data. Similarly to the Preprocessor, the run method here also returns a dictionary and every added key will create a new column in the resulting Dataframe returned by the Runner. In this tutorial we'll be testing two different clustering algorithms, the other with two different parameters, and use AMI and ARI scores for comparison. The labels are only used for calculating the scores and unlike HDBSCAN, KMeans needs to know the desired number of clusters. HDBSCAN's min_cluster_size give the smallest size of a group of datapoints that can be considered a cluster.

```
[6]: class AbstractModel(ABC):

    def __init__(self):
        self.random_state = RandomState(42)

    @abstractmethod
    def get_fitted_model(self, data, params):
        pass

    def run(self, _, data, params):
        labels = params["labels"]
        fitted_model = self.get_fitted_model(data, params)
        return {
            "AMI": round(adjusted_mutual_info_score(labels, fitted_model.labels_,
↪average_method="arithmetic"), 3),
            "ARI": round(adjusted_rand_score(labels, fitted_model.labels_), 3),
        }

class KMeansModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        labels = params["labels"]
        n_classes = len(np.unique(labels))
        return KMeans(n_clusters=n_classes, random_state=self.random_state).fit(data)

class HDBSCANModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        return HDBSCAN(
            min_samples=params["min_samples"],
            min_cluster_size=params["min_cluster_size"],
            core_dist_n_jobs=1
        ).fit(data)
```

In this step we are defining the hyperparameters used by our models. This function returns a list of model/params_list

pairs. Our Runner runs every model in this list with all different sets of hyperparameters listed in the corresponding `params_list` -element. The scores for each model/params pair will become rows in the resulting Dataframe.

```
[7]: def get_model_params_dict_list(labels):
    return [
        {"model": KMeansModel, "params_list": [{"labels": labels}]},
        {"model": HDBSCANModel, "params_list": [{"min_cluster_size": 25, "min_samples
↪": 1, "labels": labels}]},
        {"model": HDBSCANModel, "params_list": [{"min_cluster_size": 50, "min_samples
↪": 1, "labels": labels}]},
    ]
```

Next we have to choose the visualizations we want to perform on our data and results. Certainly we would like to visualize the scores for each model given the amount of error. We could also just visualize the dataset in 2D using the `reduced_data` column we added to the Dataframe.

```
[8]: def visualize(df, label_names, dataset_name, data):
    visualize_scores(
        df,
        score_names=["AMI", "ARI"],
        is_higher_score_better=[True, True],
        err_param_name="probability",
        title=f"{dataset_name} clustering scores with missing pixels",
    )
    visualize_classes(
        df,
        label_names,
        err_param_name="probability",
        reduced_data_column="reduced_data",
        labels_column="labels",
        cmap="tab10",
        title=f"{dataset_name} (n={data.shape[0]}) true classes with missing pixels"
    )
    plt.show()
```

Now basically all that is left is asking the Runner to run our pipeline. Since we are doing unsupervised learning, training data won't be needed. By default Runner creates a subprocess for each element in `err_params_list`.

```
[9]: def main():
    data, labels, label_names, dataset_name = get_data()

    df = runner.run(
        train_data=None,
        test_data=data,
        preproc=Preprocessor,
        preproc_params=None,
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(labels),
    )

    print_results_by_model(df, ["missing_value_id", "labels", "reduced_data"])
    visualize(df, label_names, dataset_name, data)
```

Let's check out the results. HDBSCAN seems to perform better with little error and KMeans the other way around. Also HDBSCAN seems to perform better on average if `min_cluster_size` isn't too small.

[10]: main()

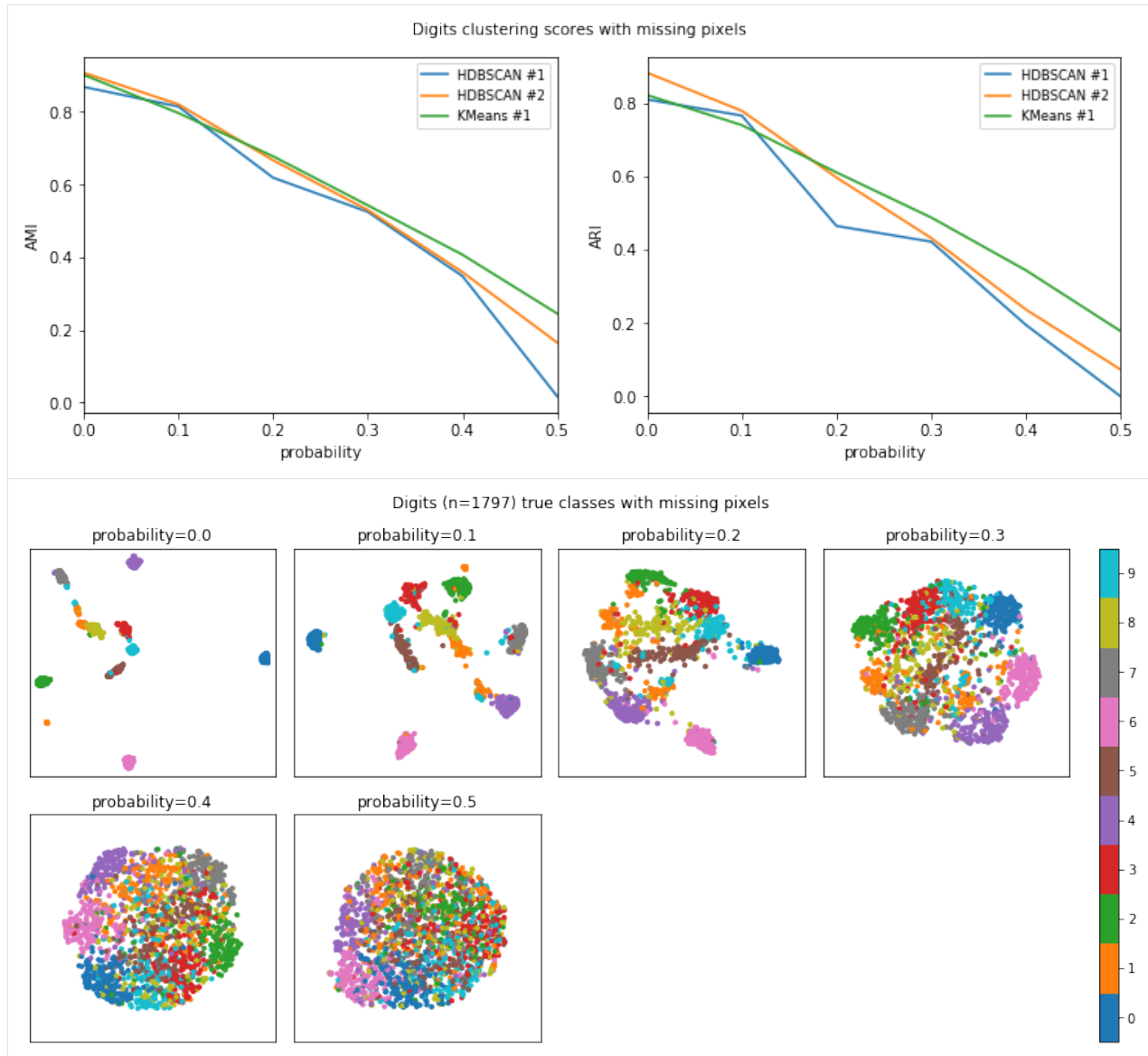
100%|| 6/6 [01:00<00:00, 11.93s/it]

HDBSCAN #1

HDBSCAN #2

KMeans #1

	AMI	ARI	probability	min_cluster_size	min_samples	time_err	time_pre	time_
↪mod								
0	0.869	0.810	0.0	25.0	1.0	0.033	38.741	0.
↪071								
1	0.815	0.766	0.1	25.0	1.0	0.029	40.570	0.
↪080								
2	0.619	0.465	0.2	25.0	1.0	0.041	42.267	0.
↪069								
3	0.525	0.422	0.3	25.0	1.0	0.069	41.331	0.
↪112								
4	0.348	0.195	0.4	25.0	1.0	0.006	19.588	0.
↪047								
5	0.017	0.000	0.5	25.0	1.0	0.008	19.054	0.
↪041								
	AMI	ARI	probability	min_cluster_size	min_samples	time_err	time_pre	time_
↪mod								
0	0.908	0.883	0.0	50.0	1.0	0.033	38.741	0.
↪063								
1	0.821	0.779	0.1	50.0	1.0	0.029	40.570	0.
↪074								
2	0.667	0.597	0.2	50.0	1.0	0.041	42.267	0.
↪064								
3	0.530	0.432	0.3	50.0	1.0	0.069	41.331	0.
↪110								
4	0.359	0.237	0.4	50.0	1.0	0.006	19.588	0.
↪041								
5	0.165	0.073	0.5	50.0	1.0	0.008	19.054	0.
↪041								
	AMI	ARI	probability	time_err	time_pre	time_mod		
0	0.902	0.822	0.0	0.033	38.741	0.112		
1	0.797	0.740	0.1	0.029	40.570	0.153		
2	0.678	0.611	0.2	0.041	42.267	0.212		
3	0.543	0.488	0.3	0.069	41.331	0.457		
4	0.407	0.344	0.4	0.006	19.588	0.213		
5	0.245	0.178	0.5	0.008	19.054	0.209		



The notebook for this tutorial can be found [here](#).

2.4 Tutorial IV: Using Runner Advanced II

The purpose of this tutorial is to learn more about Runner's advanced features and advanced visualization options. The usecase and the data are the same.

Using code from the previous tutorial:

```
[1]: %%capture --no-stdout

import warnings
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
```

(continues on next page)

(continued from previous page)

```

import numpy as np
from hdbscan import HDBSCAN
from numba.errors import NumbaWarning
from numpy.random import RandomState
from sklearn.cluster import KMeans
from sklearn.metrics import adjusted_rand_score, adjusted_mutual_info_score

from dpemu import runner
from dpemu.dataset_utils import load_digits_
from dpemu.filters.common import Missing
from dpemu.ml_utils import reduce_dimensions
from dpemu.nodes import Array
from dpemu.plotting_utils import visualize_best_model_params, visualize_scores, print_
    ↪ results_by_model

warnings.simplefilter("ignore", category=NumbaWarning)

def get_data():
    return load_digits_()

def get_err_root_node():
    err_root_node = Array()
    err_root_node.addfilter(Missing("probability", "missing_value"))
    return err_root_node

def get_err_params_list():
    p_steps = np.linspace(0, .5, num=6)
    err_params_list = [{"probability": p, "missing_value": 0} for p in p_steps]
    return err_params_list

class Preprocessor:
    def __init__(self):
        self.random_state = RandomState(42)

    def run(self, _, data, params):
        reduced_data = reduce_dimensions(data, self.random_state)
        return None, reduced_data, {"reduced_data": reduced_data}

class AbstractModel(ABC):

    def __init__(self):
        self.random_state = RandomState(42)

    @abstractmethod
    def get_fitted_model(self, data, params):
        pass

    def run(self, _, data, params):
        labels = params["labels"]
        fitted_model = self.get_fitted_model(data, params)
        return {
            "AMI": round(adjusted_mutual_info_score(labels, fitted_model.labels_,
    ↪ average_method="arithmetic"), 3),

```

(continues on next page)

(continued from previous page)

```

        "ARI": round(adjusted_rand_score(labels, fitted_model.labels_), 3),
    }

class KMeansModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        labels = params["labels"]
        n_classes = len(np.unique(labels))
        return KMeans(n_clusters=n_classes, random_state=self.random_state).fit(data)

class HDBSCANModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        return HDBSCAN(
            min_samples=params["min_samples"],
            min_cluster_size=params["min_cluster_size"],
            core_dist_n_jobs=1
        ).fit(data)

def main():
    data, labels, label_names, dataset_name = get_data()

    df = runner.run(
        train_data=None,
        test_data=data,
        preproc=Preprocessor,
        preproc_params=None,
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(labels),
    )

    print_results_by_model(df, ["missing_value_id", "labels", "reduced_data"])
    visualize(df, label_names, dataset_name, data)

```

Let's redo the step where we defined the hyperparameters used by our models. In the previous tutorial, we learned that Runner runs every model in the list defined by this function with all different sets of hyperparameters listed in the corresponding `params_list` -element. Now if we add more sets of hyperparameters to our only HDBSCAN's `param_list`, all these results will be listed under "HDBSCAN #1" in the resulting Dataframe. Using this information some of our visualizers are able to visualize hyperparameter-optimized results. Now we are also testing few different values for HDBSCAN's `min_samples`. Larger `min_samples` just means that more datapoints will be seen as noise.

```

[2]: def get_model_params_dict_list(labels):
    min_cluster_size_steps = [25, 50, 75]
    min_samples_steps = [1, 10]
    return [
        {"model": KMeansModel, "params_list": [{"labels": labels}]},

```

(continues on next page)

(continued from previous page)

```

    {"model": HDBSCANModel, "params_list": [{
        "min_cluster_size": min_cluster_size,
        "min_samples": min_samples,
        "labels": labels
    } for min_cluster_size in min_cluster_size_steps for min_samples in min_
    ↪ samples_steps]],
    ]

```

Let's also partially redo our visualizations. First we would like to visualize the hyperparameter-optimized scores for each of the models. Secondly we would like to see the best hyperparameters for HDBSCAN given the error.

```

[3]: def visualize(df, label_names, dataset_name, data):
    visualize_scores(
        df,
        score_names=["AMI", "ARI"],
        is_higher_score_better=[True, True],
        err_param_name="probability",
        title=f"{dataset_name} clustering scores with missing pixels",
    )
    visualize_best_model_params(
        df,
        model_name="HDBSCAN",
        model_params=["min_cluster_size", "min_samples"],
        score_names=["AMI", "ARI"],
        is_higher_score_better=[True, True],
        err_param_name="probability",
        title=f"Best parameters for {dataset_name} clustering"
    )
    plt.show()

```

Let's check out the results. Scores for HDBSCAN seem slightly better and smaller min_samples seems to be a better fit for data with lots of error.

```

[4]: main()
100%|| 6/6 [00:54<00:00, 10.64s/it]
HDBSCAN #1
KMeans #1

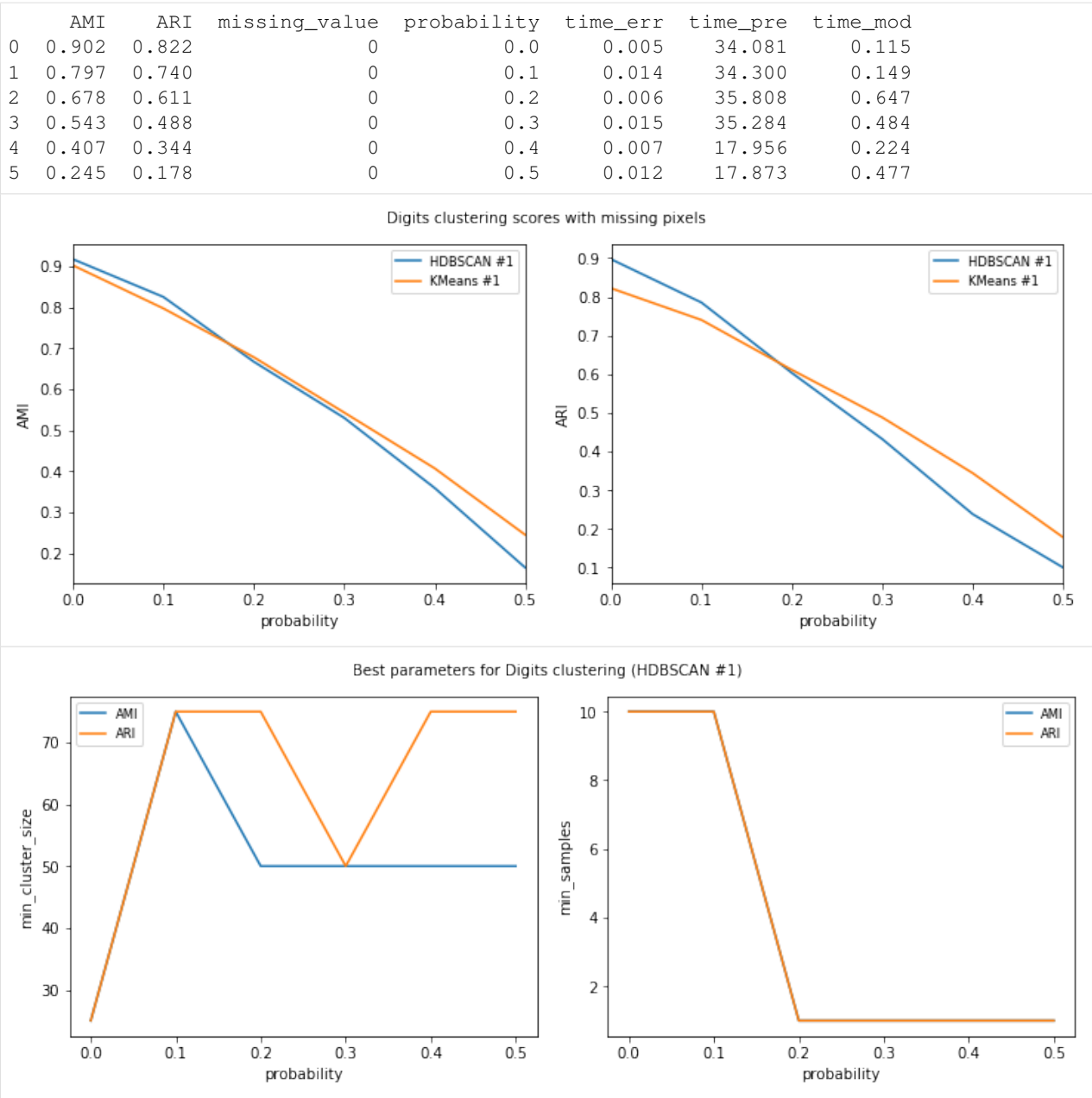
```

	AMI	ARI	missing_value	probability	min_cluster_size	min_samples	time_err_
↪ time_pre		time_mod					
0	0.869	0.810	0	0.0	25.0	1.0	0.005_
↪	34.081	0.072					
1	0.917	0.897	0	0.0	25.0	10.0	0.005_
↪	34.081	0.079					
2	0.908	0.883	0	0.0	50.0	1.0	0.005_
↪	34.081	0.069					
3	0.907	0.883	0	0.0	50.0	10.0	0.005_
↪	34.081	0.075					
4	0.908	0.883	0	0.0	75.0	1.0	0.005_
↪	34.081	0.067					
5	0.907	0.883	0	0.0	75.0	10.0	0.005_
↪	34.081	0.075					
6	0.815	0.766	0	0.1	25.0	1.0	0.014_
↪	34.300	0.074					
7	0.816	0.763	0	0.1	25.0	10.0	0.014_
↪	34.300	0.074					

(continues on next page)

(continued from previous page)

8	0.821	0.779	0	0.1	50.0	1.0	0.014	↵
↵	34.300	0.310						
9	0.807	0.751	0	0.1	50.0	10.0	0.014	↵
↵	34.300	0.214						
10	0.815	0.755	0	0.1	75.0	1.0	0.014	↵
↵	34.300	0.160						
11	0.825	0.785	0	0.1	75.0	10.0	0.014	↵
↵	34.300	0.169						
12	0.619	0.465	0	0.2	25.0	1.0	0.006	↵
↵	35.808	0.092						
13	0.630	0.518	0	0.2	25.0	10.0	0.006	↵
↵	35.808	0.080						
14	0.667	0.597	0	0.2	50.0	1.0	0.006	↵
↵	35.808	0.085						
15	0.634	0.524	0	0.2	50.0	10.0	0.006	↵
↵	35.808	0.130						
16	0.667	0.603	0	0.2	75.0	1.0	0.006	↵
↵	35.808	0.088						
17	0.631	0.509	0	0.2	75.0	10.0	0.006	↵
↵	35.808	0.076						
18	0.525	0.422	0	0.3	25.0	1.0	0.015	↵
↵	35.284	0.237						
19	0.526	0.389	0	0.3	25.0	10.0	0.015	↵
↵	35.284	0.276						
20	0.530	0.432	0	0.3	50.0	1.0	0.015	↵
↵	35.284	0.164						
21	0.524	0.390	0	0.3	50.0	10.0	0.015	↵
↵	35.284	0.084						
22	0.530	0.432	0	0.3	75.0	1.0	0.015	↵
↵	35.284	0.145						
23	0.524	0.390	0	0.3	75.0	10.0	0.015	↵
↵	35.284	0.180						
24	0.348	0.195	0	0.4	25.0	1.0	0.007	↵
↵	17.956	0.048						
25	0.196	0.061	0	0.4	25.0	10.0	0.007	↵
↵	17.956	0.052						
26	0.359	0.237	0	0.4	50.0	1.0	0.007	↵
↵	17.956	0.043						
27	0.196	0.061	0	0.4	50.0	10.0	0.007	↵
↵	17.956	0.054						
28	0.359	0.238	0	0.4	75.0	1.0	0.007	↵
↵	17.956	0.042						
29	0.196	0.061	0	0.4	75.0	10.0	0.007	↵
↵	17.956	0.049						
30	0.017	0.000	0	0.5	25.0	1.0	0.012	↵
↵	17.873	0.093						
31	0.018	0.002	0	0.5	25.0	10.0	0.012	↵
↵	17.873	0.055						
32	0.165	0.073	0	0.5	50.0	1.0	0.012	↵
↵	17.873	0.041						
33	0.084	0.021	0	0.5	50.0	10.0	0.012	↵
↵	17.873	0.050						
34	0.163	0.100	0	0.5	75.0	1.0	0.012	↵
↵	17.873	0.039						
35	0.040	0.014	0	0.5	75.0	10.0	0.012	↵
↵	17.873	0.050						



The notebook for this tutorial can be found [here](#).

3.1 Installing dpEmu

Works only with Python versions 3.6 and 3.7

To install dpEmu on your computer, execute the following commands in your terminal:

```
git clone https://github.com/dpEmu/dpEmu.git
cd dpEmu
python3 -m venv venv
source venv/bin/activate
pip install -U pip setuptools wheel
pip install -r requirements/base.txt
pip install -e "git+https://github.com/cocodataset/cocoapi.git#egg=pycocotools&
↳subdirectory=PythonAPI"
pip install -e .
```

In order to run all of the examples, you'll also need to execute the following command:

```
pip install -r requirements/examples.txt
```

Additionally, run the following command in order to locally build the documentation with Sphinx:

```
pip install -r requirements/docs.txt
```

3.1.1 Object detection example and notebooks requirements

CUDA and cuDNN installations required

Execute the following commands after all of the above:

```
git clone https://github.com/dpEmu/Detectron.git libs/Detectron
./scripts/install_detectron.sh
git clone https://github.com/dpEmu/darknet.git libs/darknet
./scripts/install_darknet.sh
```

3.2 Running dpEmu

3.2.1 Running notebooks

All jupyter notebooks provided can be opened in a browser with:

```
jupyter notebook docs/case_studies/Text_Classification_OCR_Error.ipynb
```

and remotely executed in console with:

```
jupyter nbconvert --to notebook --ExecutePreprocessor.timeout=None --inplace --  
↪execute docs/case_studies/Text_Classification_OCR_Error.ipynb
```

3.2.2 Running scripts

User defined scripts are run similarly as the predefined examples.

Run the examples from project root.

If the examples do not require command line arguments, then they can be run as follows:

```
python3 examples/filter_examples/run_saturation_example_rgb_0_to_1.py
```

If the examples require command line arguments, add them after the name of the file, each one separated by space (the argument 22 tells the angle of the counterclockwise rotation of the picture):

```
python3 examples/filter_examples/run_rotate_example.py 22
```

The interactive mode is used in some examples and is activated by writing `-i`:

```
python3 examples/run_text_classification_example.py test 4 -i
```

3.3 How dpEmu works

dpEmu consists of three components:

- A system for building error generators
- A system for running ML models with different error parameters
- Tools for visualizing the results

3.3.1 Error Generation

For a quick hands-on introduction to error generation in dpEmu, see the [Error Generation Basics tutorial](#).

Error generation in dpEmu consists of three simple steps:

- Defining the structure of the data by constructing an error generation tree.
- Attaching filters (error sources) to the tree.
- Calling the `generate_error` method on the root node of the tree.

Creating an Error Generation Tree

Error generation trees consist of tree nodes. The most common type of leaf node is the `Array`, which can represent a Numpy array (or Python list) of any dimension. Even a scalar value can be represented by an `Array` node provided that node is not the root of the tree. If the fundamental unit of your data is a tuple (as is the case with, e.g. `.wav` audio data), use a `Tuple` node as the leaf.

The simplest and most commonly used non-leaf node type is the `Series`. The `Series` represents the leftmost dimension of any unit of data passed to it. For example, you might choose to represent a matrix of data as a series of rows. In that case you would then create an `Array` node to represent a row and provide it as the argument to a `Series` node:

```
from dpemu.nodes import Array, Series

row_node = Array()
root_node = Series(row_node)
```

A `TupleSeries` represents a tuple where the first (i.e. leftmost) dimensions of the tuple elements are in some sense “the same”. For example, if we have one Numpy array, `X`, containing the input data and another, `Y`, containing each data point’s correct label, we may choose to represent `(X, Y)` as a `TupleSeries`.

There is usually more than one valid way to represent the structure of the data as a tree. For example, a 2d Numpy array can be represented as:

- a matrix, i.e. a single `Array` node
- a list of rows, i.e. a `Series` with an `Array` as its child
- a list of lists of scalars, i.e. a `Series` whose child is a `Series` whose child is an `Array`.

Adding Filters (Error Sources)

Filters can be added to leaf nodes such as `Array` or `Tuple` nodes. Dozens of filters (e.g. `Snow`, `Blur` and `SensorDrift`) are provided out of the box. They can be used to manipulate practically any kind of data, including but not limited to images, time series and sound. Users can also create their own custom error sources by subclassing the `Filter` class.

To create a filter, call the constructor and provide string identifiers for the error parameters of that filter. To attach the filter to a leaf node, call the node’s `addfilter` method with the filter object as the parameter.

Calling the `generate_error` Method

Once you have defined your error generation tree and added the desired filters, you can call the `generate_error` method of the root node of the tree. The method takes two arguments:

- the data into which the errors are to be introduced, and
- a dictionary of error parameters.

The parameter dictionary contains the error parameter values that are to be used in the error generation. The keys corresponding to the values are the error parameter identifier strings which you provided to the `Filter` constructor(s).

The `generate_error` method does not overwrite the original data but returns a copy instead.

This is an example of what the error generation process might look like:

```
1 import numpy as np
2 from dpemu.filters.common import Missing
3 from dpemu.nodes import Array, TupleSeries
4
5 # Assume our data is a tuple of the form (x, y) where x has
6 # shape (100, 10) and y has shape (100,). We can think of each
7 # row i as a data point where x_i represents the values of the
8 # explanatory variables and y_i represents the corresponding
9 # value of the response variable.
```

(continues on next page)

(continued from previous page)

```

10 x = np.random.rand(100, 10)
11 y = np.random.rand(100, 1)
12 data = (x, y)
13
14 # Build a data model tree.
15 x_node = Array()
16 y_node = Array()
17 root_node = TupleSeries([x_node, y_node])
18
19 # Suppose we want to introduce NaN values (i.e. missing data)
20 # to y only (thus keeping x intact).
21 probability = .3
22 y_node.addfilter(Missing("p", "missing_val"))
23
24 # Feed the data to the root node.
25 output = root_node.generate_error(data, {'p': probability, 'missing_val': np.nan})
26
27 print("Output type (should be tuple):", type(output))
28 print("Output length (should be 2):", len(output))
29 print("Shape of first member of output tuple (should be (100, 10)):",
30       output[0].shape)
31 print("Shape of second first member of output tuple (should be (100,)):",
32       output[1].shape)
33 print("Number of NaNs in x (should be 0):",
34       np.isnan(output[0]).sum())
35 print(f"Number of NaNs in y (should be close to {probability * y.size}):",
36       np.isnan(output[1]).sum())

```

In the example the error generation tree has a `TupleSeries` as its root node, which in turn has two `Array` nodes as its children. Then on line 19 we add a `Missing` filter to one of the children, which will transform some of the values in the 1-dimensional array `y` to NaN. The filter is given a parameter with value “*p*”, which means that the key for the probability for transforming a number into NaN is going to be “*p*” in the parameter dictionary.

We then create a `GaussianNoise` filter and attach it to `x_node`, the other child of the root node. The `GaussianNoise` filter takes two string identifier arguments, corresponding to the mean and standard deviation of the Gaussian distribution from which the noise is drawn.

Finally we call the `generate_error` method of the root node, providing it with the data and the error parameter dictionary. The method returns an errorified copy of the data. However, if you wish to run a machine learning model on the data, the ML runner – to be discussed next – will call the method for you.

ML runner system

The ML runner system, or simply runner, is a system which is used for running multiple machine learning models simultaneously with distinct filter error parameters by using multithreading. After running all the models with all desired parameter combinations the system returns a `pandas.DataFrame` object which can be used for visualizing the results.

The runner needs to be given the following values when it is run: train data, test data, a preprocessor, an error generation tree, a list of error parameters, a list of ML models and their parameters and a boolean indicating whether or not to use interactive mode.

Train data and test data

These are the original train data and test data which will be given to the ML models. A `None` value can also be passed to the runner if there is no training data.

Preprocessor

The preprocessor needs to implement a function `run(train_data, test_data)` and it returns the preprocessed train and test data. The preprocessor can return additional data as well, and it will be listed as separate columns in the `DataFrame` which the runner returns. Here is a simple example of a preprocessor, which does nothing to the original data, but returns also an array called “*negative_data*” which contains the additive inverse of each `test_data`’s element.

```

1 class Preprocessor:
2     def __init__(self):
3         self.random_state = RandomState(42)
4
5     def run(self, train_data, test_data):
6         negative_data = -test_data
7         return train_data, test_data, {"negative_data": negative_data}

```

Error generation tree

The root node of the error generation tree should be given to the runner. The structure of the error generation tree is described above.

Error parameter list

The list of error parameters is simply a list of dictionaries which contain the keys and error values for the error generation tree.

AI model parameter list

The list of AI model parameters is a list of dictionaries containing three keys: “*model*”, “*params_list*” and “*use_clean_train_data*”.

The value of “*model*” is **a class instead of an object**. The given class should implement the function `run(train_data, test_data, parameters)` which runs the model on the train data and test data with given parameters and returns a dictionary containing the scores and possibly additional data.

The value of “*params_list*” is a list of dictionaries where each dictionary contains one set of parameters for model. The model will be given these parameters when the `run(train_data, test_data, parameters)` function is called.

If the “*use_clean_train_data*” boolean is true, then no error will be added to the train data.

Here is an example AI model parameter list and a model:

```

1 from numpy.random import RandomState
2 from sklearn.cluster import KMeans
3 from sklearn.metrics import adjusted_rand_score
4 from sklearn.metrics import adjusted_mutual_info_score
5

```

(continues on next page)

(continued from previous page)

```

6 # Model
7 class KMeansModel:
8     def __init__(self):
9         self.random_state = RandomState(42)
10
11     def run(self, train_data, test_data, model_params):
12         labels = model_params["labels"]
13
14         n_classes = len(np.unique(labels))
15         fitted_model = KMeans(n_clusters=n_classes,
16                               random_state=self.random_state
17                               ).fit(test_data)
18
19         return {
20             "AMI": round(adjusted_mutual_info_score(labels,
21                                                     fitted_model.labels_,
22                                                     average_method="arithmetic"),
23                       3),
24             "ARI": round(adjusted_rand_score(labels, fitted_model.labels_), 3),
25         }
26
27 # Parameter list
28 model_params_dict_list = [
29     {"model": KMeansModel, "params_list": [{"labels": labels}]}
30 ]

```

Interactive mode

The final parameter of the runner system is a boolean indicating whether to use interactive mode or not. Some of the functions for visualizing the results require interactive mode; for others it is optional. Most visualization functions have no interactive functionality.

Basically what the interactive mode does is that it adds a column containing the modified test data to the resulting DataFrame object. The interactive visualizer functions use this data to display points of data so that e.g. the user can try to figure out why something was classified incorrectly.

3.3.2 Visualization functions

The `dpemu.plotting_utils` module contains several functions for plotting and visualizing the data.

3.3.3 A Complete Example

Here is an unrealistic but simple example which demonstrates all three components of dpEmu. In this example we are trying to predict the next value of data when we know all earlier values in the data. Our model tries to do estimate this by keeping a weighted average. In the end of the example a plot of scores is visualized.

```

1 import sys
2
3 import matplotlib.pyplot as plt
4 import numpy as np
5
6 from dpemu import runner

```

(continues on next page)

(continued from previous page)

```

7 from dpemu.plotting_utils import visualize_scores, print_results_by_model, visualize_
  ↳ best_model_params
8 from dpemu.nodes import Array
9 from dpemu.filters.common import GaussianNoise
10
11
12 class Preprocessor:
13     def run(self, train_data, test_data, params):
14         return train_data, test_data, {}
15
16
17 class PredictorModel:
18     def run(self, train_data, test_data, params):
19         # The model tries to predict the values of test_data
20         # by using a weighted average of previous values
21         estimate = 0
22         squared_error = 0
23
24         for elem in test_data:
25             # Calculate error
26             squared_error += (elem - estimate) * (elem - estimate)
27             # Update estimate
28             estimate = (1 - params["weight"]) * estimate + params["weight"] * elem
29
30         mean_squared_error = squared_error / len(test_data)
31
32         return {"MSE": mean_squared_error}
33
34
35 def get_data(argv):
36     train_data = None
37     test_data = np.arange(int(sys.argv[1]))
38     return train_data, test_data
39
40
41 def get_err_root_node():
42     # Create error generation tree that has an Array node
43     # as its root node and a GaussianNoise filter
44     err_root_node = Array()
45     err_root_node.addfilter(GaussianNoise("mean", "std"))
46     return err_root_node
47
48
49 def get_err_params_list():
50     # The standard deviation goes from 0 to 20
51     return [{"mean": 0, "std": std} for std in range(0, 21)]
52
53
54 def get_model_params_dict_list():
55     # The model is run with different weighted estimates
56     return [{
57         "model": PredictorModel,
58         "params_list": [{"weight": w} for w in [0.0, 0.05, 0.15, 0.5, 1.0]],
59         "use_clean_train_data": False
60     }]
61
62

```

(continues on next page)

(continued from previous page)

```

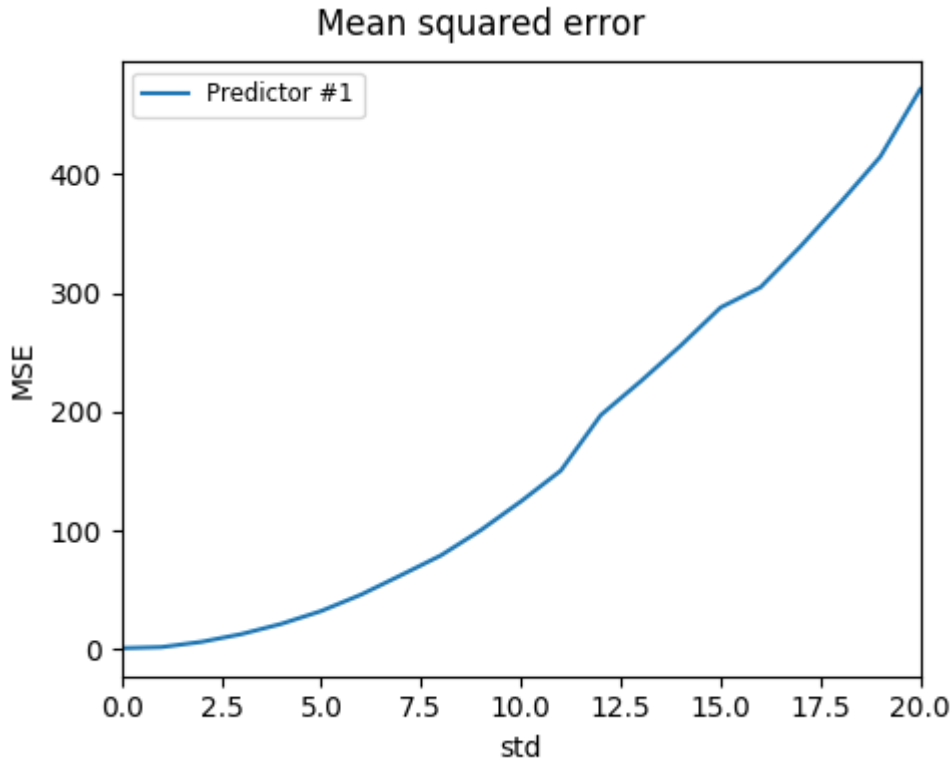
63 def visualize(df):
64     # Visualize mean squared error for all used standard deviations
65     visualize_scores(
66         df=df,
67         score_names=["MSE"],
68         is_higher_score_better=[False],
69         err_param_name="std",
70         title="Mean squared error"
71     )
72     visualize_best_model_params(
73         df=df,
74         model_name="Predictor #1",
75         model_params=["weight"],
76         score_names=["MSE"],
77         is_higher_score_better=[False],
78         err_param_name="std",
79         title=f"Best model params"
80     )
81
82     plt.show()
83
84
85 def main(argv):
86     # Create some fake data
87     if len(argv) == 2:
88         train_data, test_data = get_data(argv)
89     else:
90         exit(0)
91
92     # Run the whole thing and get DataFrame for visualization
93     df = runner.run(
94         train_data=train_data,
95         test_data=test_data,
96         preproc=Preprocessor,
97         preproc_params=None,
98         err_root_node=get_err_root_node(),
99         err_params_list=get_err_params_list(),
100         model_params_dict_list=get_model_params_dict_list()
101     )
102
103     print_results_by_model(df)
104     visualize(df)
105
106
107 if __name__ == "__main__":
108     main(sys.argv)

```

Run the program with the command:

```
python3 examples/run_manual_predictor_example 1000
```

Here's what the resulting image should look like:



3.4 Installation on University of Helsinki clusters (Ukko2 and Kale)

First you need to have access rights to the clusters. See instructions for who can get access rights to [Kale](#) or to [Ukko2](#).

To install dpEmu on Kale or Ukko2 clusters, first establish a ssh connection to the cluster:

```
ssh ukko2.cs.helsinki.fi
```

Or:

```
ssh kale.grid.helsinki.fi
```

To install dpEmu without the ability of running all of the examples, execute the following commands in remote terminal:

```
module load Python/3.7.0-intel-2018b
export SCIKIT_LEARN_DATA=$TMPDIR

cd $WRKDIR
git clone https://github.com/dpEmu/dpEmu.git
cd dpEmu
python3 -m venv venv
source venv/bin/activate
pip install -U pip setuptools wheel --cache-dir $TMPDIR
pip install -r requirements/base.txt --cache-dir $TMPDIR
pip install -e "git+https://github.com/cocodataset/cocoapi.git#egg=pycocotools&
↳subdirectory=PythonAPI" --cache-dir $TMPDIR
pip install -e . --cache-dir $TMPDIR
```

In order to run all of the examples, you'll also need to execute the following command:

```
pip install -r requirements/examples.txt --cache-dir $TMPDIR
```

3.4.1 Object detection example and notebooks requirements

Further installation steps are needed to run the object detection example or notebooks. Execute the following commands after all of the above:

```
module load CUDA/10.0.130
module load cuDNN/7.5.0.56-CUDA-10.0.130

git clone https://github.com/dpEmu/Detectron.git libs/Detectron
./scripts/install_detectron.sh
git clone https://github.com/dpEmu/darknet.git libs/darknet
./scripts/install_darknet.sh
```

Instructions and examples for running jobs on Kale or Ukko2.

3.5 Instructions and examples for running jobs on Kale or Ukko2

3.5.1 Official instructions

Kale

Ukko2

3.5.2 Example jobs

The following commands need to be run every time you log in to one of the clusters:

```
module load Python/3.7.0-intel-2018b
export SCIKIT_LEARN_DATA=$TMPDIR

cd $WRKDIR/dpEmu
source venv/bin/activate
```

Running text classification example

Create the batch file for the job:

```
nano text_classification.job
```

Then write the following content to it and save the file. **Remember to put your username in place of <username>:**

```
#!/bin/bash
#SBATCH -J dpEmu
#SBATCH --workdir=/wrk/users/<username>/dpEmu/
#SBATCH -o text_classification_results.txt
#SBATCH -c 8
#SBATCH --mem=64G
```

(continues on next page)

(continued from previous page)

```
#SBATCH -t 10:00

srun python3 examples/run_text_classification_example.py all 10
srun sleep 60
```

Submit the batch job to be run:

```
sbatch text_classification.job
```

You can view the execution of the code as if it was executed on your home terminal with:

```
tail -f text_classification_results.txt
```

The resulting images will be saved to the dpEmu/out directory.

Running object detection example

First remember to load the required modules and install the object detection example requirements while in the virtual environment, if not done already: *Object detection example and notebooks requirements*.

Create the batch file for the job:

```
nano object_detection.job
```

Then write the following content to it and save the file. **Remember to put your username in place of <username>:**

```
#!/bin/bash
#SBATCH -J dpEmu
#SBATCH --workdir=/wrk/users/<username>/dpEmu/
#SBATCH -o object_detection_results.txt
#SBATCH -c 4
#SBATCH --mem=32G
#SBATCH -p gpu
#SBATCH --gres=gpu:1
#SBATCH -t 10:00:00

srun python3 examples/run_object_detection_example.py
srun sleep 60
```

Running this example can take a lot of time. You could try to disable some of the slowest models i.e. FasterRCNN and RetinaNet. To further speed up the job on Kale, by using the latest GPUs, add the following line to the batch file:

```
#SBATCH --constraint=v100
```

Submit the batch job to be run:

```
sbatch object_detection.job
```

You can view the execution of the code as if it was executed on your home terminal with:

```
tail -f object_detection_results.txt
```

The resulting images will be saved to the dpEmu/out directory.

Running object detection notebook

In the batch file replace:

```
srun python3 examples/run_object_detection_example.py
```

with for example:

```
srun jupyter nbconvert --to notebook --ExecutePreprocessor.timeout=None --inplace --  
↪execute docs/case_studies/Object_Detection_JPEG_Compression.ipynb
```

3.6 License

The MIT License (MIT)

Copyright (c) 2019 Tuomas Halvari, Juha Harviainen, Juha Mylläri, Antti Röyskö, Juuso Silvennoinen

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

CASE STUDIES

4.1 Image Clustering: Added noise

Warning: Agglomerative clustering scales badly with dataset size. This leads to high memory usage (about 360 GB on Kale).

```
[1]: import warnings
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
import numpy as np
from hdbscan import HDBSCAN
from numba.errors import NumbaDeprecationWarning, NumbaWarning
from numpy.random import RandomState
from sklearn.cluster import KMeans, AgglomerativeClustering
from sklearn.metrics import adjusted_rand_score, adjusted_mutual_info_score

from dpemu import runner
from dpemu.dataset_utils import load_mnist_unsplit
from dpemu.filters.common import GaussianNoise, Clip
from dpemu.ml_utils import reduce_dimensions
from dpemu.nodes import Array
from dpemu.nodes.series import Series
from dpemu.plotting_utils import visualize_best_model_params, visualize_scores, \
    visualize_classes, \
    print_results_by_model

warnings.simplefilter("ignore", category=NumbaDeprecationWarning)
warnings.simplefilter("ignore", category=NumbaWarning)

/wrk/users/thalvari/dpEmu/venv/lib/python3.7/site-packages/sklearn/externals/six.
→py:31: DeprecationWarning: The module is deprecated in version 0.21 and will be
→removed in version 0.23 since we've dropped support for Python 2.7. Please rely on
→the official version of six (https://pypi.org/project/six/).
→ " (https://pypi.org/project/six/).", DeprecationWarning)
/wrk/users/thalvari/dpEmu/venv/lib/python3.7/site-packages/sklearn/externals/joblib/
→__init__.py:15: DeprecationWarning: sklearn.externals.joblib is deprecated in 0.21
→and will be removed in 0.23. Please import this functionality directly from joblib,
→which can be installed with: pip install joblib. If this warning is raised when
→loading pickled models, you may need to re-serialize those models with scikit-learn
→0.21+.
    warnings.warn(msg, category=DeprecationWarning)
```

```
[2]: def get_data():
    return load_mnist_unsplit(70000)
```

```
[3]: def get_err_root_node():
    err_img_node = Array(reshape=(28, 28))
    err_root_node = Series(err_img_node)
    err_img_node.addfilter(GaussianNoise("mean", "std"))
    err_img_node.addfilter(Clip("min_val", "max_val"))
    return err_root_node
```

```
[4]: def get_err_params_list(data):
    min_val = np.amin(data)
    max_val = np.amax(data)
    std_steps = np.linspace(0, max_val, num=8)
    err_params_list = [{"mean": 0, "std": std, "min_val": min_val, "max_val": max_val}
    ↪ for std in std_steps]
    return err_params_list
```

```
[5]: class Preprocessor:
    def __init__(self):
        self.random_state = RandomState(42)

    def run(self, _, data, params):
        reduced_data = reduce_dimensions(data, self.random_state)
        return None, reduced_data, {"reduced_data": reduced_data}
```

```
[6]: class AbstractModel(ABC):

    def __init__(self):
        self.random_state = RandomState(42)

    @abstractmethod
    def get_fitted_model(self, data, params):
        pass

    def run(self, _, data, params):
        labels = params["labels"]
        fitted_model = self.get_fitted_model(data, params)
        return {
            "AMI": round(adjusted_mutual_info_score(labels, fitted_model.labels_, ↪
    ↪ average_method="arithmetic"), 3),
            "ARI": round(adjusted_rand_score(labels, fitted_model.labels_), 3),
        }

class KMeansModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        labels = params["labels"]
        n_classes = len(np.unique(labels))
        return KMeans(n_clusters=n_classes, random_state=self.random_state).fit(data)
```

(continues on next page)

(continued from previous page)

```

class AgglomerativeModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        labels = params["labels"]
        n_classes = len(np.unique(labels))
        return AgglomerativeClustering(n_clusters=n_classes).fit(data)

class HDBSCANModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        return HDBSCAN(
            min_samples=params["min_samples"],
            min_cluster_size=params["min_cluster_size"],
            core_dist_n_jobs=1
        ).fit(data)

```

```

[7]: def get_model_params_dict_list(data, labels):
    n_data = data.shape[0]
    divs = [12, 15, 20, 30, 45, 65, 90]
    min_cluster_size_steps = [round(n_data / div) for div in divs]
    min_samples_steps = [1, 5, 10]
    return [
        {"model": KMeansModel, "params_list": [{"labels": labels}]},
        {"model": AgglomerativeModel, "params_list": [{"labels": labels}]},
        {"model": HDBSCANModel, "params_list": [{
            "min_cluster_size": min_cluster_size,
            "min_samples": min_samples,
            "labels": labels
        } for min_cluster_size in min_cluster_size_steps for min_samples in min_
↪ samples_steps]},
    ]

```

```

[8]: def visualize(df, label_names, dataset_name, data):
    visualize_scores(
        df,
        score_names=["AMI", "ARI"],
        is_higher_score_better=[True, True],
        err_param_name="std",
        title=f"{dataset_name} clustering scores with added noise",
    )
    visualize_best_model_params(
        df,
        model_name="HDBSCAN",
        model_params=["min_cluster_size", "min_samples"],
        score_names=["AMI", "ARI"],
        is_higher_score_better=[True, True],
        err_param_name="std",
        title=f"Best parameters for {dataset_name} clustering"
    )

```

(continues on next page)

(continued from previous page)

```

visualize_classes(
    df,
    label_names,
    err_param_name="std",
    reduced_data_column="reduced_data",
    labels_column="labels",
    cmap="tab10",
    title=f"{dataset_name} (n={data.shape[0]}) true classes with added noise"
)
plt.show()

```

```

[9]: def main():
    data, labels, label_names, dataset_name = get_data()

    df = runner.run(
        train_data=None,
        test_data=data,
        preproc=Preprocessor,
        preproc_params=None,
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(data),
        model_params_dict_list=get_model_params_dict_list(data, labels),
    )

    print_results_by_model(df, ["mean", "min_val", "max_val", "reduced_data", "labels
↪"])
    visualize(df, label_names, dataset_name, data)

```

[10]: main()

```
100%|| 8/8 [16:28<00:00, 123.56s/it]
```

```
Agglomerative #1
```

	AMI	ARI	std	time_err	time_pre	time_mod
0	0.886	0.819	0.000000	6.254	398.776	180.600
1	0.883	0.816	36.428571	6.359	413.132	203.826
2	0.872	0.806	72.857143	6.636	421.233	178.826
3	0.818	0.728	109.285714	6.558	427.936	376.926
4	0.726	0.630	145.714286	6.380	422.852	376.770
5	0.646	0.538	182.142857	6.632	437.173	338.948
6	0.558	0.451	218.571429	9.686	466.863	458.426
7	0.445	0.360	255.000000	6.564	462.590	417.582

```
HDBSCAN #1
```

	AMI	ARI	std	min_cluster_size	min_samples	time_err	time_pre	↪
↪time_mod								
0	0.899	0.850	0.000000	5833.0	1.0	6.254	398.776	↪
↪2.248								
1	0.900	0.850	0.000000	5833.0	5.0	6.254	398.776	↪
↪2.495								
2	0.899	0.850	0.000000	5833.0	10.0	6.254	398.776	↪
↪2.641								
3	0.899	0.850	0.000000	4667.0	1.0	6.254	398.776	↪
↪2.269								
4	0.900	0.850	0.000000	4667.0	5.0	6.254	398.776	↪
↪2.530								

(continues on next page)

(continued from previous page)

5	0.899	0.850	0.000000	4667.0	10.0	6.254	398.776	└
↪2.648								
6	0.899	0.850	0.000000	3500.0	1.0	6.254	398.776	└
↪2.277								
7	0.900	0.850	0.000000	3500.0	5.0	6.254	398.776	└
↪2.528								
8	0.899	0.850	0.000000	3500.0	10.0	6.254	398.776	└
↪2.691								
9	0.899	0.850	0.000000	2333.0	1.0	6.254	398.776	└
↪2.314								
10	0.900	0.850	0.000000	2333.0	5.0	6.254	398.776	└
↪2.546								
11	0.899	0.850	0.000000	2333.0	10.0	6.254	398.776	└
↪2.676								
12	0.899	0.850	0.000000	1556.0	1.0	6.254	398.776	└
↪2.301								
13	0.900	0.850	0.000000	1556.0	5.0	6.254	398.776	└
↪2.563								
14	0.899	0.850	0.000000	1556.0	10.0	6.254	398.776	└
↪2.674								
15	0.899	0.850	0.000000	1077.0	1.0	6.254	398.776	└
↪2.301								
16	0.900	0.850	0.000000	1077.0	5.0	6.254	398.776	└
↪2.579								
17	0.899	0.850	0.000000	1077.0	10.0	6.254	398.776	└
↪2.673								
18	0.899	0.850	0.000000	778.0	1.0	6.254	398.776	└
↪2.321								
19	0.900	0.850	0.000000	778.0	5.0	6.254	398.776	└
↪2.557								
20	0.899	0.850	0.000000	778.0	10.0	6.254	398.776	└
↪2.677								
21	0.896	0.848	36.428571	5833.0	1.0	6.359	413.132	└
↪2.298								
22	0.897	0.848	36.428571	5833.0	5.0	6.359	413.132	└
↪2.551								
23	0.896	0.848	36.428571	5833.0	10.0	6.359	413.132	└
↪2.745								
24	0.896	0.848	36.428571	4667.0	1.0	6.359	413.132	└
↪2.314								
25	0.897	0.848	36.428571	4667.0	5.0	6.359	413.132	└
↪2.579								
26	0.896	0.848	36.428571	4667.0	10.0	6.359	413.132	└
↪2.769								
27	0.896	0.848	36.428571	3500.0	1.0	6.359	413.132	└
↪2.340								
28	0.897	0.848	36.428571	3500.0	5.0	6.359	413.132	└
↪2.601								
29	0.896	0.848	36.428571	3500.0	10.0	6.359	413.132	└
↪2.775								
30	0.896	0.848	36.428571	2333.0	1.0	6.359	413.132	└
↪2.318								
31	0.897	0.848	36.428571	2333.0	5.0	6.359	413.132	└
↪2.581								
32	0.896	0.848	36.428571	2333.0	10.0	6.359	413.132	└
↪2.735								
33	0.896	0.848	36.428571	1556.0	1.0	6.359	413.132	└
↪2.313								

(continues on next page)

(continued from previous page)

34	0.897	0.848	36.428571	1556.0	5.0	6.359	413.132	└
↪2.582								
35	0.896	0.848	36.428571	1556.0	10.0	6.359	413.132	└
↪2.743								
36	0.896	0.848	36.428571	1077.0	1.0	6.359	413.132	└
↪2.324								
37	0.897	0.848	36.428571	1077.0	5.0	6.359	413.132	└
↪2.601								
38	0.896	0.848	36.428571	1077.0	10.0	6.359	413.132	└
↪2.750								
39	0.896	0.848	36.428571	778.0	1.0	6.359	413.132	└
↪2.340								
40	0.897	0.848	36.428571	778.0	5.0	6.359	413.132	└
↪2.623								
41	0.896	0.848	36.428571	778.0	10.0	6.359	413.132	└
↪2.764								
42	0.824	0.653	72.857143	5833.0	1.0	6.636	421.233	└
↪2.241								
43	0.825	0.653	72.857143	5833.0	5.0	6.636	421.233	└
↪2.539								
44	0.825	0.653	72.857143	5833.0	10.0	6.636	421.233	└
↪2.688								
45	0.824	0.653	72.857143	4667.0	1.0	6.636	421.233	└
↪2.256								
46	0.825	0.653	72.857143	4667.0	5.0	6.636	421.233	└
↪2.560								
47	0.825	0.653	72.857143	4667.0	10.0	6.636	421.233	└
↪2.694								
48	0.824	0.653	72.857143	3500.0	1.0	6.636	421.233	└
↪2.244								
49	0.881	0.835	72.857143	3500.0	5.0	6.636	421.233	└
↪2.569								
50	0.880	0.834	72.857143	3500.0	10.0	6.636	421.233	└
↪2.720								
51	0.824	0.653	72.857143	2333.0	1.0	6.636	421.233	└
↪2.252								
52	0.881	0.835	72.857143	2333.0	5.0	6.636	421.233	└
↪2.569								
53	0.880	0.834	72.857143	2333.0	10.0	6.636	421.233	└
↪2.761								
54	0.824	0.653	72.857143	1556.0	1.0	6.636	421.233	└
↪2.251								
55	0.881	0.835	72.857143	1556.0	5.0	6.636	421.233	└
↪2.481								
56	0.880	0.834	72.857143	1556.0	10.0	6.636	421.233	└
↪2.629								
57	0.824	0.653	72.857143	1077.0	1.0	6.636	421.233	└
↪2.166								
58	0.881	0.835	72.857143	1077.0	5.0	6.636	421.233	└
↪2.606								
59	0.880	0.834	72.857143	1077.0	10.0	6.636	421.233	└
↪2.737								
60	0.824	0.653	72.857143	778.0	1.0	6.636	421.233	└
↪2.300								
61	0.881	0.835	72.857143	778.0	5.0	6.636	421.233	└
↪2.605								
62	0.880	0.834	72.857143	778.0	10.0	6.636	421.233	└
↪2.807								

(continues on next page)

(continued from previous page)

63	0.802	0.640	109.285714	5833.0	1.0	6.558	427.936	└
↪2.186								
64	0.803	0.640	109.285714	5833.0	5.0	6.558	427.936	└
↪2.511								
65	0.803	0.640	109.285714	5833.0	10.0	6.558	427.936	└
↪2.617								
66	0.802	0.640	109.285714	4667.0	1.0	6.558	427.936	└
↪2.196								
67	0.803	0.640	109.285714	4667.0	5.0	6.558	427.936	└
↪2.515								
68	0.803	0.640	109.285714	4667.0	10.0	6.558	427.936	└
↪2.668								
69	0.802	0.640	109.285714	3500.0	1.0	6.558	427.936	└
↪2.198								
70	0.803	0.640	109.285714	3500.0	5.0	6.558	427.936	└
↪2.502								
71	0.803	0.640	109.285714	3500.0	10.0	6.558	427.936	└
↪2.675								
72	0.802	0.640	109.285714	2333.0	1.0	6.558	427.936	└
↪2.199								
73	0.803	0.640	109.285714	2333.0	5.0	6.558	427.936	└
↪2.522								
74	0.803	0.640	109.285714	2333.0	10.0	6.558	427.936	└
↪2.676								
75	0.802	0.640	109.285714	1556.0	1.0	6.558	427.936	└
↪2.219								
76	0.803	0.640	109.285714	1556.0	5.0	6.558	427.936	└
↪2.570								
77	0.803	0.640	109.285714	1556.0	10.0	6.558	427.936	└
↪2.682								
78	0.802	0.640	109.285714	1077.0	1.0	6.558	427.936	└
↪2.232								
79	0.803	0.640	109.285714	1077.0	5.0	6.558	427.936	└
↪2.579								
80	0.803	0.640	109.285714	1077.0	10.0	6.558	427.936	└
↪2.675								
81	0.802	0.640	109.285714	778.0	1.0	6.558	427.936	└
↪2.250								
82	0.803	0.640	109.285714	778.0	5.0	6.558	427.936	└
↪2.573								
83	0.803	0.640	109.285714	778.0	10.0	6.558	427.936	└
↪2.726								
84	0.724	0.519	145.714286	5833.0	1.0	6.380	422.852	└
↪2.209								
85	0.724	0.519	145.714286	5833.0	5.0	6.380	422.852	└
↪2.585								
86	0.760	0.611	145.714286	5833.0	10.0	6.380	422.852	└
↪2.730								
87	0.724	0.519	145.714286	4667.0	1.0	6.380	422.852	└
↪2.201								
88	0.724	0.519	145.714286	4667.0	5.0	6.380	422.852	└
↪2.585								
89	0.760	0.611	145.714286	4667.0	10.0	6.380	422.852	└
↪2.770								
90	0.724	0.519	145.714286	3500.0	1.0	6.380	422.852	└
↪2.212								
91	0.724	0.519	145.714286	3500.0	5.0	6.380	422.852	└
↪2.594								

(continues on next page)

(continued from previous page)

92	0.760	0.611	145.714286	3500.0	10.0	6.380	422.852	└
↪2.770								
93	0.724	0.519	145.714286	2333.0	1.0	6.380	422.852	└
↪2.228								
94	0.724	0.519	145.714286	2333.0	5.0	6.380	422.852	└
↪2.595								
95	0.760	0.611	145.714286	2333.0	10.0	6.380	422.852	└
↪2.793								
96	0.724	0.519	145.714286	1556.0	1.0	6.380	422.852	└
↪2.234								
97	0.724	0.519	145.714286	1556.0	5.0	6.380	422.852	└
↪2.658								
98	0.760	0.611	145.714286	1556.0	10.0	6.380	422.852	└
↪2.801								
99	0.724	0.519	145.714286	1077.0	1.0	6.380	422.852	└
↪2.235								
100	0.724	0.519	145.714286	1077.0	5.0	6.380	422.852	└
↪2.636								
101	0.760	0.611	145.714286	1077.0	10.0	6.380	422.852	└
↪2.826								
102	0.724	0.519	145.714286	778.0	1.0	6.380	422.852	└
↪2.251								
103	0.724	0.519	145.714286	778.0	5.0	6.380	422.852	└
↪2.614								
104	0.760	0.611	145.714286	778.0	10.0	6.380	422.852	└
↪2.794								
105	0.662	0.489	182.142857	5833.0	1.0	6.632	437.173	└
↪2.337								
106	0.661	0.488	182.142857	5833.0	5.0	6.632	437.173	└
↪2.726								
107	0.665	0.483	182.142857	5833.0	10.0	6.632	437.173	└
↪2.898								
108	0.662	0.489	182.142857	4667.0	1.0	6.632	437.173	└
↪2.312								
109	0.661	0.488	182.142857	4667.0	5.0	6.632	437.173	└
↪2.711								
110	0.665	0.483	182.142857	4667.0	10.0	6.632	437.173	└
↪2.659								
111	0.662	0.489	182.142857	3500.0	1.0	6.632	437.173	└
↪2.011								
112	0.661	0.488	182.142857	3500.0	5.0	6.632	437.173	└
↪2.358								
113	0.665	0.483	182.142857	3500.0	10.0	6.632	437.173	└
↪2.593								
114	0.662	0.489	182.142857	2333.0	1.0	6.632	437.173	└
↪2.032								
115	0.661	0.488	182.142857	2333.0	5.0	6.632	437.173	└
↪2.380								
116	0.665	0.483	182.142857	2333.0	10.0	6.632	437.173	└
↪2.529								
117	0.679	0.548	182.142857	1556.0	1.0	6.632	437.173	└
↪2.036								
118	0.661	0.488	182.142857	1556.0	5.0	6.632	437.173	└
↪2.365								
119	0.665	0.483	182.142857	1556.0	10.0	6.632	437.173	└
↪2.509								
120	0.662	0.489	182.142857	1077.0	1.0	6.632	437.173	└
↪2.029								

(continues on next page)

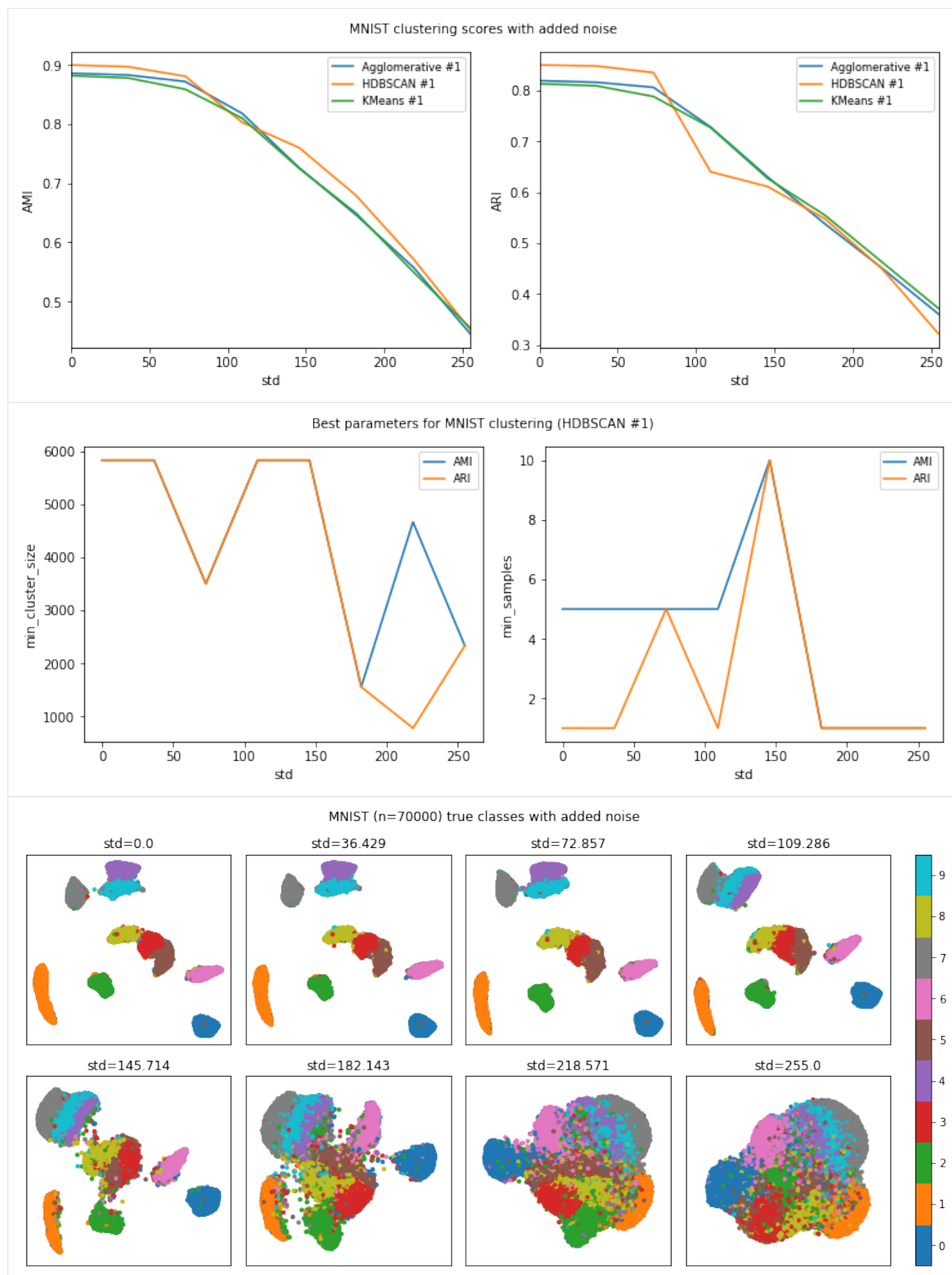
(continued from previous page)

121	0.661	0.488	182.142857	1077.0	5.0	6.632	437.173	└
↪2.382								
122	0.665	0.483	182.142857	1077.0	10.0	6.632	437.173	└
↪2.506								
123	0.662	0.489	182.142857	778.0	1.0	6.632	437.173	└
↪2.048								
124	0.661	0.488	182.142857	778.0	5.0	6.632	437.173	└
↪2.364								
125	0.665	0.483	182.142857	778.0	10.0	6.632	437.173	└
↪2.520								
126	0.558	0.425	218.571429	5833.0	1.0	9.686	466.863	└
↪2.015								
127	0.554	0.420	218.571429	5833.0	5.0	9.686	466.863	└
↪2.406								
128	0.554	0.423	218.571429	5833.0	10.0	9.686	466.863	└
↪2.589								
129	0.572	0.439	218.571429	4667.0	1.0	9.686	466.863	└
↪1.994								
130	0.570	0.438	218.571429	4667.0	5.0	9.686	466.863	└
↪2.410								
131	0.569	0.437	218.571429	4667.0	10.0	9.686	466.863	└
↪2.588								
132	0.572	0.439	218.571429	3500.0	1.0	9.686	466.863	└
↪2.027								
133	0.570	0.438	218.571429	3500.0	5.0	9.686	466.863	└
↪2.426								
134	0.569	0.437	218.571429	3500.0	10.0	9.686	466.863	└
↪2.573								
135	0.572	0.439	218.571429	2333.0	1.0	9.686	466.863	└
↪2.042								
136	0.570	0.438	218.571429	2333.0	5.0	9.686	466.863	└
↪2.442								
137	0.569	0.437	218.571429	2333.0	10.0	9.686	466.863	└
↪2.585								
138	0.572	0.439	218.571429	1556.0	1.0	9.686	466.863	└
↪2.061								
139	0.570	0.438	218.571429	1556.0	5.0	9.686	466.863	└
↪2.469								
140	0.569	0.437	218.571429	1556.0	10.0	9.686	466.863	└
↪2.583								
141	0.572	0.439	218.571429	1077.0	1.0	9.686	466.863	└
↪2.067								
142	0.570	0.438	218.571429	1077.0	5.0	9.686	466.863	└
↪2.475								
143	0.569	0.437	218.571429	1077.0	10.0	9.686	466.863	└
↪2.607								
144	0.570	0.450	218.571429	778.0	1.0	9.686	466.863	└
↪2.068								
145	0.418	0.151	218.571429	778.0	5.0	9.686	466.863	└
↪2.443								
146	0.569	0.437	218.571429	778.0	10.0	9.686	466.863	└
↪2.620								
147	0.427	0.300	255.000000	5833.0	1.0	6.564	462.590	└
↪1.966								
148	0.411	0.280	255.000000	5833.0	5.0	6.564	462.590	└
↪2.493								
149	0.282	0.098	255.000000	5833.0	10.0	6.564	462.590	└
↪2.765								

(continues on next page)

(continued from previous page)

150	0.444	0.316	255.000000	4667.0	1.0	6.564	462.590	└
↪2.008								
151	0.432	0.302	255.000000	4667.0	5.0	6.564	462.590	└
↪2.563								
152	0.431	0.305	255.000000	4667.0	10.0	6.564	462.590	└
↪2.761								
153	0.444	0.316	255.000000	3500.0	1.0	6.564	462.590	└
↪2.055								
154	0.432	0.302	255.000000	3500.0	5.0	6.564	462.590	└
↪2.568								
155	0.431	0.305	255.000000	3500.0	10.0	6.564	462.590	└
↪2.762								
156	0.452	0.321	255.000000	2333.0	1.0	6.564	462.590	└
↪2.059								
157	0.432	0.302	255.000000	2333.0	5.0	6.564	462.590	└
↪2.561								
158	0.431	0.305	255.000000	2333.0	10.0	6.564	462.590	└
↪2.777								
159	0.172	0.036	255.000000	1556.0	1.0	6.564	462.590	└
↪2.018								
160	0.449	0.311	255.000000	1556.0	5.0	6.564	462.590	└
↪2.571								
161	0.433	0.314	255.000000	1556.0	10.0	6.564	462.590	└
↪2.784								
162	0.172	0.036	255.000000	1077.0	1.0	6.564	462.590	└
↪2.016								
163	0.449	0.311	255.000000	1077.0	5.0	6.564	462.590	└
↪2.573								
164	0.433	0.314	255.000000	1077.0	10.0	6.564	462.590	└
↪2.763								
165	0.172	0.036	255.000000	778.0	1.0	6.564	462.590	└
↪2.005								
166	0.449	0.311	255.000000	778.0	5.0	6.564	462.590	└
↪2.553								
167	0.433	0.314	255.000000	778.0	10.0	6.564	462.590	└
↪2.734								
KMeans #1								
	AMI	ARI	std	time_err	time_pre	time_mod		
0	0.882	0.813	0.000000	6.254	398.776	1.225		
1	0.878	0.809	36.428571	6.359	413.132	1.282		
2	0.859	0.788	72.857143	6.636	421.233	1.285		
3	0.810	0.727	109.285714	6.558	427.936	1.473		
4	0.725	0.627	145.714286	6.380	422.852	1.563		
5	0.649	0.555	182.142857	6.632	437.173	1.785		
6	0.550	0.463	218.571429	9.686	466.863	3.467		
7	0.455	0.371	255.000000	6.564	462.590	3.420		



The notebook for this case study can be found [here](#).

4.2 Image Clustering: With rotation

Warning: Agglomerative clustering scales badly with dataset size. This leads to high memory usage (about 360 GB on Kale).

```
[1]: import warnings
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
import numpy as np
from hdbscan import HDBSCAN
from numba.errors import NumbaDeprecationWarning, NumbaWarning
from numpy.random import RandomState
from sklearn.cluster import KMeans, AgglomerativeClustering
from sklearn.metrics import adjusted_rand_score, adjusted_mutual_info_score

from dpemu import runner
from dpemu.dataset_utils import load_mnist_unsplit
from dpemu.filters.image import Rotation
from dpemu.ml_utils import reduce_dimensions
from dpemu.nodes import Array
from dpemu.nodes.series import Series
from dpemu.plotting_utils import visualize_best_model_params, visualize_scores, \
    visualize_classes, \
    print_results_by_model

warnings.simplefilter("ignore", category=NumbaDeprecationWarning)
warnings.simplefilter("ignore", category=NumbaWarning)

/wrk/users/thalvari/dpEmu/venv/lib/python3.7/site-packages/sklearn/externals/six.
↳py:31: DeprecationWarning: The module is deprecated in version 0.21 and will be
↳removed in version 0.23 since we've dropped support for Python 2.7. Please rely on
↳the official version of six (https://pypi.org/project/six/).
    "(https://pypi.org/project/six/).", DeprecationWarning)
/wrk/users/thalvari/dpEmu/venv/lib/python3.7/site-packages/sklearn/externals/joblib/
↳__init__.py:15: DeprecationWarning: sklearn.externals.joblib is deprecated in 0.21
↳and will be removed in 0.23. Please import this functionality directly from joblib,
↳which can be installed with: pip install joblib. If this warning is raised when
↳loading pickled models, you may need to re-serialize those models with scikit-learn
↳0.21+.
    warnings.warn(msg, category=DeprecationWarning)
```

```
[2]: def get_data():
    return load_mnist_unsplit(70000)
```

```
[3]: def get_err_root_node():
    err_img_node = Array(reshape=(28, 28))
    err_root_node = Series(err_img_node)
    err_img_node.addfilter(Rotation("min_angle", "max_angle"))
    return err_root_node
```



```
[4]: def get_err_params_list():
    angle_steps = np.linspace(0, 84, num=8)
    err_params_list = [{"min_angle": -a, "max_angle": a} for a in angle_steps]
    return err_params_list

[5]: class Preprocessor:
    def __init__(self):
        self.random_state = RandomState(42)

    def run(self, _, data, params):
        reduced_data = reduce_dimensions(data, self.random_state)
        return None, reduced_data, {"reduced_data": reduced_data}

[6]: class AbstractModel(ABC):

    def __init__(self):
        self.random_state = RandomState(42)

    @abstractmethod
    def get_fitted_model(self, data, params):
        pass

    def run(self, _, data, params):
        labels = params["labels"]
        fitted_model = self.get_fitted_model(data, params)
        return {
            "AMI": round(adjusted_mutual_info_score(labels, fitted_model.labels_,
→average_method="arithmetic"), 3),
            "ARI": round(adjusted_rand_score(labels, fitted_model.labels_), 3),
        }

class KMeansModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        labels = params["labels"]
        n_classes = len(np.unique(labels))
        return KMeans(n_clusters=n_classes, random_state=self.random_state).fit(data)

class AgglomerativeModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, data, params):
        labels = params["labels"]
        n_classes = len(np.unique(labels))
        return AgglomerativeClustering(n_clusters=n_classes).fit(data)

class HDBSCANModel(AbstractModel):

    def __init__(self):
```

(continues on next page)

(continued from previous page)

```

    super().__init__()

    def get_fitted_model(self, data, params):
        return HDBSCAN(
            min_samples=params["min_samples"],
            min_cluster_size=params["min_cluster_size"],
            core_dist_n_jobs=1
        ).fit(data)

```

```

[7]: def get_model_params_dict_list(data, labels):
    n_data = data.shape[0]
    divs = [12, 15, 20, 30, 45, 65, 90]
    min_cluster_size_steps = [round(n_data / div) for div in divs]
    min_samples_steps = [1, 5, 10]
    return [
        {"model": KMeansModel, "params_list": [{"labels": labels}]},
        {"model": AgglomerativeModel, "params_list": [{"labels": labels}]},
        {"model": HDBSCANModel, "params_list": [
            {
                "min_cluster_size": min_cluster_size,
                "min_samples": min_samples,
                "labels": labels
            } for min_cluster_size in min_cluster_size_steps for min_samples in min_
            samples_steps]},
    ]

```

```

[8]: def visualize(df, label_names, dataset_name, data):
    visualize_scores(
        df,
        score_names=["AMI", "ARI"],
        is_higher_score_better=[True, True],
        err_param_name="max_angle",
        title=f"{dataset_name} clustering scores with rotation",
    )
    visualize_best_model_params(
        df,
        model_name="HDBSCAN",
        model_params=["min_cluster_size", "min_samples"],
        score_names=["AMI", "ARI"],
        is_higher_score_better=[True, True],
        err_param_name="max_angle",
        title=f"Best parameters for {dataset_name} clustering"
    )
    visualize_classes(
        df,
        label_names,
        err_param_name="max_angle",
        reduced_data_column="reduced_data",
        labels_column="labels",
        cmap="tab10",
        title=f"{dataset_name} (n={data.shape[0]}) true classes with rotation"
    )
    plt.show()

```

```

[9]: def main():
    data, labels, label_names, dataset_name = get_data()

```

(continues on next page)

(continued from previous page)

```

df = runner.run(
    train_data=None,
    test_data=data,
    preproc=Preprocessor,
    preproc_params=None,
    err_root_node=get_err_root_node(),
    err_params_list=get_err_params_list(),
    model_params_dict_list=get_model_params_dict_list(data, labels),
)

print_results_by_model(df, ["min_angle", "reduced_data", "labels"])
visualize(df, label_names, dataset_name, data)

```

[10]: main()

100%|| 8/8 [11:45<00:00, 88.18s/it]

Agglomerative #1

	AMI	ARI	max_angle	time_err	time_pre	time_mod
0	0.888	0.819	0.0	5.404	445.732	185.305
1	0.871	0.804	12.0	5.886	465.016	183.202
2	0.805	0.707	24.0	6.258	448.705	184.230
3	0.701	0.542	36.0	5.708	454.822	186.735
4	0.689	0.530	48.0	6.299	461.554	185.352
5	0.607	0.419	60.0	16.806	448.424	185.581
6	0.580	0.398	72.0	6.675	452.017	188.154
7	0.522	0.342	84.0	16.912	447.885	186.312

HDBSCAN #1

	AMI	ARI	max_angle	min_cluster_size	min_samples	time_err	time_pre	time_
↪mod								
0	0.902	0.853	0.0	5833.0	1.0	5.404	445.732	2.
↪263								
1	0.901	0.852	0.0	5833.0	5.0	5.404	445.732	2.
↪498								
2	0.903	0.853	0.0	5833.0	10.0	5.404	445.732	2.
↪572								
3	0.902	0.853	0.0	4667.0	1.0	5.404	445.732	2.
↪222								
4	0.901	0.852	0.0	4667.0	5.0	5.404	445.732	2.
↪456								
5	0.903	0.853	0.0	4667.0	10.0	5.404	445.732	2.
↪530								
6	0.902	0.853	0.0	3500.0	1.0	5.404	445.732	2.
↪171								
7	0.901	0.852	0.0	3500.0	5.0	5.404	445.732	2.
↪416								
8	0.903	0.853	0.0	3500.0	10.0	5.404	445.732	2.
↪516								
9	0.902	0.853	0.0	2333.0	1.0	5.404	445.732	2.
↪170								
10	0.901	0.852	0.0	2333.0	5.0	5.404	445.732	2.
↪448								
11	0.903	0.853	0.0	2333.0	10.0	5.404	445.732	2.
↪535								
12	0.902	0.853	0.0	1556.0	1.0	5.404	445.732	2.
↪201								

(continues on next page)

(continued from previous page)

13	0.901	0.852	0.0	1556.0	5.0	5.404	445.732	2.
↪449								
14	0.903	0.853	0.0	1556.0	10.0	5.404	445.732	2.
↪591								
15	0.902	0.853	0.0	1077.0	1.0	5.404	445.732	2.
↪204								
16	0.901	0.852	0.0	1077.0	5.0	5.404	445.732	2.
↪453								
17	0.903	0.853	0.0	1077.0	10.0	5.404	445.732	2.
↪555								
18	0.902	0.853	0.0	778.0	1.0	5.404	445.732	2.
↪230								
19	0.901	0.852	0.0	778.0	5.0	5.404	445.732	2.
↪473								
20	0.903	0.853	0.0	778.0	10.0	5.404	445.732	2.
↪615								
21	0.883	0.837	12.0	5833.0	1.0	5.886	465.016	2.
↪122								
22	0.882	0.836	12.0	5833.0	5.0	5.886	465.016	2.
↪423								
23	0.882	0.836	12.0	5833.0	10.0	5.886	465.016	2.
↪611								
24	0.883	0.837	12.0	4667.0	1.0	5.886	465.016	2.
↪132								
25	0.882	0.836	12.0	4667.0	5.0	5.886	465.016	2.
↪400								
26	0.882	0.836	12.0	4667.0	10.0	5.886	465.016	2.
↪570								
27	0.883	0.837	12.0	3500.0	1.0	5.886	465.016	2.
↪133								
28	0.882	0.836	12.0	3500.0	5.0	5.886	465.016	2.
↪412								
29	0.882	0.836	12.0	3500.0	10.0	5.886	465.016	2.
↪596								
30	0.883	0.837	12.0	2333.0	1.0	5.886	465.016	2.
↪151								
31	0.882	0.836	12.0	2333.0	5.0	5.886	465.016	2.
↪423								
32	0.882	0.836	12.0	2333.0	10.0	5.886	465.016	2.
↪555								
33	0.883	0.837	12.0	1556.0	1.0	5.886	465.016	2.
↪166								
34	0.882	0.836	12.0	1556.0	5.0	5.886	465.016	2.
↪457								
35	0.882	0.836	12.0	1556.0	10.0	5.886	465.016	2.
↪475								
36	0.883	0.837	12.0	1077.0	1.0	5.886	465.016	2.
↪070								
37	0.882	0.836	12.0	1077.0	5.0	5.886	465.016	2.
↪362								
38	0.882	0.836	12.0	1077.0	10.0	5.886	465.016	2.
↪484								
39	0.883	0.837	12.0	778.0	1.0	5.886	465.016	2.
↪096								
40	0.882	0.836	12.0	778.0	5.0	5.886	465.016	2.
↪468								
41	0.882	0.836	12.0	778.0	10.0	5.886	465.016	2.
↪350								

(continues on next page)

(continued from previous page)

42	0.806	0.642	24.0	5833.0	1.0	6.258	448.705	2.
↪170								
43	0.865	0.819	24.0	5833.0	5.0	6.258	448.705	2.
↪429								
44	0.806	0.642	24.0	5833.0	10.0	6.258	448.705	2.
↪599								
45	0.836	0.747	24.0	4667.0	1.0	6.258	448.705	2.
↪149								
46	0.865	0.819	24.0	4667.0	5.0	6.258	448.705	2.
↪403								
47	0.862	0.816	24.0	4667.0	10.0	6.258	448.705	2.
↪559								
48	0.836	0.747	24.0	3500.0	1.0	6.258	448.705	2.
↪152								
49	0.865	0.819	24.0	3500.0	5.0	6.258	448.705	2.
↪415								
50	0.862	0.816	24.0	3500.0	10.0	6.258	448.705	2.
↪576								
51	0.836	0.747	24.0	2333.0	1.0	6.258	448.705	2.
↪147								
52	0.865	0.819	24.0	2333.0	5.0	6.258	448.705	2.
↪420								
53	0.862	0.816	24.0	2333.0	10.0	6.258	448.705	2.
↪571								
54	0.836	0.747	24.0	1556.0	1.0	6.258	448.705	2.
↪153								
55	0.865	0.819	24.0	1556.0	5.0	6.258	448.705	2.
↪437								
56	0.862	0.816	24.0	1556.0	10.0	6.258	448.705	2.
↪589								
57	0.836	0.747	24.0	1077.0	1.0	6.258	448.705	2.
↪161								
58	0.865	0.819	24.0	1077.0	5.0	6.258	448.705	2.
↪460								
59	0.862	0.816	24.0	1077.0	10.0	6.258	448.705	2.
↪608								
60	0.836	0.747	24.0	778.0	1.0	6.258	448.705	2.
↪177								
61	0.865	0.819	24.0	778.0	5.0	6.258	448.705	2.
↪453								
62	0.862	0.816	24.0	778.0	10.0	6.258	448.705	2.
↪617								
63	0.738	0.497	36.0	5833.0	1.0	5.708	454.822	2.
↪295								
64	0.738	0.497	36.0	5833.0	5.0	5.708	454.822	2.
↪409								
65	0.739	0.498	36.0	5833.0	10.0	5.708	454.822	2.
↪567								
66	0.738	0.497	36.0	4667.0	1.0	5.708	454.822	2.
↪174								
67	0.818	0.736	36.0	4667.0	5.0	5.708	454.822	2.
↪459								
68	0.820	0.736	36.0	4667.0	10.0	5.708	454.822	2.
↪735								
69	0.738	0.497	36.0	3500.0	1.0	5.708	454.822	2.
↪196								
70	0.818	0.736	36.0	3500.0	5.0	5.708	454.822	2.
↪468								

(continues on next page)

(continued from previous page)

71	0.820	0.736	36.0	3500.0	10.0	5.708	454.822	2.
↪626								
72	0.819	0.736	36.0	2333.0	1.0	5.708	454.822	2.
↪199								
73	0.835	0.788	36.0	2333.0	5.0	5.708	454.822	2.
↪614								
74	0.837	0.789	36.0	2333.0	10.0	5.708	454.822	2.
↪623								
75	0.738	0.497	36.0	1556.0	1.0	5.708	454.822	2.
↪195								
76	0.835	0.788	36.0	1556.0	5.0	5.708	454.822	2.
↪496								
77	0.820	0.736	36.0	1556.0	10.0	5.708	454.822	2.
↪597								
78	0.738	0.497	36.0	1077.0	1.0	5.708	454.822	2.
↪295								
79	0.835	0.788	36.0	1077.0	5.0	5.708	454.822	2.
↪558								
80	0.837	0.789	36.0	1077.0	10.0	5.708	454.822	2.
↪604								
81	0.738	0.497	36.0	778.0	1.0	5.708	454.822	2.
↪198								
82	0.835	0.788	36.0	778.0	5.0	5.708	454.822	2.
↪481								
83	0.820	0.736	36.0	778.0	10.0	5.708	454.822	2.
↪514								
84	0.731	0.492	48.0	5833.0	1.0	6.299	461.554	2.
↪113								
85	0.730	0.492	48.0	5833.0	5.0	6.299	461.554	2.
↪343								
86	0.731	0.492	48.0	5833.0	10.0	6.299	461.554	2.
↪464								
87	0.731	0.492	48.0	4667.0	1.0	6.299	461.554	2.
↪112								
88	0.730	0.492	48.0	4667.0	5.0	6.299	461.554	2.
↪343								
89	0.731	0.492	48.0	4667.0	10.0	6.299	461.554	2.
↪496								
90	0.731	0.492	48.0	3500.0	1.0	6.299	461.554	2.
↪143								
91	0.805	0.726	48.0	3500.0	5.0	6.299	461.554	2.
↪385								
92	0.731	0.492	48.0	3500.0	10.0	6.299	461.554	2.
↪468								
93	0.812	0.761	48.0	2333.0	1.0	6.299	461.554	2.
↪150								
94	0.805	0.726	48.0	2333.0	5.0	6.299	461.554	2.
↪421								
95	0.731	0.492	48.0	2333.0	10.0	6.299	461.554	2.
↪498								
96	0.812	0.761	48.0	1556.0	1.0	6.299	461.554	2.
↪181								
97	0.802	0.739	48.0	1556.0	5.0	6.299	461.554	2.
↪432								
98	0.731	0.492	48.0	1556.0	10.0	6.299	461.554	2.
↪454								
99	0.731	0.492	48.0	1077.0	1.0	6.299	461.554	2.
↪055								

(continues on next page)

(continued from previous page)

100	0.730	0.492	48.0	1077.0	5.0	6.299	461.554	2.
↪327								
101	0.731	0.492	48.0	1077.0	10.0	6.299	461.554	2.
↪436								
102	0.731	0.492	48.0	778.0	1.0	6.299	461.554	2.
↪093								
103	0.730	0.492	48.0	778.0	5.0	6.299	461.554	2.
↪391								
104	0.731	0.492	48.0	778.0	10.0	6.299	461.554	2.
↪548								
105	0.721	0.487	60.0	5833.0	1.0	16.806	448.424	2.
↪179								
106	0.722	0.487	60.0	5833.0	5.0	16.806	448.424	2.
↪342								
107	0.720	0.486	60.0	5833.0	10.0	16.806	448.424	2.
↪457								
108	0.721	0.487	60.0	4667.0	1.0	16.806	448.424	2.
↪133								
109	0.722	0.487	60.0	4667.0	5.0	16.806	448.424	2.
↪374								
110	0.720	0.486	60.0	4667.0	10.0	16.806	448.424	2.
↪480								
111	0.721	0.487	60.0	3500.0	1.0	16.806	448.424	2.
↪138								
112	0.722	0.487	60.0	3500.0	5.0	16.806	448.424	2.
↪361								
113	0.720	0.486	60.0	3500.0	10.0	16.806	448.424	2.
↪494								
114	0.721	0.487	60.0	2333.0	1.0	16.806	448.424	2.
↪151								
115	0.722	0.487	60.0	2333.0	5.0	16.806	448.424	2.
↪388								
116	0.720	0.486	60.0	2333.0	10.0	16.806	448.424	2.
↪544								
117	0.721	0.487	60.0	1556.0	1.0	16.806	448.424	2.
↪147								
118	0.722	0.487	60.0	1556.0	5.0	16.806	448.424	2.
↪396								
119	0.720	0.486	60.0	1556.0	10.0	16.806	448.424	2.
↪487								
120	0.721	0.487	60.0	1077.0	1.0	16.806	448.424	2.
↪228								
121	0.722	0.487	60.0	1077.0	5.0	16.806	448.424	2.
↪435								
122	0.720	0.486	60.0	1077.0	10.0	16.806	448.424	2.
↪519								
123	0.721	0.487	60.0	778.0	1.0	16.806	448.424	2.
↪202								
124	0.722	0.487	60.0	778.0	5.0	16.806	448.424	2.
↪333								
125	0.720	0.486	60.0	778.0	10.0	16.806	448.424	2.
↪435								
126	0.710	0.482	72.0	5833.0	1.0	6.675	452.017	2.
↪060								
127	0.709	0.481	72.0	5833.0	5.0	6.675	452.017	2.
↪336								
128	0.708	0.481	72.0	5833.0	10.0	6.675	452.017	2.
↪434								

(continues on next page)

(continued from previous page)

129	0.710	0.482	72.0	4667.0	1.0	6.675	452.017	2.
↪066								
130	0.709	0.481	72.0	4667.0	5.0	6.675	452.017	2.
↪339								
131	0.708	0.481	72.0	4667.0	10.0	6.675	452.017	2.
↪443								
132	0.699	0.463	72.0	3500.0	1.0	6.675	452.017	2.
↪076								
133	0.698	0.462	72.0	3500.0	5.0	6.675	452.017	2.
↪344								
134	0.698	0.462	72.0	3500.0	10.0	6.675	452.017	2.
↪471								
135	0.700	0.463	72.0	2333.0	1.0	6.675	452.017	2.
↪103								
136	0.700	0.463	72.0	2333.0	5.0	6.675	452.017	2.
↪363								
137	0.700	0.463	72.0	2333.0	10.0	6.675	452.017	2.
↪460								
138	0.700	0.463	72.0	1556.0	1.0	6.675	452.017	2.
↪113								
139	0.700	0.463	72.0	1556.0	5.0	6.675	452.017	2.
↪382								
140	0.700	0.463	72.0	1556.0	10.0	6.675	452.017	2.
↪488								
141	0.700	0.463	72.0	1077.0	1.0	6.675	452.017	2.
↪098								
142	0.700	0.463	72.0	1077.0	5.0	6.675	452.017	2.
↪382								
143	0.700	0.463	72.0	1077.0	10.0	6.675	452.017	2.
↪477								
144	0.679	0.431	72.0	778.0	1.0	6.675	452.017	2.
↪098								
145	0.700	0.463	72.0	778.0	5.0	6.675	452.017	2.
↪358								
146	0.700	0.463	72.0	778.0	10.0	6.675	452.017	2.
↪470								
147	0.405	0.111	84.0	5833.0	1.0	16.912	447.885	2.
↪167								
148	0.405	0.111	84.0	5833.0	5.0	16.912	447.885	2.
↪388								
149	0.405	0.111	84.0	5833.0	10.0	16.912	447.885	2.
↪536								
150	0.405	0.111	84.0	4667.0	1.0	16.912	447.885	2.
↪131								
151	0.405	0.111	84.0	4667.0	5.0	16.912	447.885	2.
↪422								
152	0.405	0.111	84.0	4667.0	10.0	16.912	447.885	2.
↪556								
153	0.405	0.111	84.0	3500.0	1.0	16.912	447.885	2.
↪146								
154	0.405	0.111	84.0	3500.0	5.0	16.912	447.885	2.
↪409								
155	0.405	0.111	84.0	3500.0	10.0	16.912	447.885	2.
↪557								
156	0.405	0.111	84.0	2333.0	1.0	16.912	447.885	2.
↪147								
157	0.405	0.111	84.0	2333.0	5.0	16.912	447.885	2.
↪461								

(continues on next page)

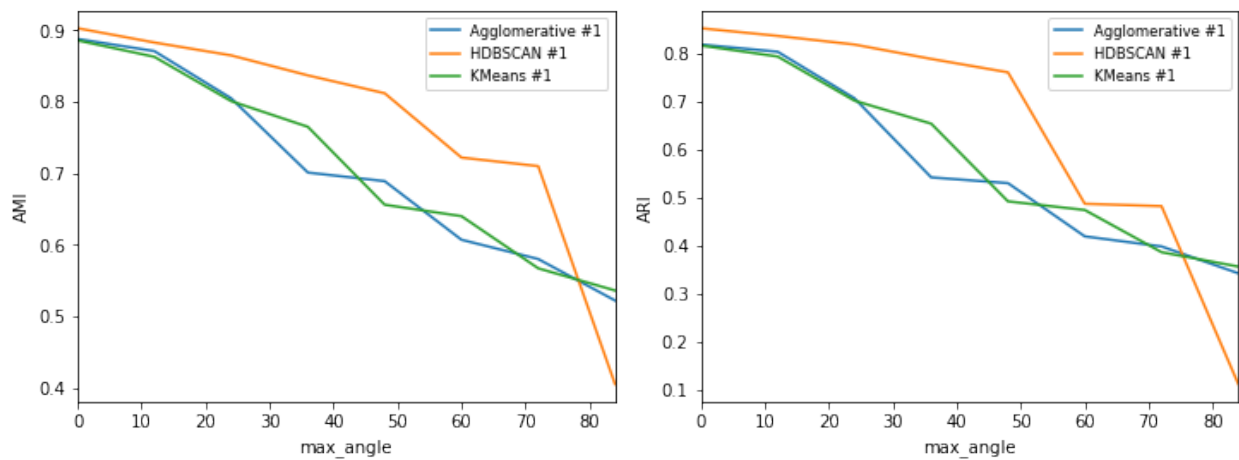
(continued from previous page)

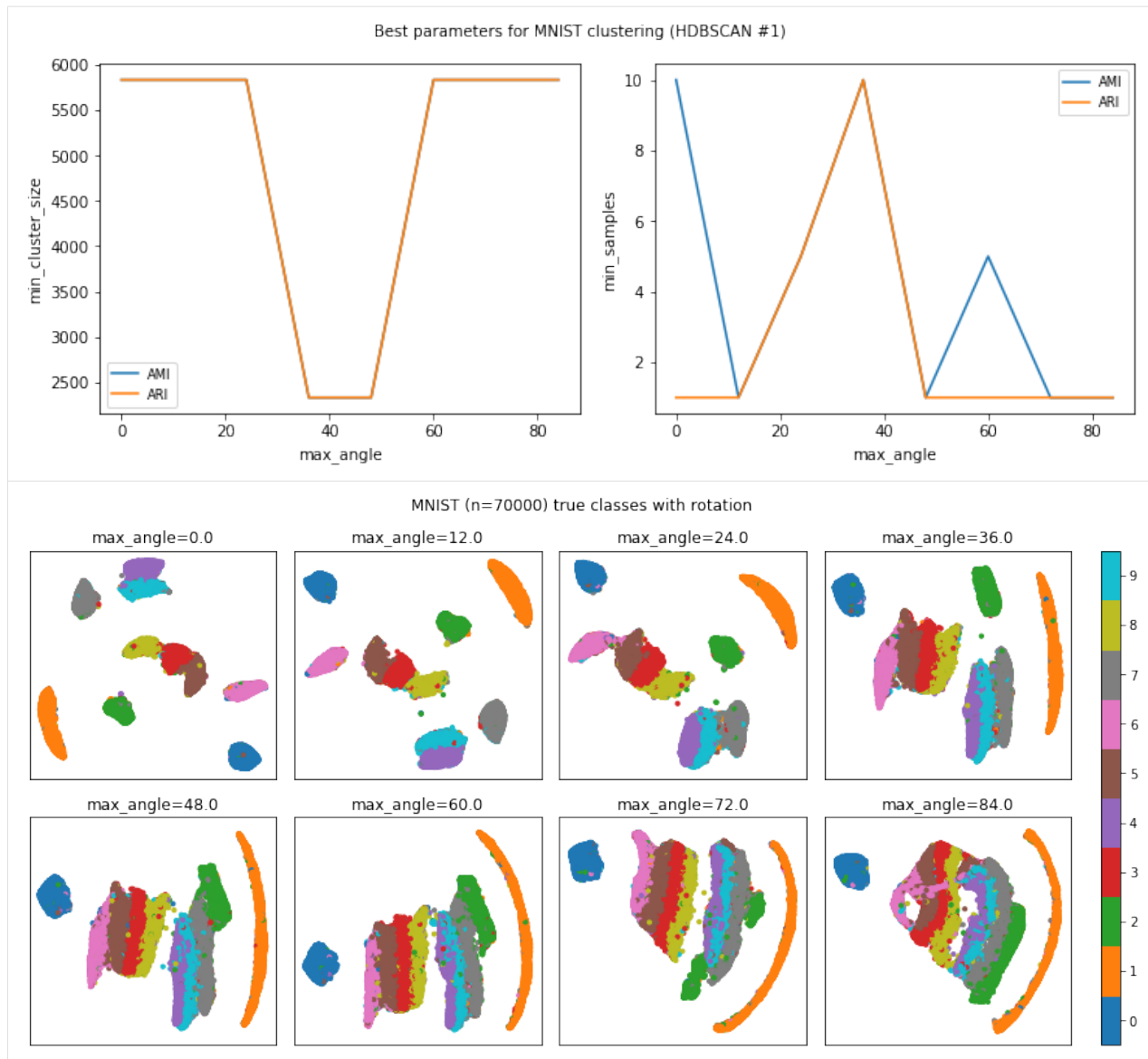
158	0.405	0.111	84.0	2333.0	10.0	16.912	447.885	2.
↪579								
159	0.405	0.111	84.0	1556.0	1.0	16.912	447.885	2.
↪178								
160	0.405	0.111	84.0	1556.0	5.0	16.912	447.885	2.
↪457								
161	0.405	0.111	84.0	1556.0	10.0	16.912	447.885	2.
↪562								
162	0.405	0.111	84.0	1077.0	1.0	16.912	447.885	2.
↪110								
163	0.405	0.111	84.0	1077.0	5.0	16.912	447.885	2.
↪349								
164	0.405	0.111	84.0	1077.0	10.0	16.912	447.885	2.
↪478								
165	0.405	0.111	84.0	778.0	1.0	16.912	447.885	2.
↪105								
166	0.405	0.111	84.0	778.0	5.0	16.912	447.885	2.
↪381								
167	0.405	0.111	84.0	778.0	10.0	16.912	447.885	2.
↪368								

KMeans #1

	AMI	ARI	max_angle	time_err	time_pre	time_mod
0	0.886	0.817	0.0	5.404	445.732	1.163
1	0.863	0.794	12.0	5.886	465.016	1.379
2	0.801	0.702	24.0	6.258	448.705	1.382
3	0.765	0.654	36.0	5.708	454.822	1.749
4	0.656	0.492	48.0	6.299	461.554	2.179
5	0.640	0.474	60.0	16.806	448.424	2.301
6	0.567	0.386	72.0	6.675	452.017	2.628
7	0.536	0.356	84.0	16.912	447.885	2.123

MNIST clustering scores with rotation





The notebook for this case study can be found [here](#).

4.3 Text classification: Missing areas

```
[1]: import warnings
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
import numpy as np
from numba import NumbaDeprecationWarning, NumbaWarning
from numpy.random import RandomState
from sklearn.exceptions import ConvergenceWarning
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import confusion_matrix
```

(continues on next page)

(continued from previous page)

```

from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.svm import LinearSVC

from dpemu import runner
from dpemu.dataset_utils import load_newsgroups
from dpemu.filters.text import MissingArea
from dpemu.ml_utils import reduce_dimensions_sparse
from dpemu.nodes.array import Array
from dpemu.plotting_utils import visualize_best_model_params, visualize_scores, \
    visualize_classes, \
    print_results_by_model, visualize_confusion_matrices
from dpemu.radius_generators import GaussianRadiusGenerator

warnings.simplefilter("ignore", category=ConvergenceWarning)
warnings.simplefilter("ignore", category=NumbaDeprecationWarning)
warnings.simplefilter("ignore", category=NumbaWarning)

```

```

[2]: def get_data():
    data, labels, label_names, dataset_name = load_newsgroups("all", 10)
    train_data, test_data, train_labels, test_labels = train_test_split(data, labels, \
        test_size=.2,
        random_state=RandomState(42))
    return train_data, test_data, train_labels, test_labels, label_names, dataset_name

```

```

[3]: def get_err_root_node():
    err_root_node = Array()
    err_root_node.addfilter(MissingArea("p", "radius_generator", "missing_value"))
    return err_root_node

```

```

[4]: def get_err_params_list():
    p_steps = np.linspace(0, .28, num=8)
    err_params_list = [{
        "p": p,
        "radius_generator": GaussianRadiusGenerator(0, 1),
        "missing_value": " "
    } for p in p_steps]
    return err_params_list

```

```

[5]: class Preprocessor:
    def __init__(self):
        self.random_state = RandomState(0)

    def run(self, train_data, test_data, _):
        vectorizer = TfidfVectorizer(max_df=0.5, min_df=2, stop_words="english")
        vectorized_train_data = vectorizer.fit_transform(train_data)
        vectorized_test_data = vectorizer.transform(test_data)

        reduced_test_data = reduce_dimensions_sparse(vectorized_test_data, self.
            random_state)

        return vectorized_train_data, vectorized_test_data, {"reduced_test_data": \
            reduced_test_data}

```

```
[6]: class AbstractModel(ABC):

    def __init__(self):
        self.random_state = RandomState(42)

    @abstractmethod
    def get_fitted_model(self, train_data, train_labels, params):
        pass

    def run(self, train_data, test_data, params):
        train_labels = params["train_labels"]
        test_labels = params["test_labels"]

        fitted_model = self.get_fitted_model(train_data, train_labels, params)

        predicted_test_labels = fitted_model.predict(test_data)
        cm = confusion_matrix(test_labels, predicted_test_labels)
        return {
            "confusion_matrix": cm,
            "predicted_test_labels": predicted_test_labels,
            "test_mean_accuracy": round(np.mean(predicted_test_labels == test_labels),
→ 3),
            "train_mean_accuracy": fitted_model.score(train_data, train_labels),
        }

class MultinomialNBModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, train_data, train_labels, params):
        return MultinomialNB(params["alpha"]).fit(train_data, train_labels)

class LinearSVCModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, train_data, train_labels, params):
        return LinearSVC(C=params["C"], random_state=self.random_state).fit(train_
→data, train_labels)

[7]: def get_model_params_dict_list(train_labels, test_labels):
    alpha_steps = [10 ** i for i in range(-4, 1)]
    C_steps = [10 ** k for k in range(-3, 2)]
    model_params_base = {"train_labels": train_labels, "test_labels": test_labels}
    return [
        {
            "model": MultinomialNBModel,
            "params_list": [{"alpha": alpha, **model_params_base} for alpha in alpha_
→steps],
            "use_clean_train_data": False
        },
        {
            "model": MultinomialNBModel,
```

(continues on next page)

(continued from previous page)

```

        "params_list": [{"alpha": alpha, **model_params_base} for alpha in alpha_
→steps],
        "use_clean_train_data": True
    },
    {
        "model": LinearSVCModel,
        "params_list": [{"C": C, **model_params_base} for C in C_steps],
        "use_clean_train_data": False
    },
    {
        "model": LinearSVCModel,
        "params_list": [{"C": C, **model_params_base} for C in C_steps],
        "use_clean_train_data": True
    },
]

```

```

[8]: def visualize(df, dataset_name, label_names, test_data):
    visualize_scores(
        df,
        score_names=["test_mean_accuracy", "train_mean_accuracy"],
        is_higher_score_better=[True, True],
        err_param_name="p",
        title=f"{dataset_name} classification scores with added error"
    )
    visualize_best_model_params(
        df,
        "MultinomialNB",
        model_params=["alpha"],
        score_names=["test_mean_accuracy"],
        is_higher_score_better=[True],
        err_param_name="p",
        title=f"Best parameters for {dataset_name} classification",
        y_log=True
    )
    visualize_best_model_params(
        df,
        "LinearSVC",
        model_params=["C"],
        score_names=["test_mean_accuracy"],
        is_higher_score_better=[True],
        err_param_name="p",
        title=f"Best parameters for {dataset_name} classification",
        y_log=True
    )
    visualize_classes(
        df,
        label_names,
        err_param_name="p",
        reduced_data_column="reduced_test_data",
        labels_column="test_labels",
        cmap="tab20",
        title=f"{dataset_name} test set (n={len(test_data)}) true classes with added_
→error"
    )
    visualize_confusion_matrices(
        df,

```

(continues on next page)

(continued from previous page)

```

    label_names,
    score_name="test_mean_accuracy",
    is_higher_score_better=True,
    err_param_name="p",
    labels_col="test_labels",
    predictions_col="predicted_test_labels",
)
plt.show()

```

```

[9]: def main():
    train_data, test_data, train_labels, test_labels, label_names, dataset_name = get_
    ↪data()

    df = runner.run(
        train_data=train_data,
        test_data=test_data,
        preproc=Preprocessor,
        preproc_params=None,
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(train_labels, test_labels),
    )

    print_results_by_model(df, dropped_columns=[
        "train_labels", "test_labels", "reduced_test_data", "confusion_matrix",
    ↪"predicted_test_labels",
        "radius_generator", "missing_value"
    ])
    visualize(df, dataset_name, label_names, test_data)

```

Models LinearSVCclean and MultinomialNBClean have been trained with clean data and LinearSVC and MultinomialNB with erroneous data.

```

[10]: main()
100%|| 8/8 [02:28<00:00, 17.53s/it]
LinearSVC #1

```

	test_mean_accuracy	train_mean_accuracy	p	C	time_err	time_pre	time_
↪mod							
0	0.720	0.792343	0.00	0.001	2.937	19.858	0.
↪198							
1	0.809	0.884620	0.00	0.010	2.937	19.858	0.
↪250							
2	0.842	0.946658	0.00	0.100	2.937	19.858	0.
↪292							
3	0.846	0.973005	0.00	1.000	2.937	19.858	0.
↪468							
4	0.835	0.974043	0.00	10.000	2.937	19.858	2.
↪291							
5	0.633	0.749513	0.04	0.001	27.034	20.458	0.
↪244							
6	0.734	0.875016	0.04	0.010	27.034	20.458	0.
↪282							
7	0.780	0.948215	0.04	0.100	27.034	20.458	0.
↪347							
8	0.787	0.973653	0.04	1.000	27.034	20.458	0.
↪586							

(continues on next page)

(continued from previous page)

9	0.772	0.974173	0.04	10.000	27.034	20.458	2.
↪ 702							
10	0.537	0.695652	0.08	0.001	40.465	19.995	0.
↪ 345							
11	0.649	0.842440	0.08	0.010	40.465	19.995	0.
↪ 321							
12	0.704	0.947696	0.08	0.100	40.465	19.995	0.
↪ 336							
13	0.706	0.973783	0.08	1.000	40.465	19.995	0.
↪ 589							
14	0.686	0.974822	0.08	10.000	40.465	19.995	2.
↪ 649							
15	0.472	0.642440	0.12	0.001	58.590	19.594	0.
↪ 198							
16	0.589	0.823361	0.12	0.010	58.590	19.594	0.
↪ 235							
17	0.652	0.941856	0.12	0.100	58.590	19.594	0.
↪ 298							
18	0.652	0.973913	0.12	1.000	58.590	19.594	0.
↪ 528							
19	0.617	0.974432	0.12	10.000	58.590	19.594	2.
↪ 288							
20	0.420	0.578326	0.16	0.001	82.425	20.489	0.
↪ 190							
21	0.524	0.765347	0.16	0.010	82.425	20.489	0.
↪ 237							
22	0.598	0.924335	0.16	0.100	82.425	20.489	0.
↪ 302							
23	0.601	0.972745	0.16	1.000	82.425	20.489	0.
↪ 579							
24	0.558	0.974043	0.16	10.000	82.425	20.489	2.
↪ 857							
25	0.352	0.537573	0.20	0.001	93.475	21.156	0.
↪ 148							
26	0.449	0.724854	0.20	0.010	93.475	21.156	0.
↪ 188							
27	0.520	0.903310	0.20	0.100	93.475	21.156	0.
↪ 237							
28	0.516	0.968592	0.20	1.000	93.475	21.156	0.
↪ 473							
29	0.475	0.972745	0.20	10.000	93.475	21.156	2.
↪ 220							
30	0.310	0.489552	0.24	0.001	94.942	20.544	0.
↪ 131							
31	0.376	0.655029	0.24	0.010	94.942	20.544	0.
↪ 171							
32	0.461	0.872550	0.24	0.100	94.942	20.544	0.
↪ 208							
33	0.473	0.962362	0.24	1.000	94.942	20.544	0.
↪ 459							
34	0.420	0.971188	0.24	10.000	94.942	20.544	2.
↪ 434							
35	0.275	0.442959	0.28	0.001	125.194	17.138	0.
↪ 115							
36	0.349	0.600130	0.28	0.010	125.194	17.138	0.
↪ 147							
37	0.420	0.838676	0.28	0.100	125.194	17.138	0.
↪ 186							

(continues on next page)

(continued from previous page)

38	0.393	0.950162	0.28	1.000	125.194	17.138	0.
↪425							
39	0.361	0.967683	0.28	10.000	125.194	17.138	2.
↪252							
LinearSVCClean #1							
	test_mean_accuracy	train_mean_accuracy	p	C	time_err	time_pre	time_
↪mod							
0	0.720	0.792343	0.00	0.001	2.937	19.858	0.
↪199							
1	0.809	0.884620	0.00	0.010	2.937	19.858	0.
↪253							
2	0.842	0.946658	0.00	0.100	2.937	19.858	0.
↪289							
3	0.846	0.973005	0.00	1.000	2.937	19.858	0.
↪462							
4	0.835	0.974043	0.00	10.000	2.937	19.858	2.
↪554							
5	0.658	0.792343	0.04	0.001	27.034	20.458	0.
↪198							
6	0.750	0.884620	0.04	0.010	27.034	20.458	0.
↪251							
7	0.778	0.946658	0.04	0.100	27.034	20.458	0.
↪290							
8	0.771	0.973005	0.04	1.000	27.034	20.458	0.
↪463							
9	0.745	0.974043	0.04	10.000	27.034	20.458	2.
↪313							
10	0.567	0.792343	0.08	0.001	40.465	19.995	0.
↪230							
11	0.677	0.884620	0.08	0.010	40.465	19.995	0.
↪242							
12	0.706	0.946658	0.08	0.100	40.465	19.995	0.
↪278							
13	0.664	0.973005	0.08	1.000	40.465	19.995	0.
↪443							
14	0.625	0.974043	0.08	10.000	40.465	19.995	2.
↪165							
15	0.475	0.792343	0.12	0.001	58.590	19.594	0.
↪182							
16	0.583	0.884620	0.12	0.010	58.590	19.594	0.
↪232							
17	0.608	0.946658	0.12	0.100	58.590	19.594	0.
↪269							
18	0.555	0.973005	0.12	1.000	58.590	19.594	0.
↪427							
19	0.519	0.974043	0.12	10.000	58.590	19.594	2.
↪482							
20	0.396	0.792343	0.16	0.001	82.425	20.489	0.
↪203							
21	0.502	0.884620	0.16	0.010	82.425	20.489	0.
↪259							
22	0.509	0.946658	0.16	0.100	82.425	20.489	0.
↪298							
23	0.451	0.973005	0.16	1.000	82.425	20.489	0.
↪478							
24	0.413	0.974043	0.16	10.000	82.425	20.489	3.
↪299							

(continues on next page)

(continued from previous page)

25	0.321	0.792343	0.20	0.001	93.475	21.156	0.
↪181							
26	0.409	0.884620	0.20	0.010	93.475	21.156	0.
↪229							
27	0.411	0.946658	0.20	0.100	93.475	21.156	0.
↪265							
28	0.349	0.973005	0.20	1.000	93.475	21.156	0.
↪425							
29	0.312	0.974043	0.20	10.000	93.475	21.156	2.
↪086							
30	0.263	0.792343	0.24	0.001	94.942	20.544	0.
↪189							
31	0.344	0.884620	0.24	0.010	94.942	20.544	0.
↪232							
32	0.348	0.946658	0.24	0.100	94.942	20.544	0.
↪268							
33	0.297	0.973005	0.24	1.000	94.942	20.544	0.
↪430							
34	0.269	0.974043	0.24	10.000	94.942	20.544	2.
↪112							
35	0.220	0.792343	0.28	0.001	125.194	17.138	0.
↪180							
36	0.273	0.884620	0.28	0.010	125.194	17.138	0.
↪230							
37	0.280	0.946658	0.28	0.100	125.194	17.138	0.
↪265							
38	0.243	0.973005	0.28	1.000	125.194	17.138	0.
↪426							
39	0.232	0.974043	0.28	10.000	125.194	17.138	2.
↪078							
MultinomialNB #1							
	test_mean_accuracy	train_mean_accuracy	p	alpha	time_err	time_pre	time_
↪mod							
0	0.834	0.961583	0.00	0.0001	2.937	19.858	0.
↪035							
1	0.843	0.960675	0.00	0.0010	2.937	19.858	0.
↪032							
2	0.851	0.958209	0.00	0.0100	2.937	19.858	0.
↪032							
3	0.852	0.951979	0.00	0.1000	2.937	19.858	0.
↪032							
4	0.832	0.925892	0.00	1.0000	2.937	19.858	0.
↪032							
5	0.741	0.965737	0.04	0.0001	27.034	20.458	0.
↪046							
6	0.762	0.965217	0.04	0.0010	27.034	20.458	0.
↪040							
7	0.782	0.964698	0.04	0.0100	27.034	20.458	0.
↪039							
8	0.796	0.959507	0.04	0.1000	27.034	20.458	0.
↪039							
9	0.770	0.930305	0.04	1.0000	27.034	20.458	0.
↪039							
10	0.671	0.964309	0.08	0.0001	40.465	19.995	0.
↪083							
11	0.698	0.963920	0.08	0.0010	40.465	19.995	0.
↪078							

(continues on next page)

(continued from previous page)

12	0.730	0.963790	0.08	0.0100	40.465	19.995	0.
↪079							
13	0.750	0.960545	0.08	0.1000	40.465	19.995	0.
↪078							
14	0.687	0.921739	0.08	1.0000	40.465	19.995	0.
↪056							
15	0.631	0.961324	0.12	0.0001	58.590	19.594	0.
↪036							
16	0.659	0.961454	0.12	0.0010	58.590	19.594	0.
↪029							
17	0.691	0.961064	0.12	0.0100	58.590	19.594	0.
↪029							
18	0.705	0.957820	0.12	0.1000	58.590	19.594	0.
↪029							
19	0.629	0.908371	0.12	1.0000	58.590	19.594	0.
↪029							
20	0.551	0.954835	0.16	0.0001	82.425	20.489	0.
↪031							
21	0.584	0.954835	0.16	0.0010	82.425	20.489	0.
↪029							
22	0.614	0.954445	0.16	0.0100	82.425	20.489	0.
↪029							
23	0.631	0.952888	0.16	0.1000	82.425	20.489	0.
↪029							
24	0.562	0.882154	0.16	1.0000	82.425	20.489	0.
↪029							
25	0.475	0.944452	0.20	0.0001	93.475	21.156	0.
↪025							
26	0.500	0.944322	0.20	0.0010	93.475	21.156	0.
↪023							
27	0.532	0.943933	0.20	0.0100	93.475	21.156	0.
↪023							
28	0.553	0.939779	0.20	0.1000	93.475	21.156	0.
↪023							
29	0.475	0.852304	0.20	1.0000	93.475	21.156	0.
↪023							
30	0.434	0.930435	0.24	0.0001	94.942	20.544	0.
↪026							
31	0.441	0.930565	0.24	0.0010	94.942	20.544	0.
↪020							
32	0.479	0.930954	0.24	0.0100	94.942	20.544	0.
↪020							
33	0.499	0.927709	0.24	0.1000	94.942	20.544	0.
↪020							
34	0.419	0.814666	0.24	1.0000	94.942	20.544	0.
↪020							
35	0.365	0.910578	0.28	0.0001	125.194	17.138	0.
↪018							
36	0.385	0.910578	0.28	0.0010	125.194	17.138	0.
↪017							
37	0.404	0.909799	0.28	0.0100	125.194	17.138	0.
↪017							
38	0.419	0.905775	0.28	0.1000	125.194	17.138	0.
↪017							
39	0.367	0.771317	0.28	1.0000	125.194	17.138	0.
↪017							

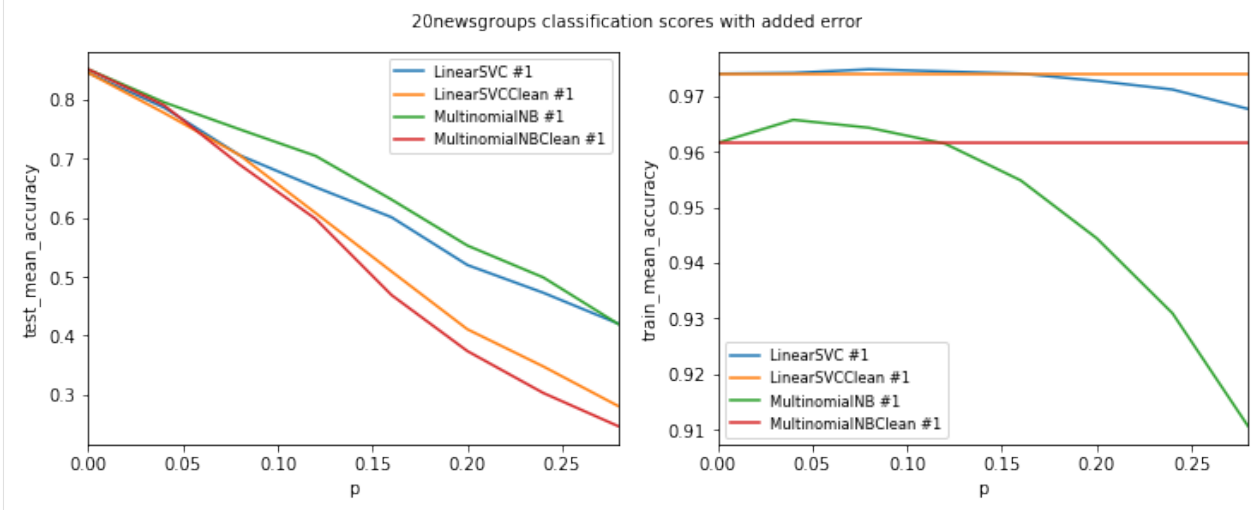
MultinomialNBClean #1							
	test_mean_accuracy	train_mean_accuracy	p	alpha	time_err	time_pre	time_
↪mod							
0	0.834	0.961583	0.00	0.0001	2.937	19.858	0.
↪035							
1	0.843	0.960675	0.00	0.0010	2.937	19.858	0.
↪033							
2	0.851	0.958209	0.00	0.0100	2.937	19.858	0.
↪033							
3	0.852	0.951979	0.00	0.1000	2.937	19.858	0.
↪032							
4	0.832	0.925892	0.00	1.0000	2.937	19.858	0.
↪032							
5	0.707	0.961583	0.04	0.0001	27.034	20.458	0.
↪035							
6	0.732	0.960675	0.04	0.0010	27.034	20.458	0.
↪033							
7	0.769	0.958209	0.04	0.0100	27.034	20.458	0.
↪033							
8	0.791	0.951979	0.04	0.1000	27.034	20.458	0.
↪032							
9	0.782	0.925892	0.04	1.0000	27.034	20.458	0.
↪033							
10	0.520	0.961583	0.08	0.0001	40.465	19.995	0.
↪033							
11	0.559	0.960675	0.08	0.0010	40.465	19.995	0.
↪038							
12	0.596	0.958209	0.08	0.0100	40.465	19.995	0.
↪045							
13	0.665	0.951979	0.08	0.1000	40.465	19.995	0.
↪053							
14	0.690	0.925892	0.08	1.0000	40.465	19.995	0.
↪062							
15	0.393	0.961583	0.12	0.0001	58.590	19.594	0.
↪030							
16	0.417	0.960675	0.12	0.0010	58.590	19.594	0.
↪028							
17	0.455	0.958209	0.12	0.0100	58.590	19.594	0.
↪027							
18	0.522	0.951979	0.12	0.1000	58.590	19.594	0.
↪028							
19	0.598	0.925892	0.12	1.0000	58.590	19.594	0.
↪028							
20	0.295	0.961583	0.16	0.0001	82.425	20.489	0.
↪034							
21	0.304	0.960675	0.16	0.0010	82.425	20.489	0.
↪031							
22	0.326	0.958209	0.16	0.0100	82.425	20.489	0.
↪031							
23	0.382	0.951979	0.16	0.1000	82.425	20.489	0.
↪031							
24	0.469	0.925892	0.16	1.0000	82.425	20.489	0.
↪031							
25	0.235	0.961583	0.20	0.0001	93.475	21.156	0.
↪030							
26	0.242	0.960675	0.20	0.0010	93.475	21.156	0.
↪028							

(continues on next page)

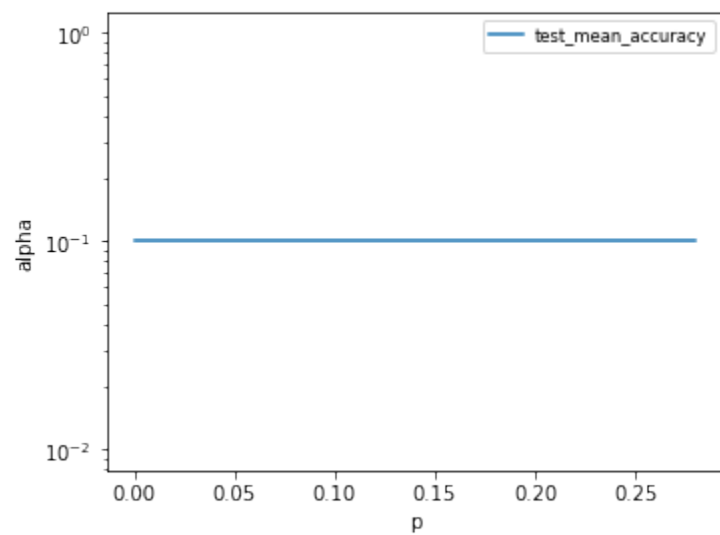
(continued from previous page)

27	0.246	0.958209	0.20	0.0100	93.475	21.156	0.
↪ 028							
28	0.293	0.951979	0.20	0.1000	93.475	21.156	0.
↪ 028							
29	0.374	0.925892	0.20	1.0000	93.475	21.156	0.
↪ 028							
30	0.189	0.961583	0.24	0.0001	94.942	20.544	0.
↪ 031							
31	0.195	0.960675	0.24	0.0010	94.942	20.544	0.
↪ 028							
32	0.214	0.958209	0.24	0.0100	94.942	20.544	0.
↪ 028							
33	0.242	0.951979	0.24	0.1000	94.942	20.544	0.
↪ 028							
34	0.303	0.925892	0.24	1.0000	94.942	20.544	0.
↪ 028							
35	0.174	0.961583	0.28	0.0001	125.194	17.138	0.
↪ 031							
36	0.175	0.960675	0.28	0.0010	125.194	17.138	0.
↪ 028							
37	0.182	0.958209	0.28	0.0100	125.194	17.138	0.
↪ 028							
38	0.211	0.951979	0.28	0.1000	125.194	17.138	0.
↪ 028							
39	0.246	0.925892	0.28	1.0000	125.194	17.138	0.
↪ 028							

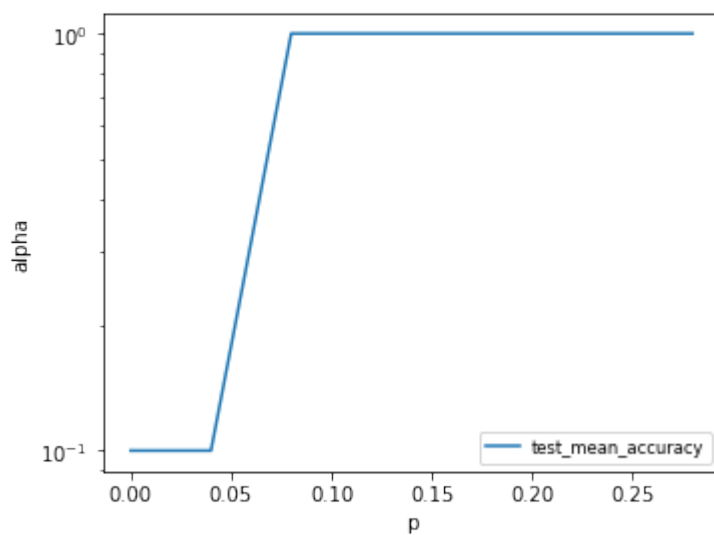
```
/wrk/users/thalvari/dpEmu/dpemu/plotting_utils.py:299: RuntimeWarning: More than 20
figures have been opened. Figures created through the pyplot interface (`matplotlib.
pyplot.figure`) are retained until explicitly closed and may consume too much
memory. (To control this warning, see the rcParam `figure.max_open_warning`).
fig, ax = plt.subplots(figsize=(10, 8))
```



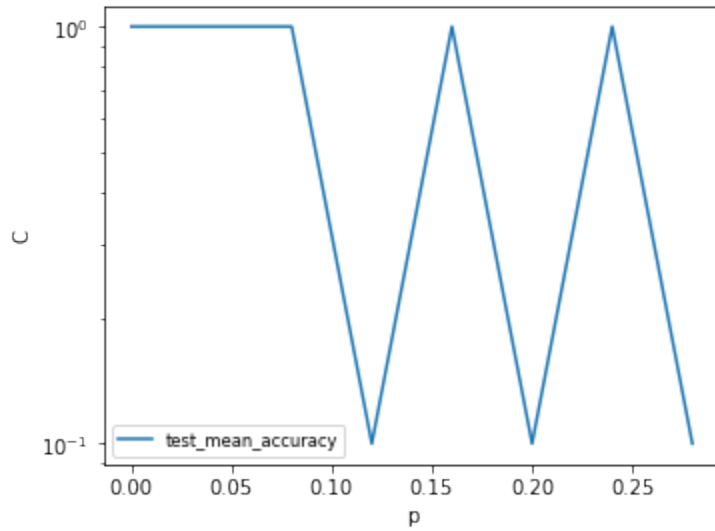
Best parameters for 20newsgroups classification (MultinomialNB #1)



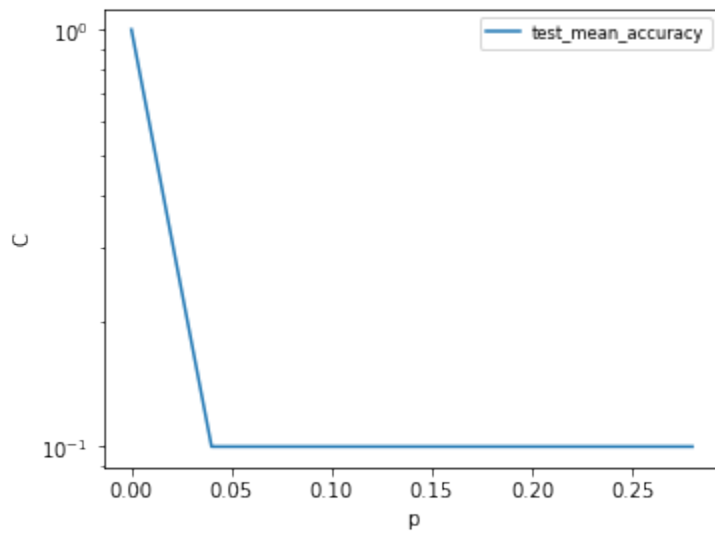
Best parameters for 20newsgroups classification (MultinomialNBClean #1)



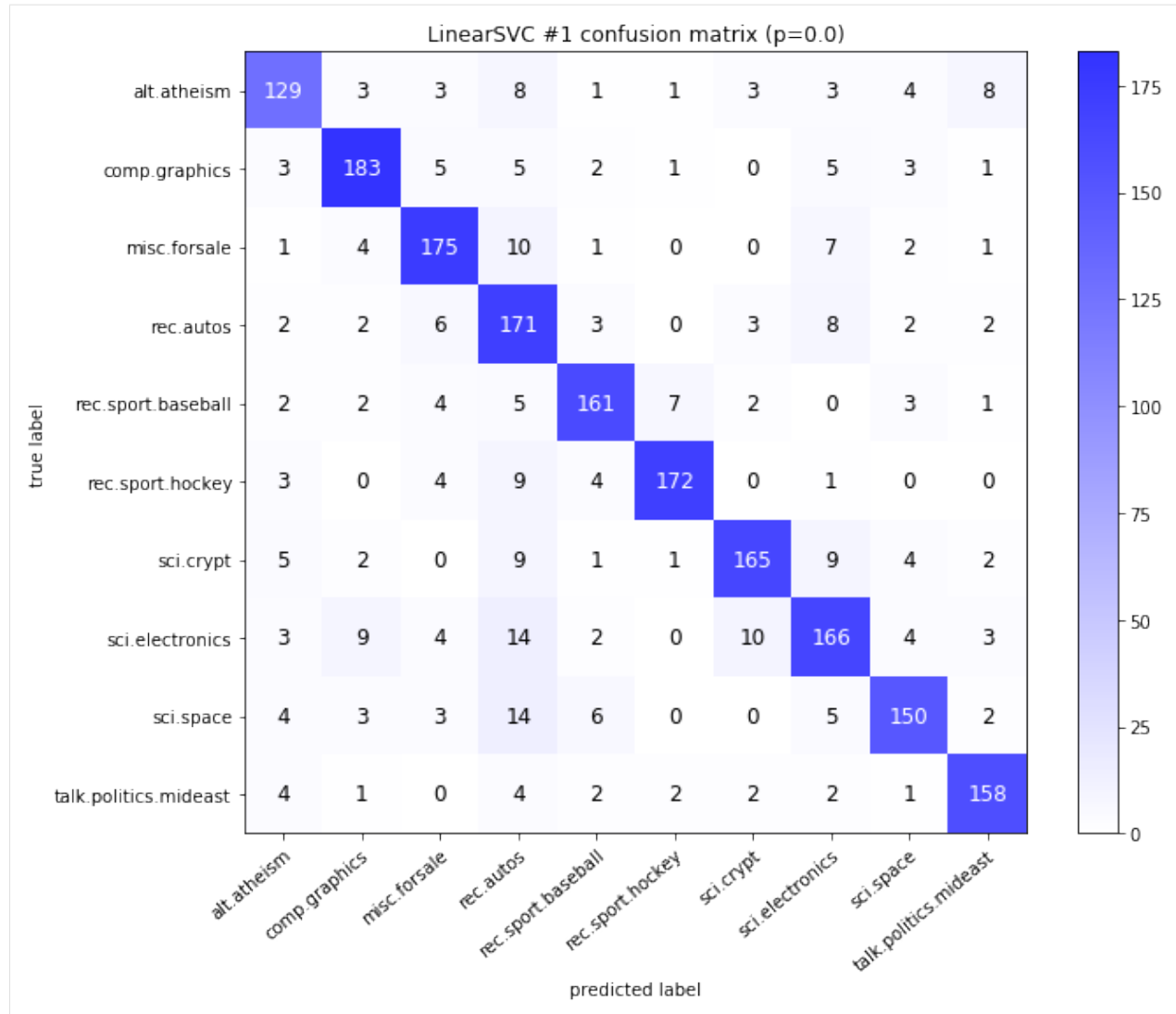
Best parameters for 20newsgroups classification (LinearSVC #1)

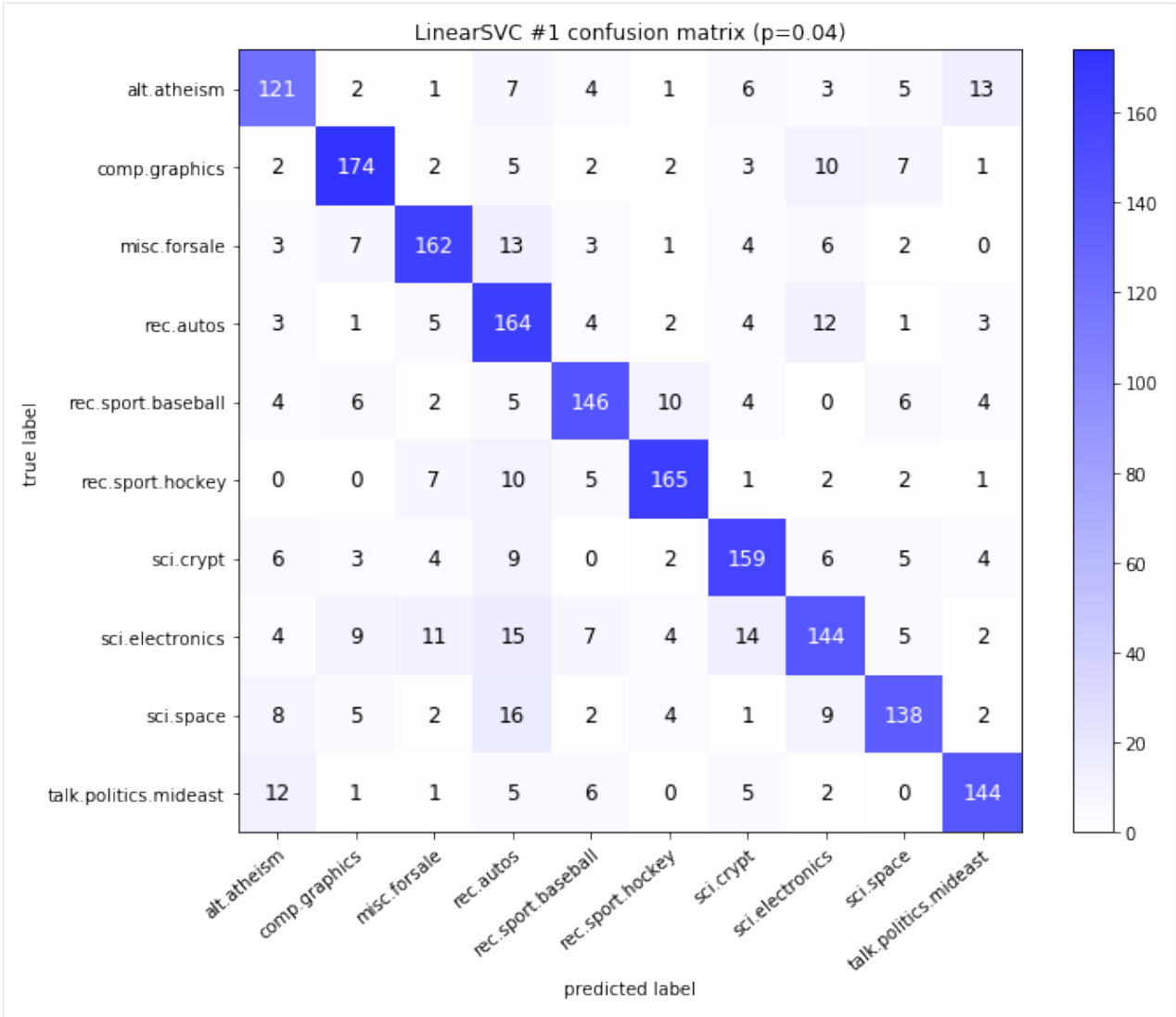


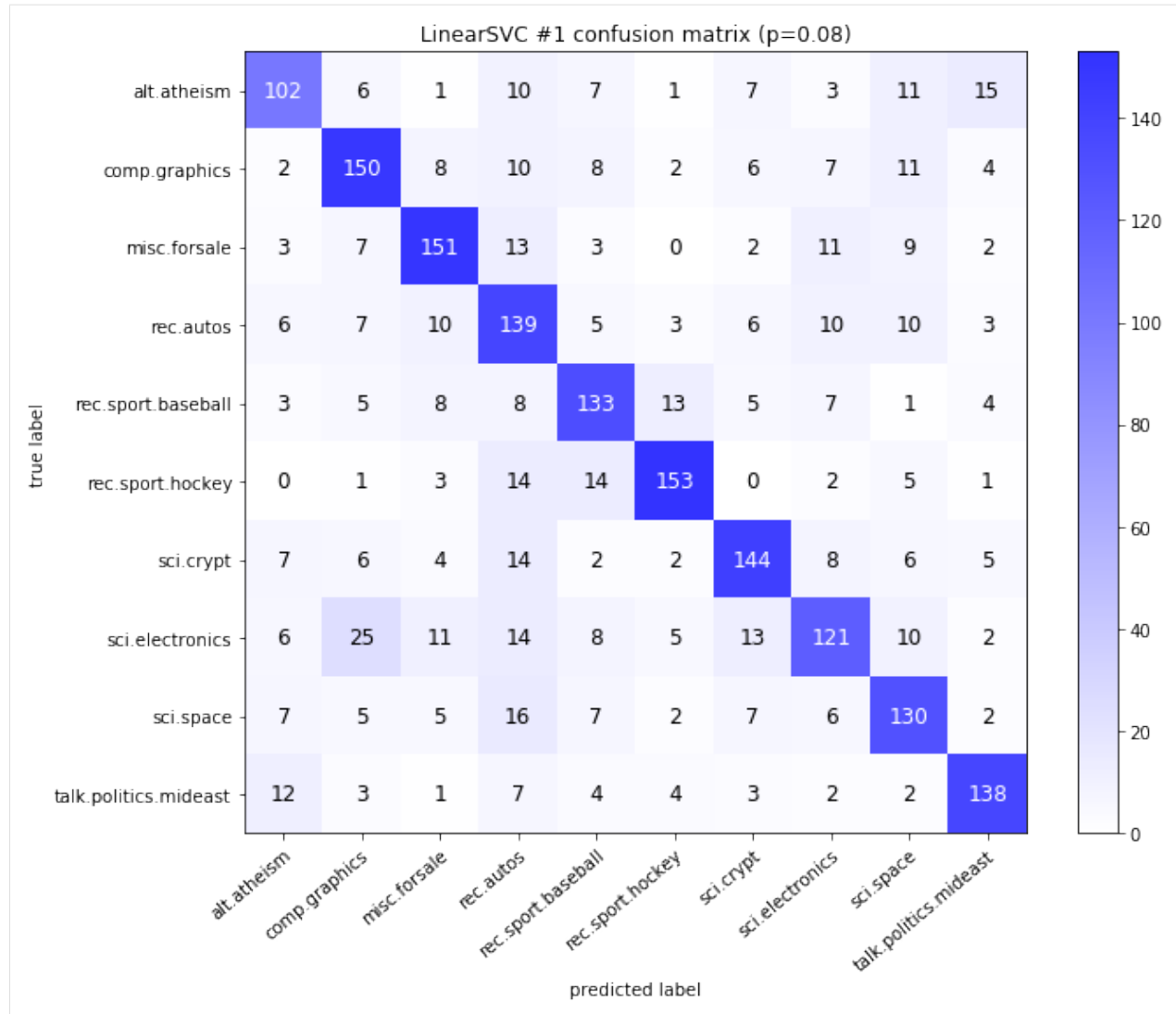
Best parameters for 20newsgroups classification (LinearSVCClean #1)

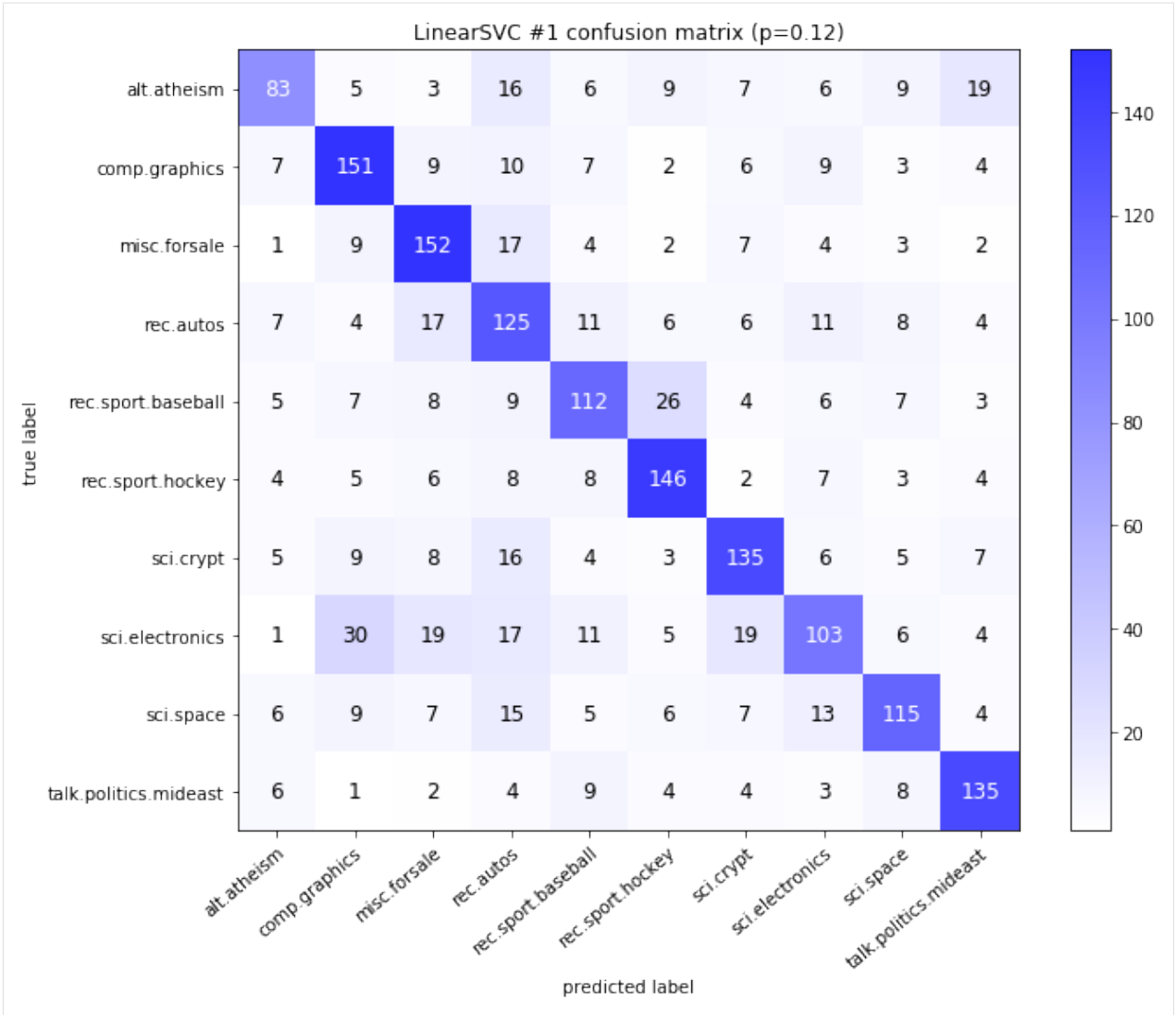


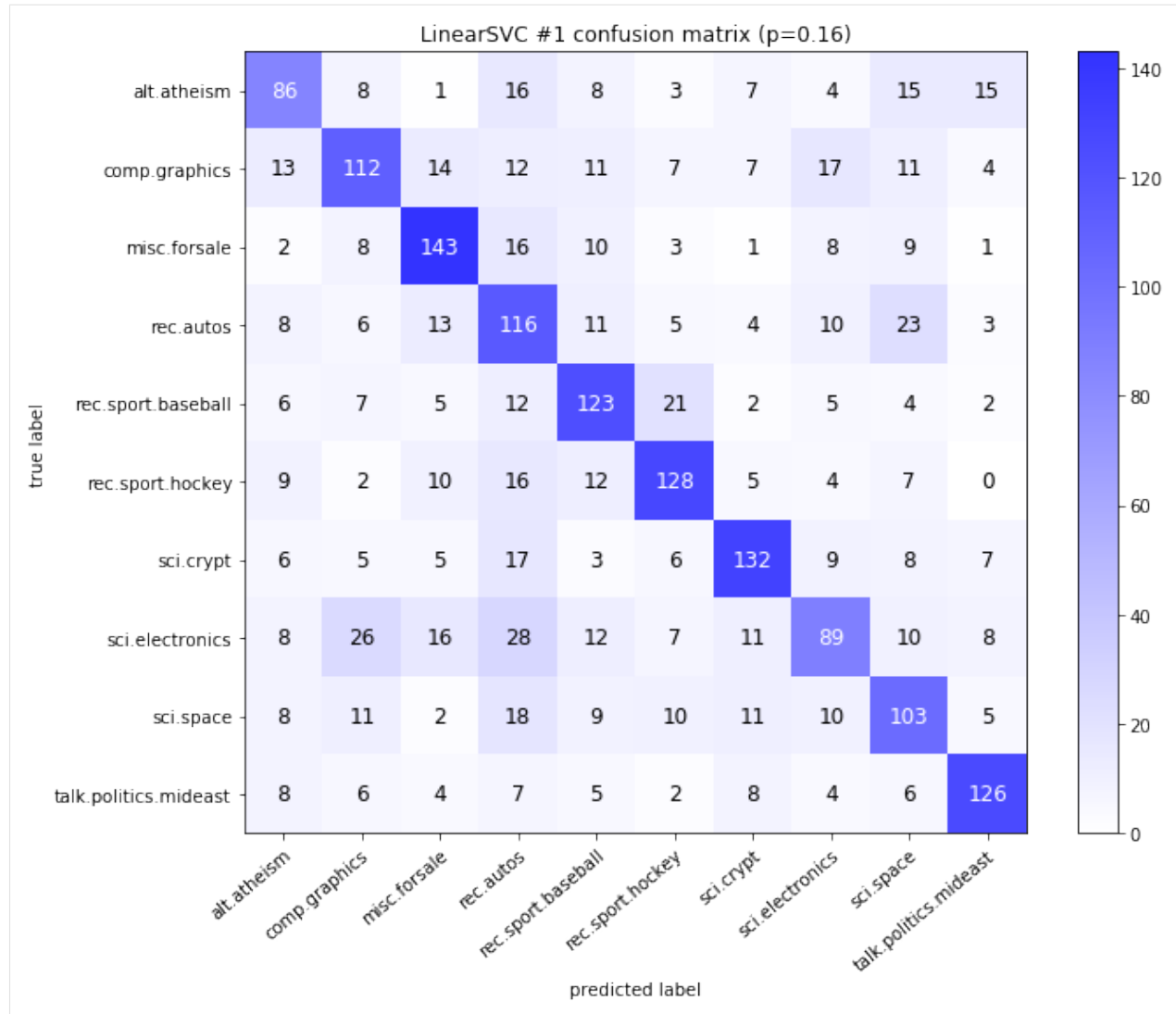


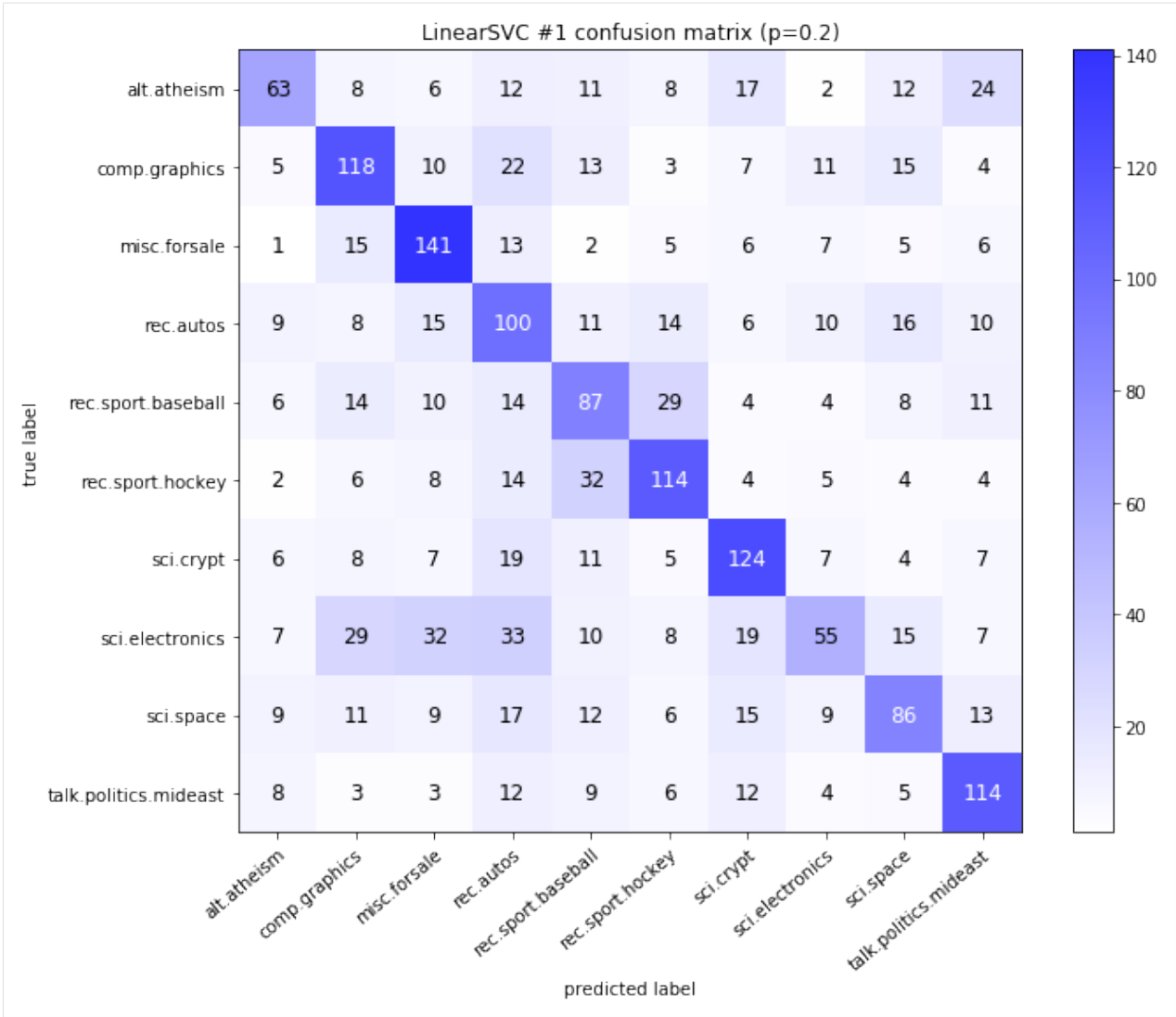


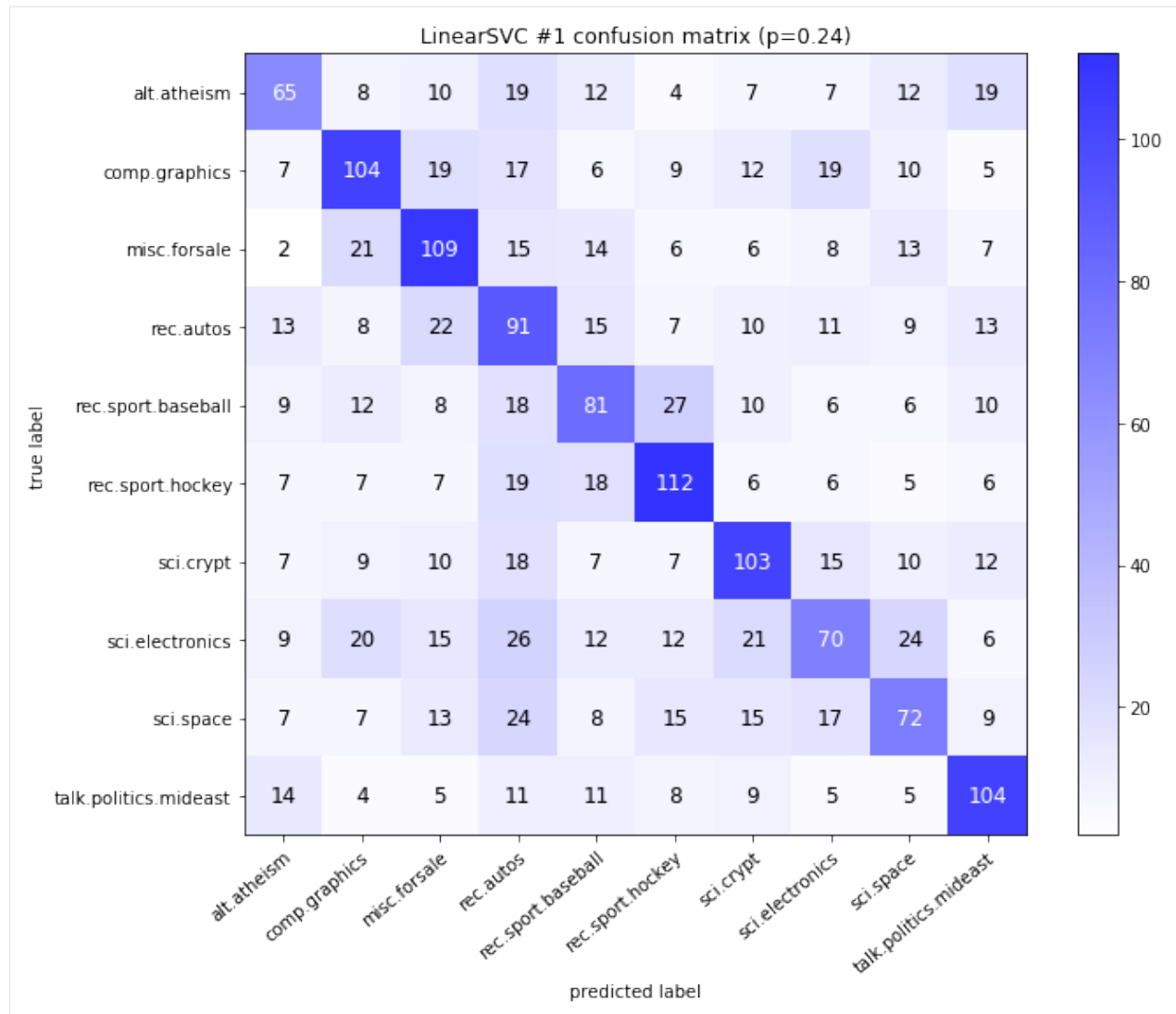


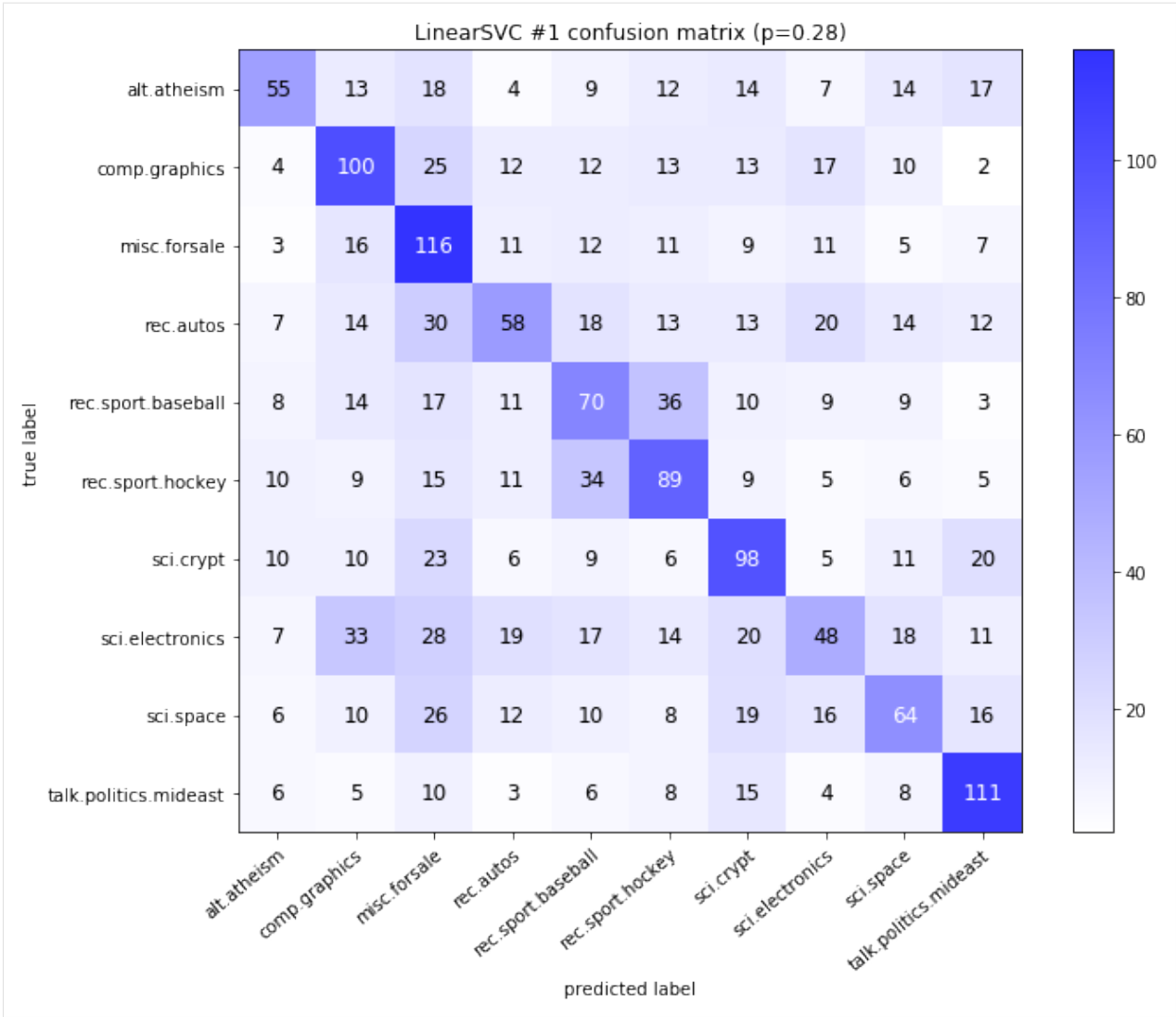


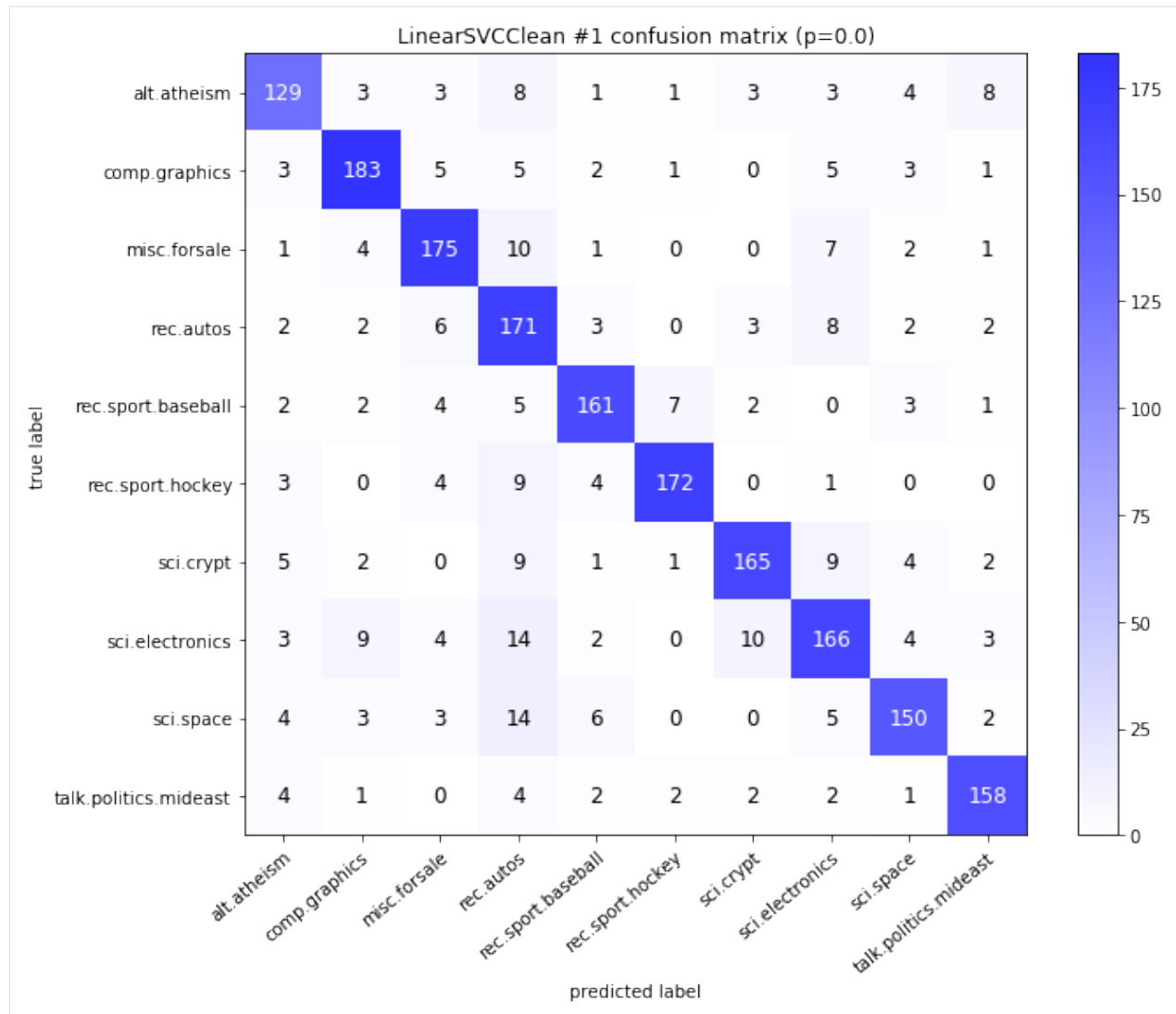


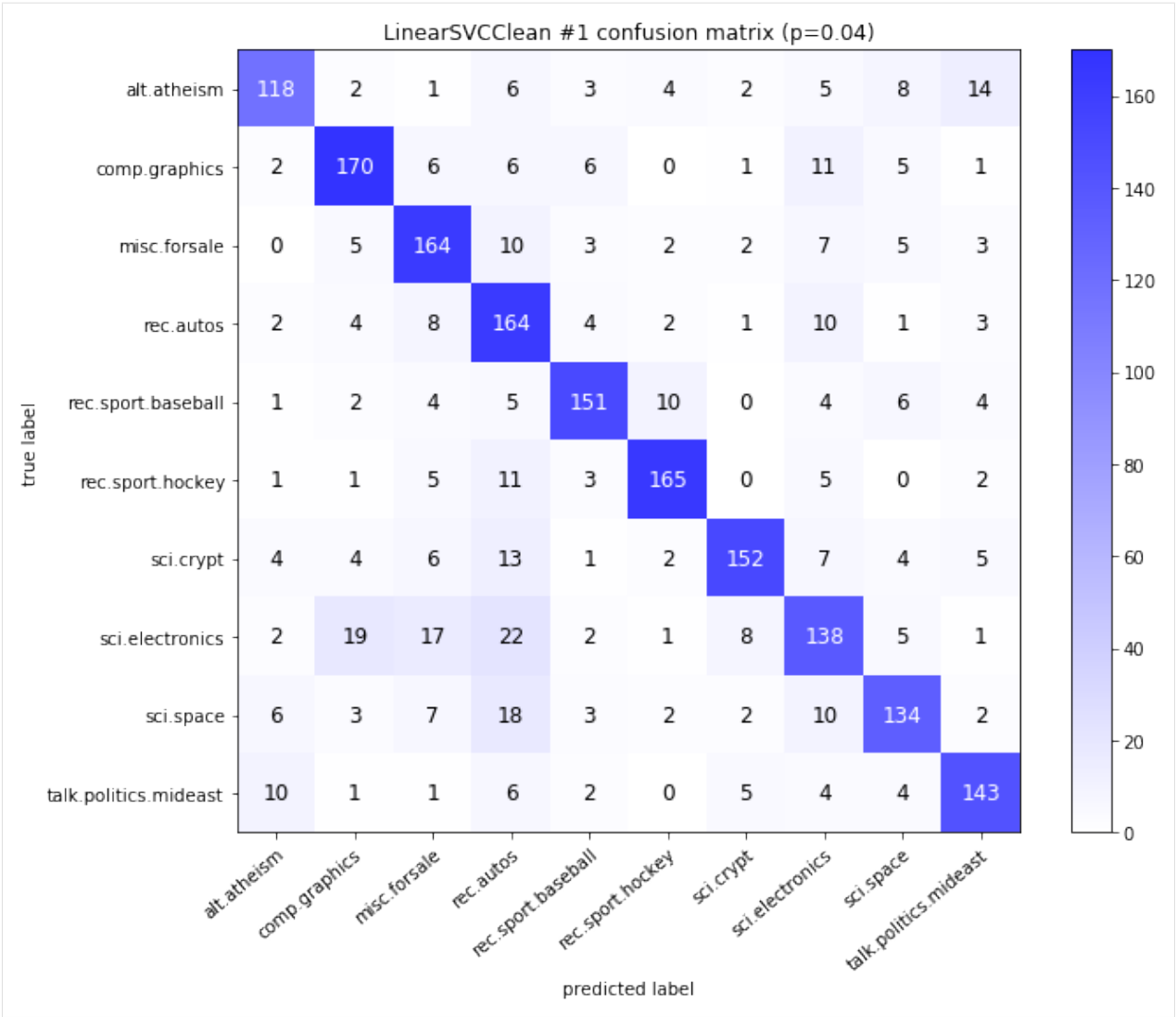


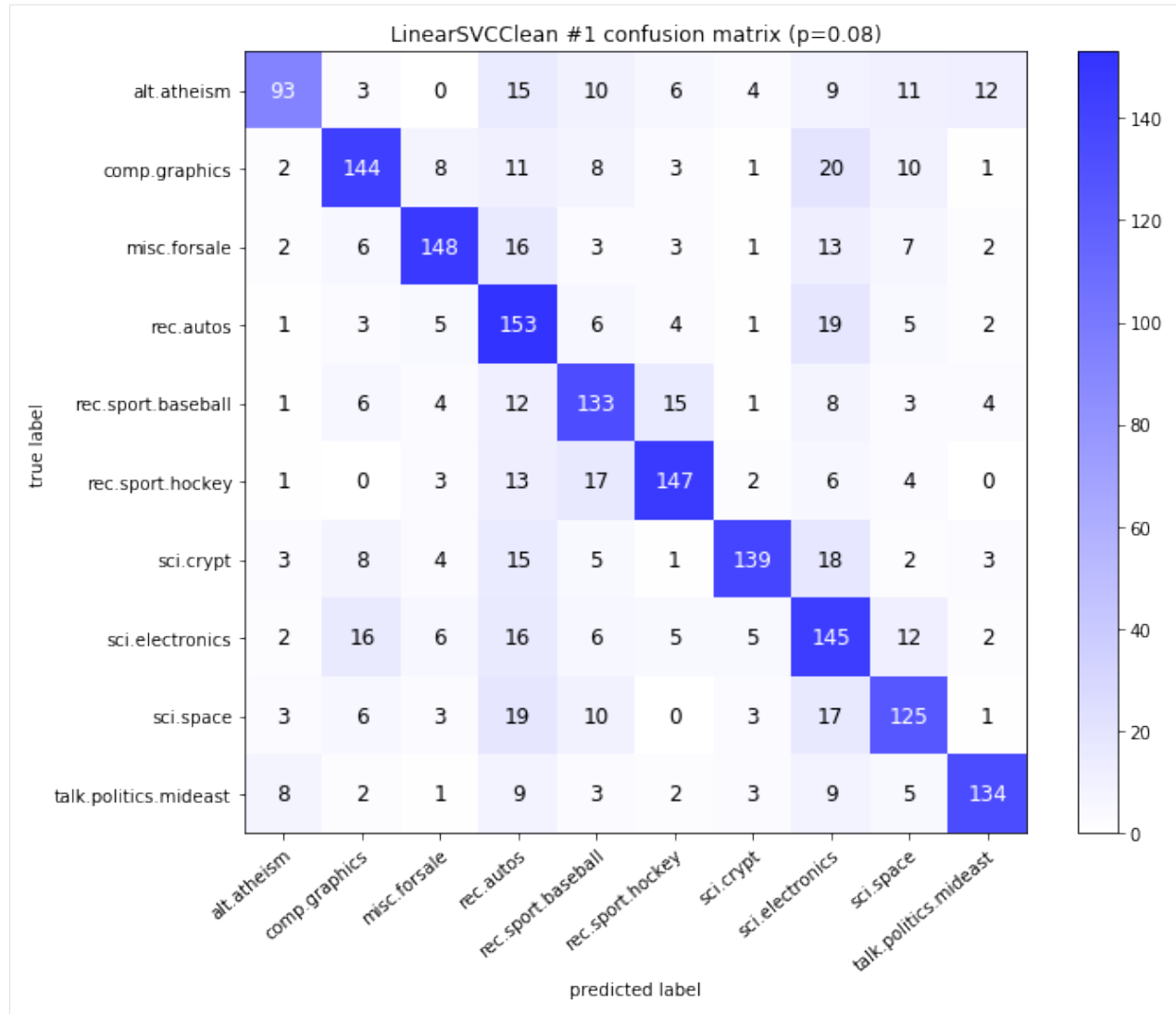


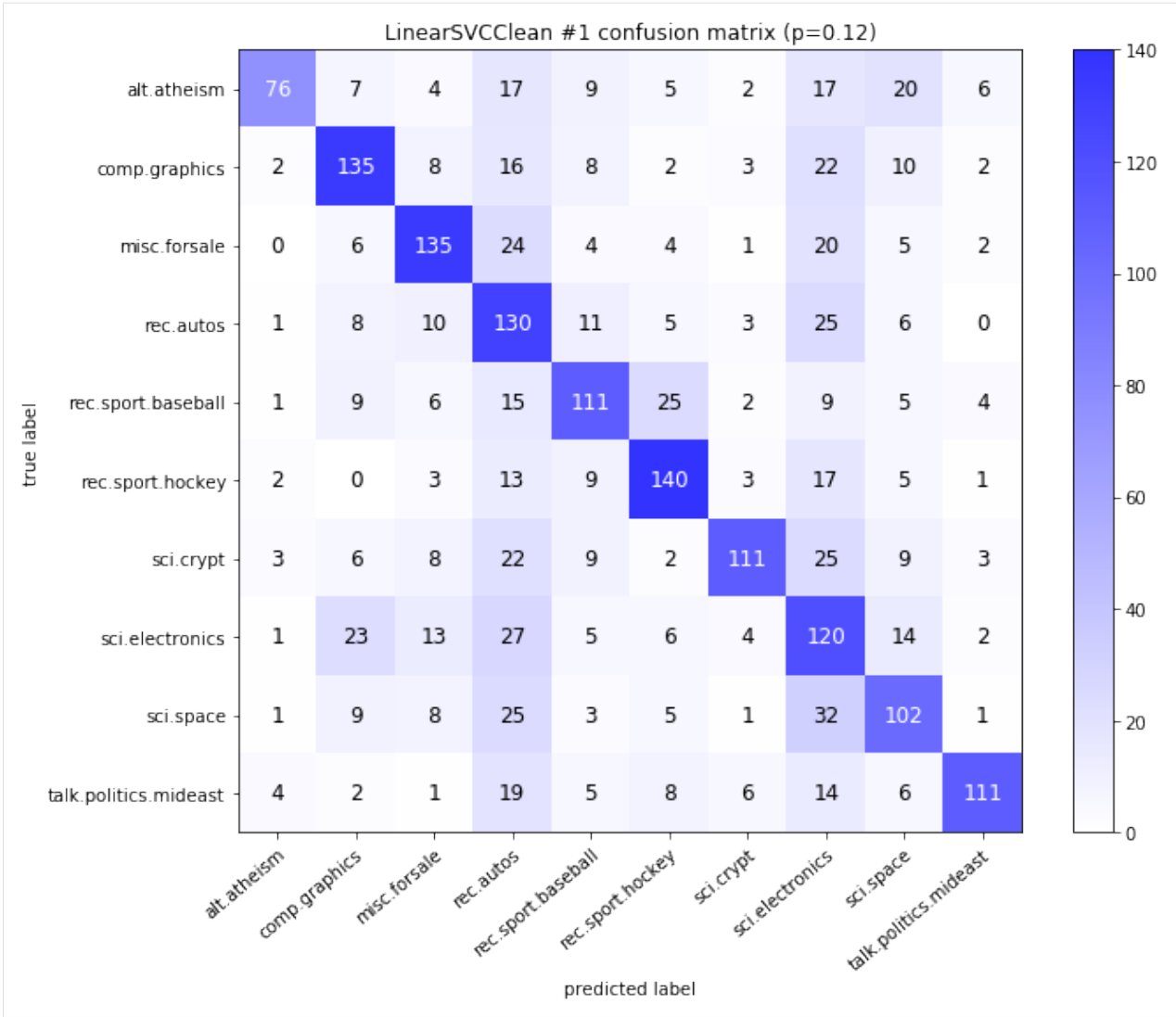


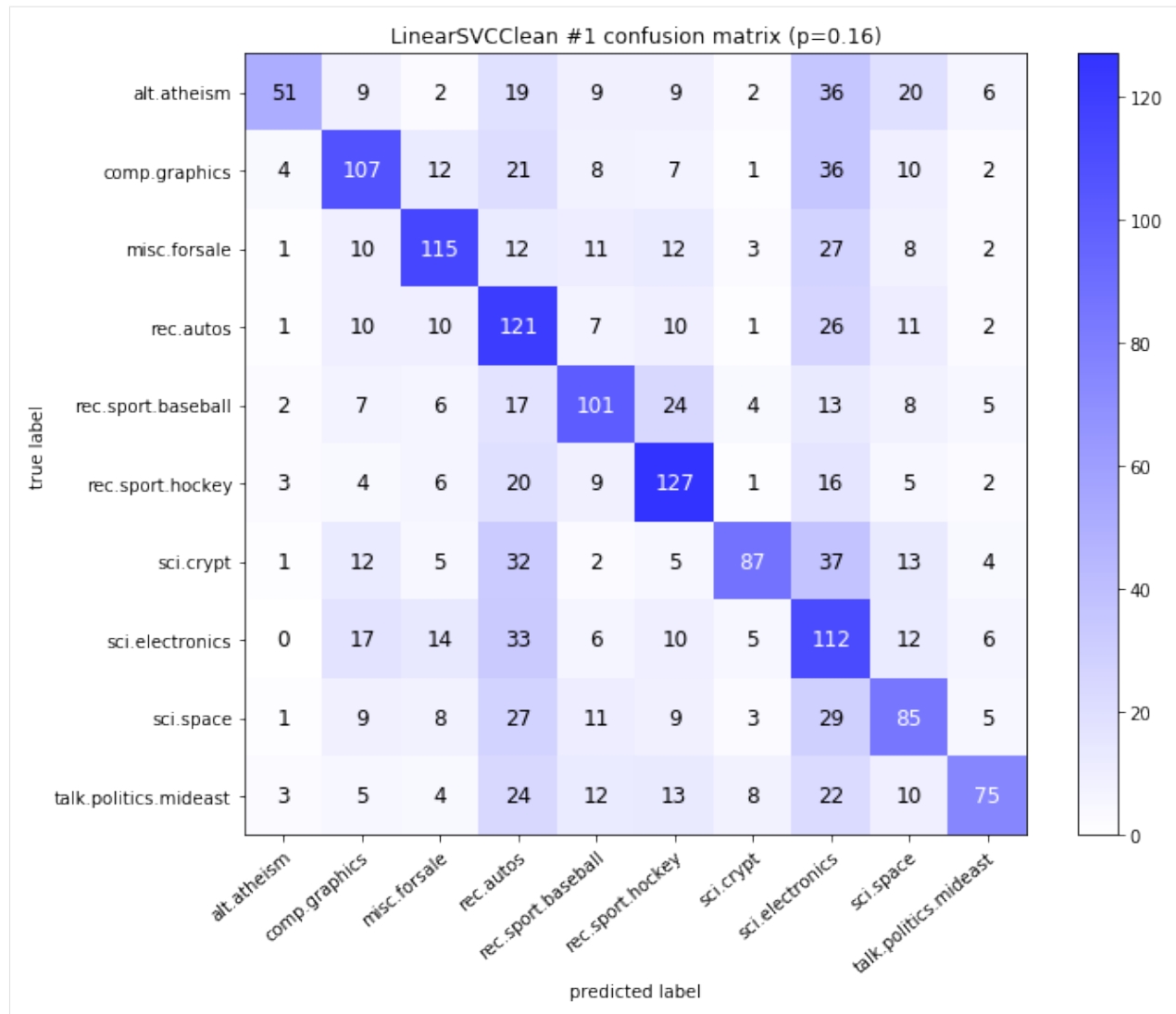


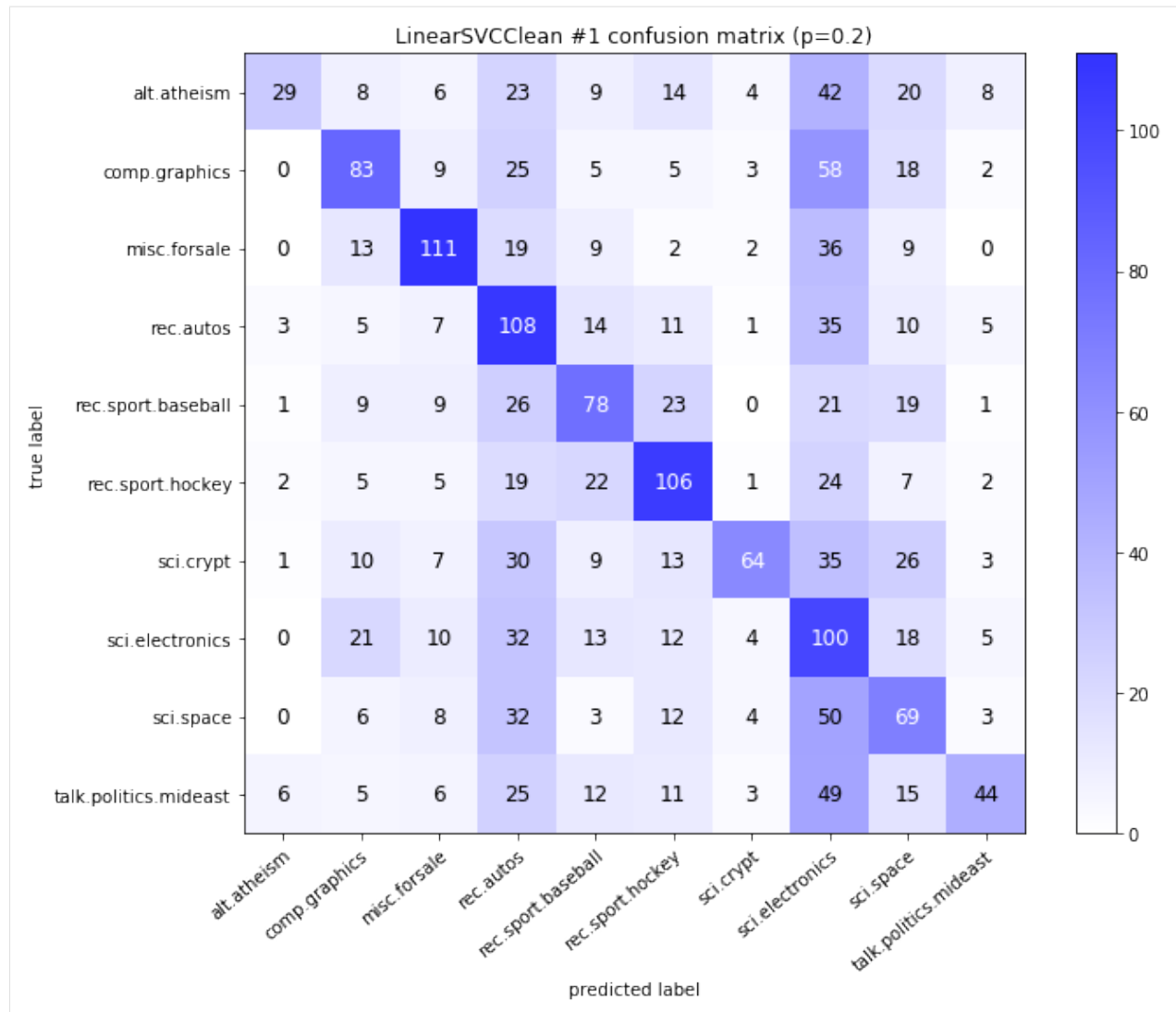


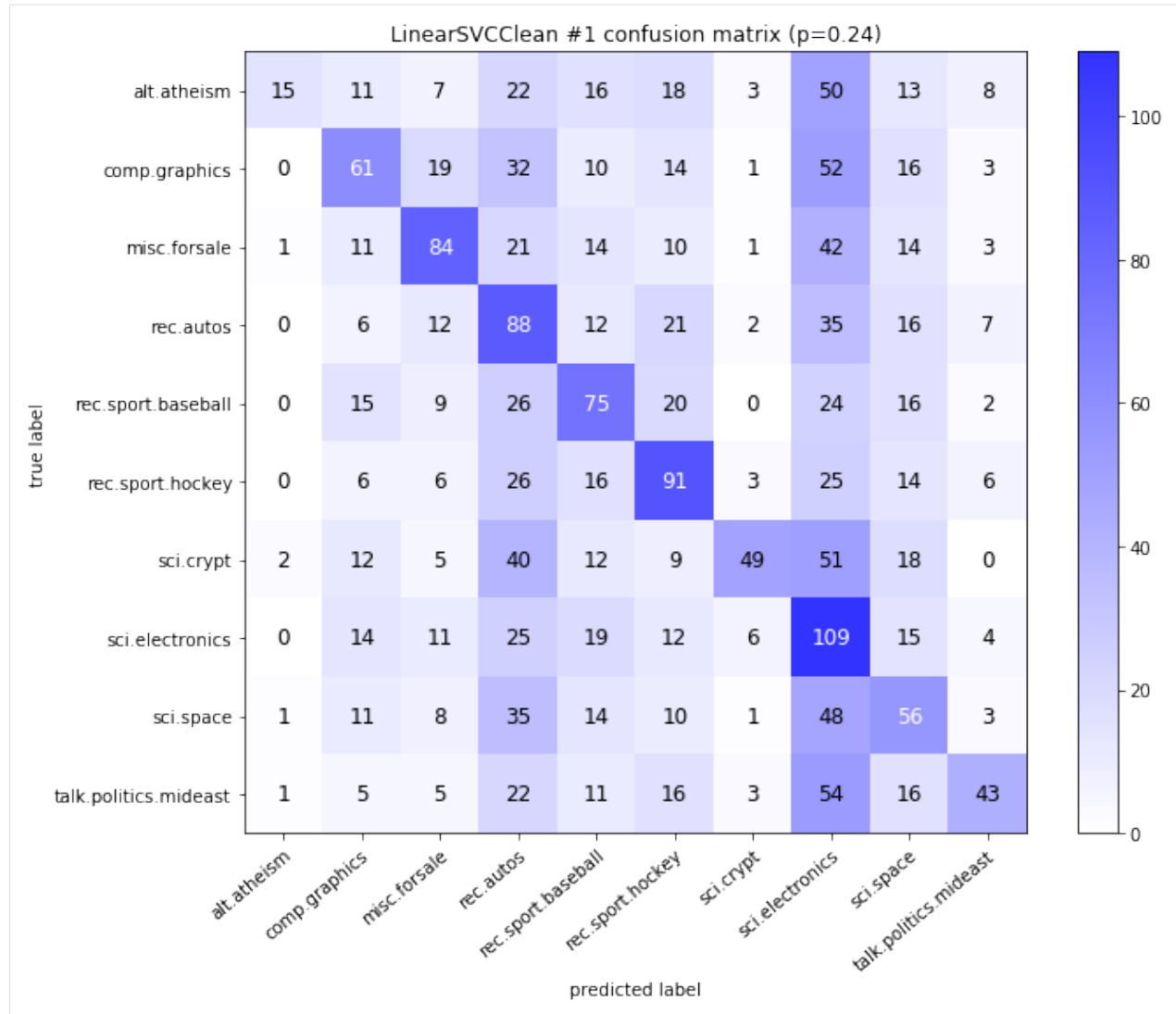


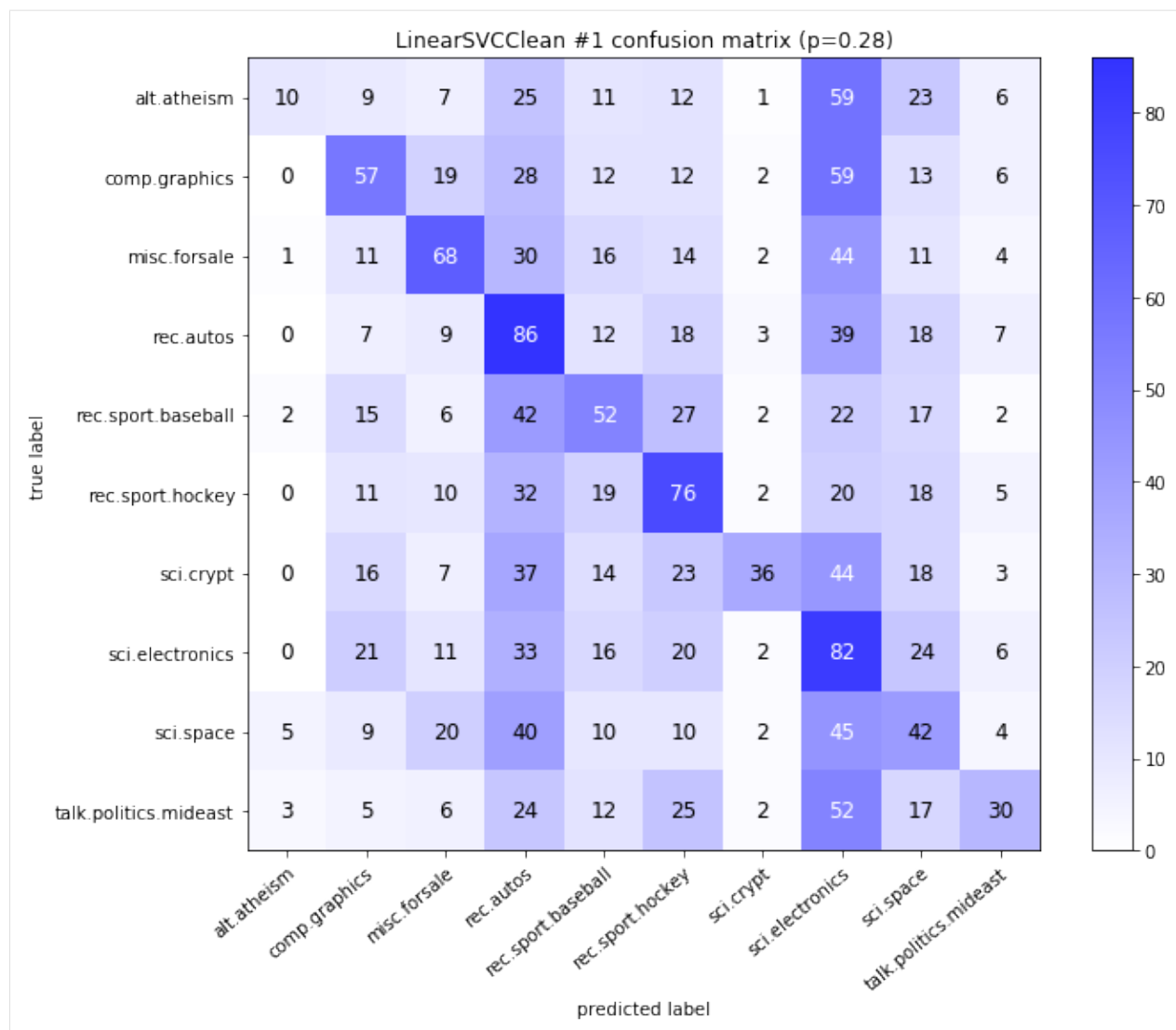


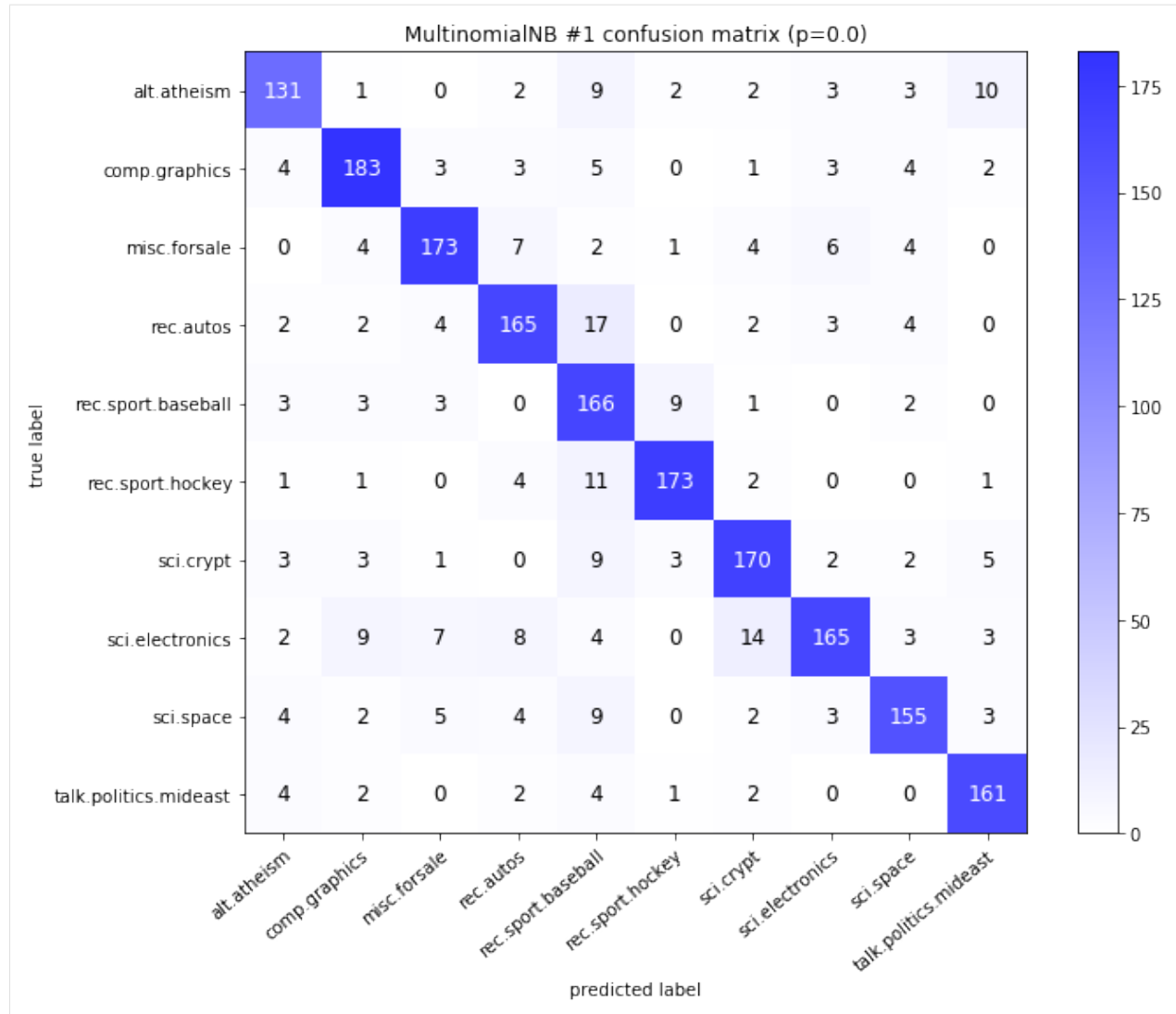


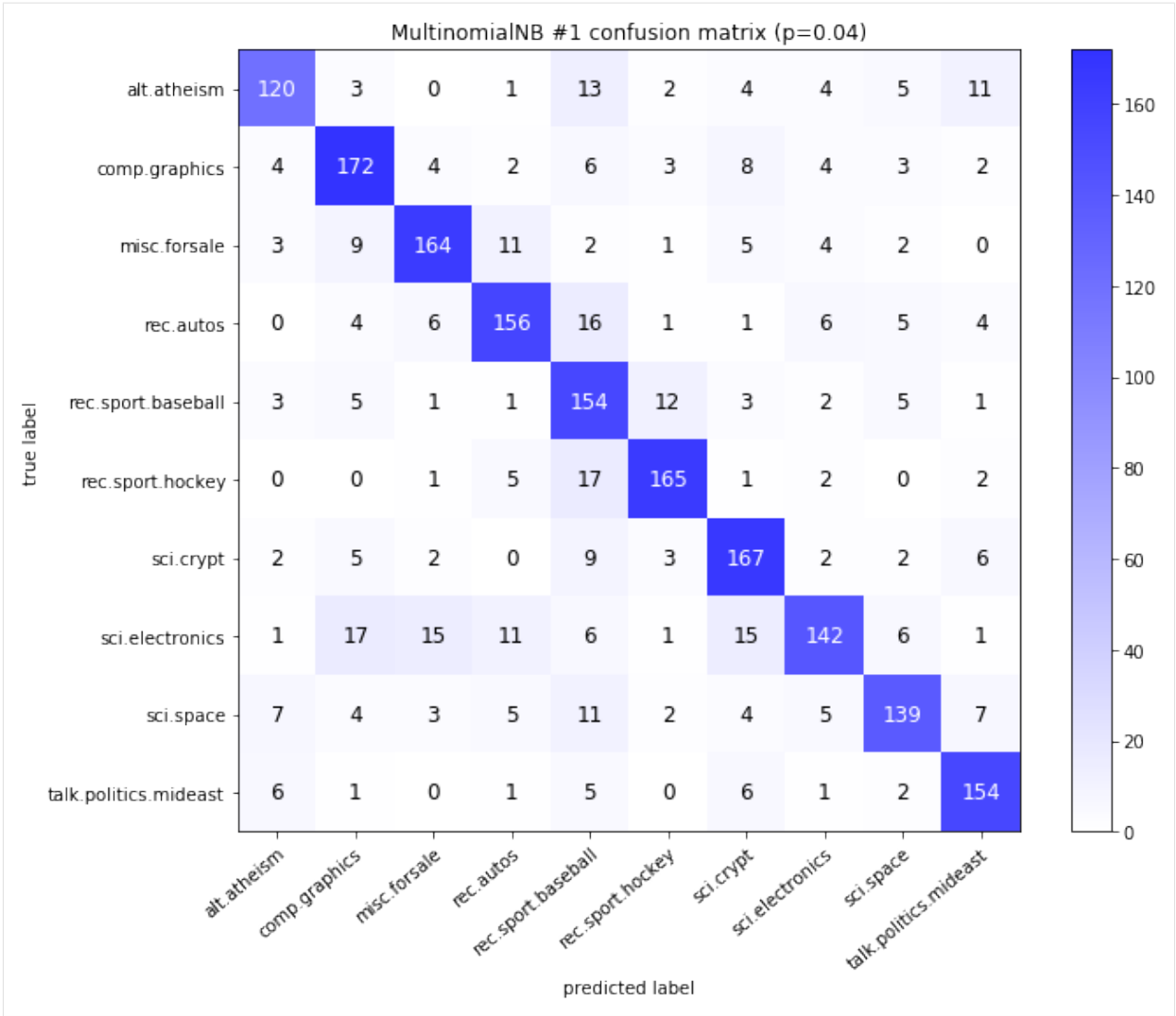


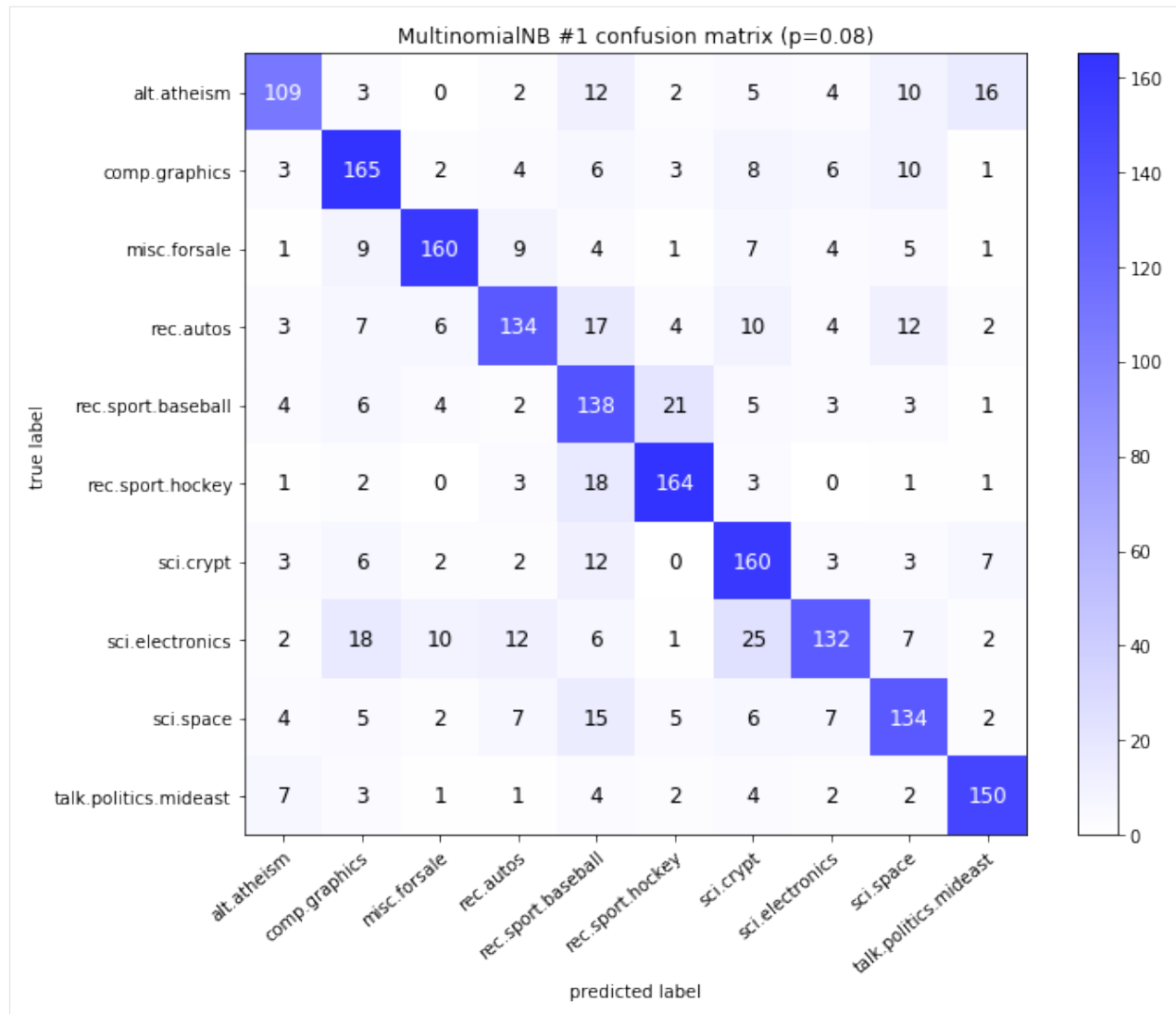


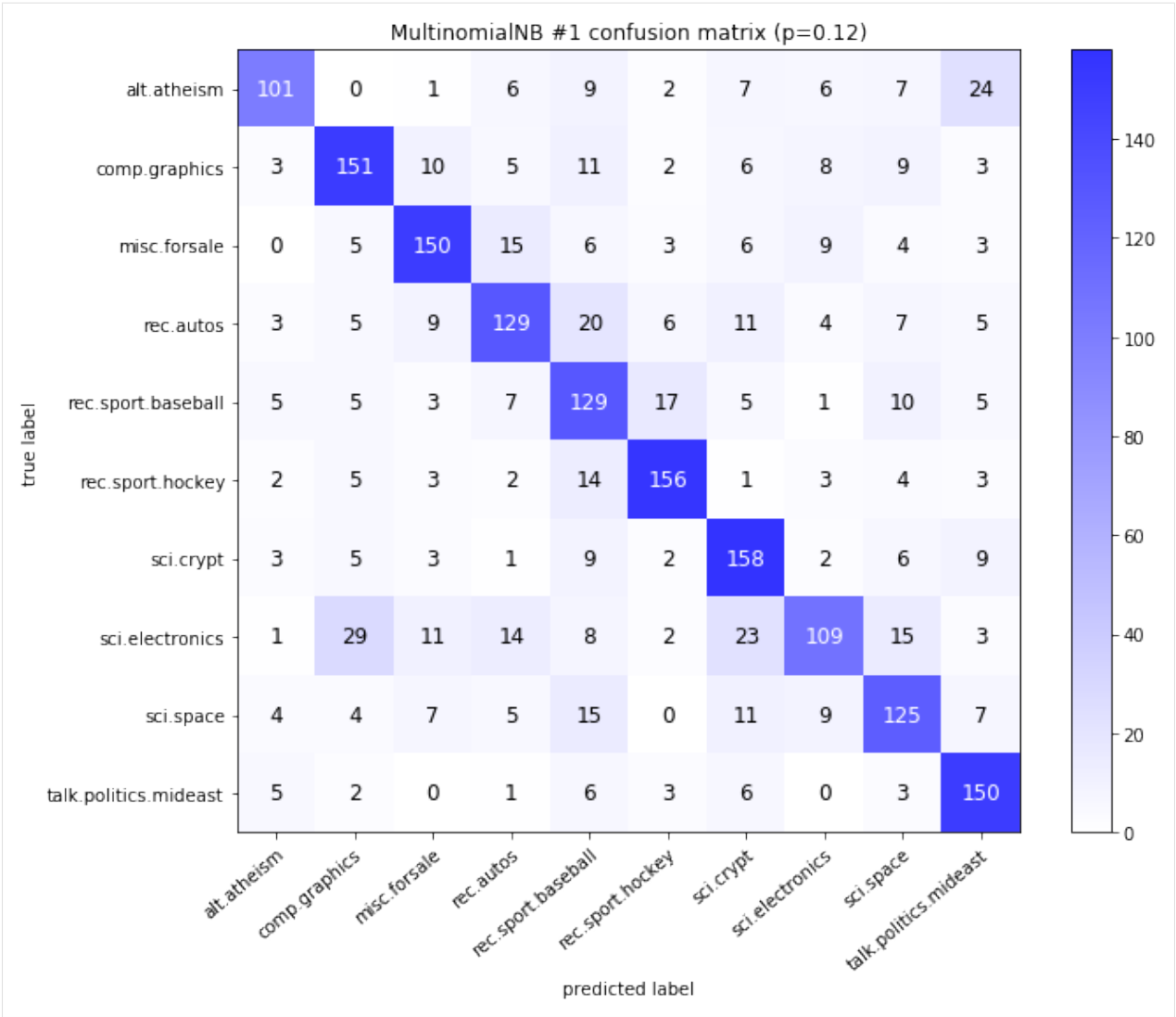


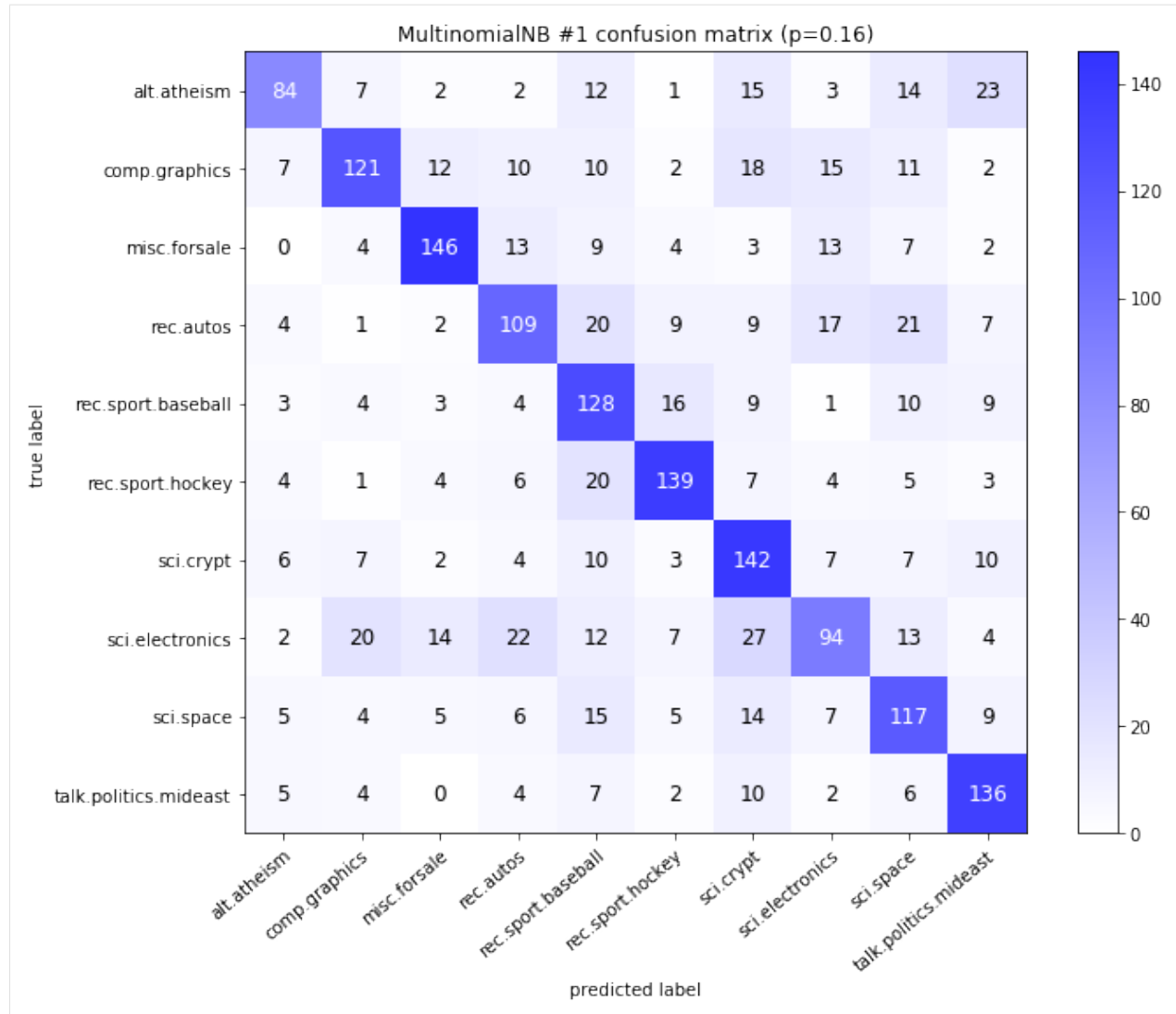


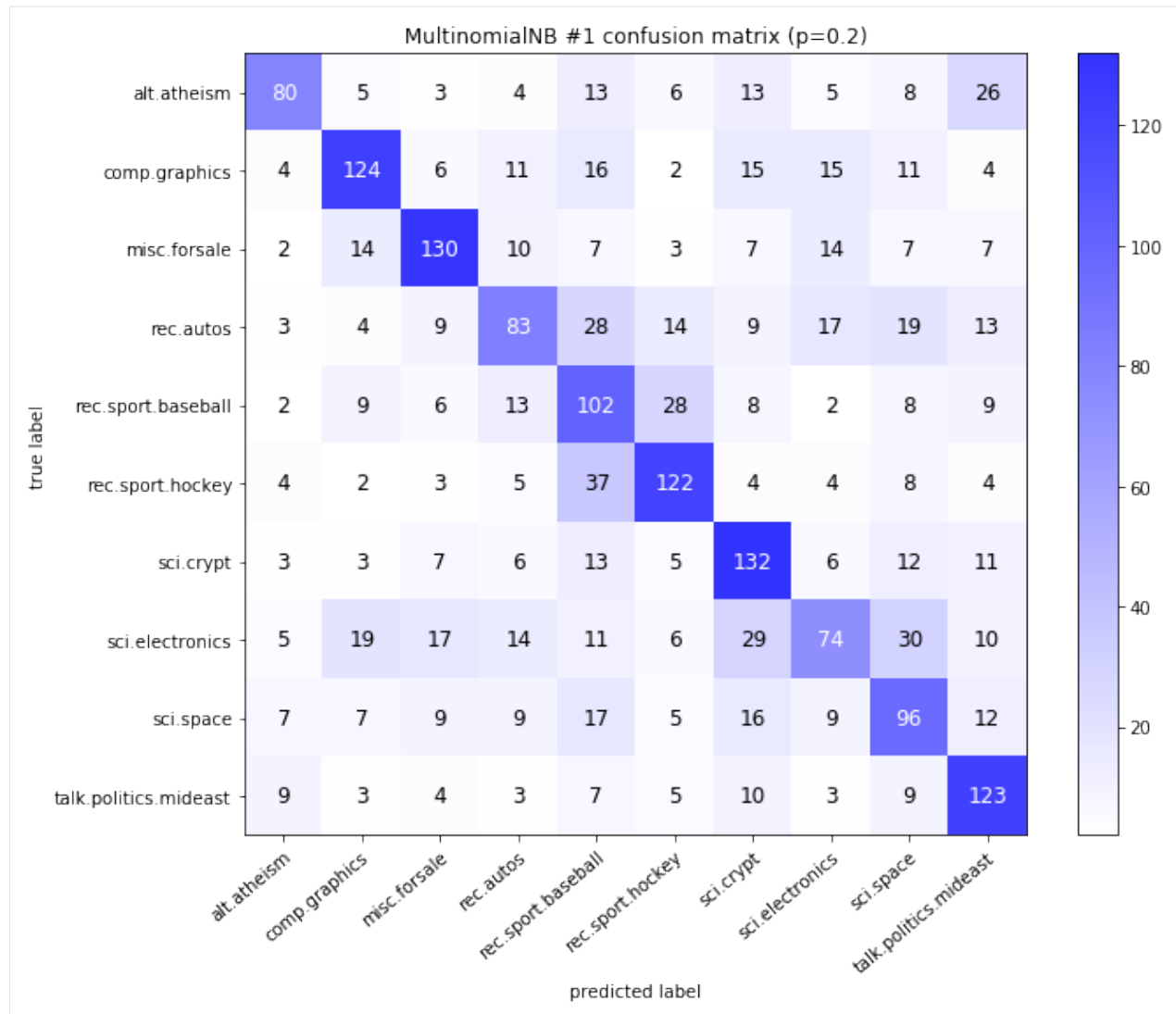


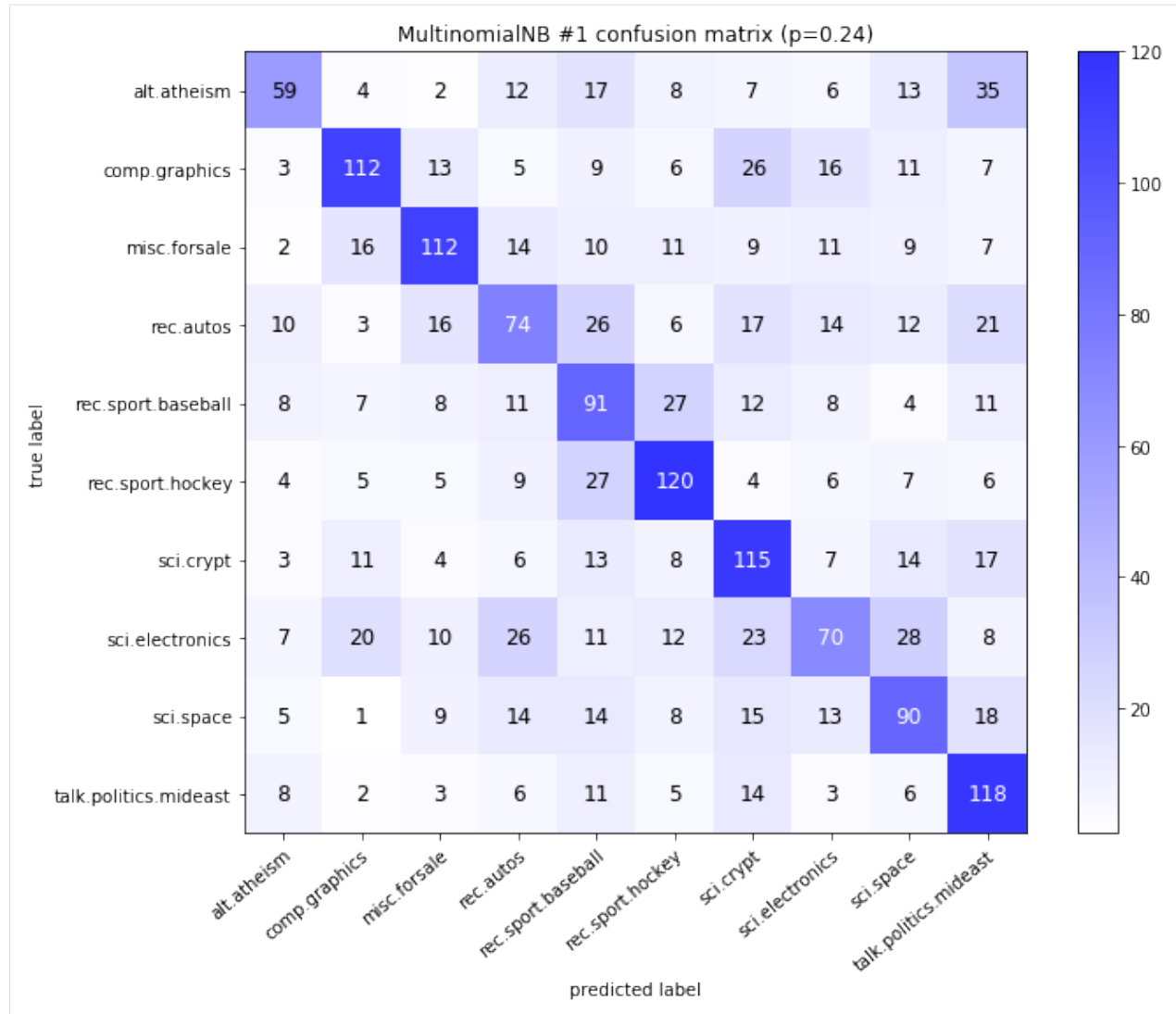


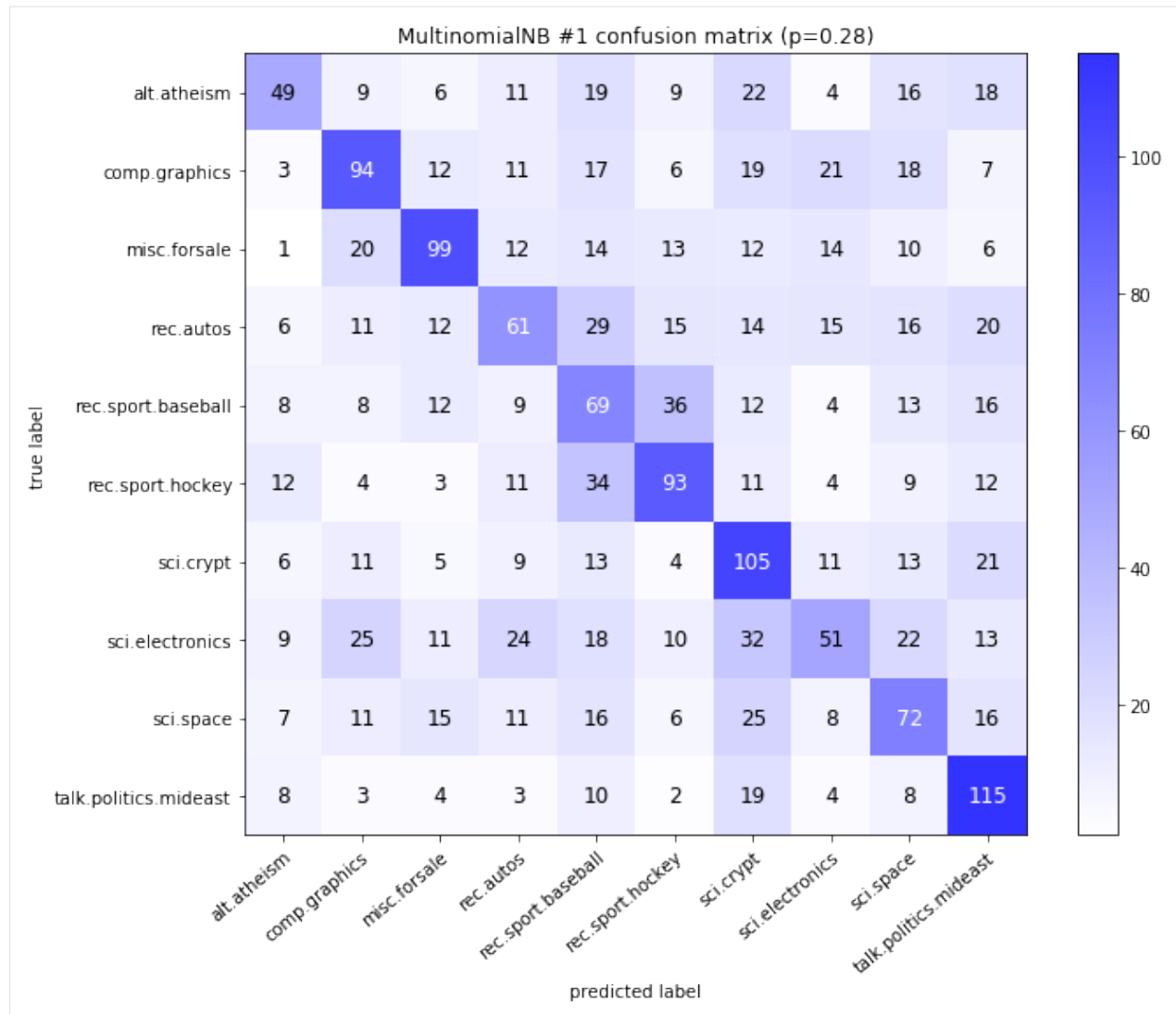


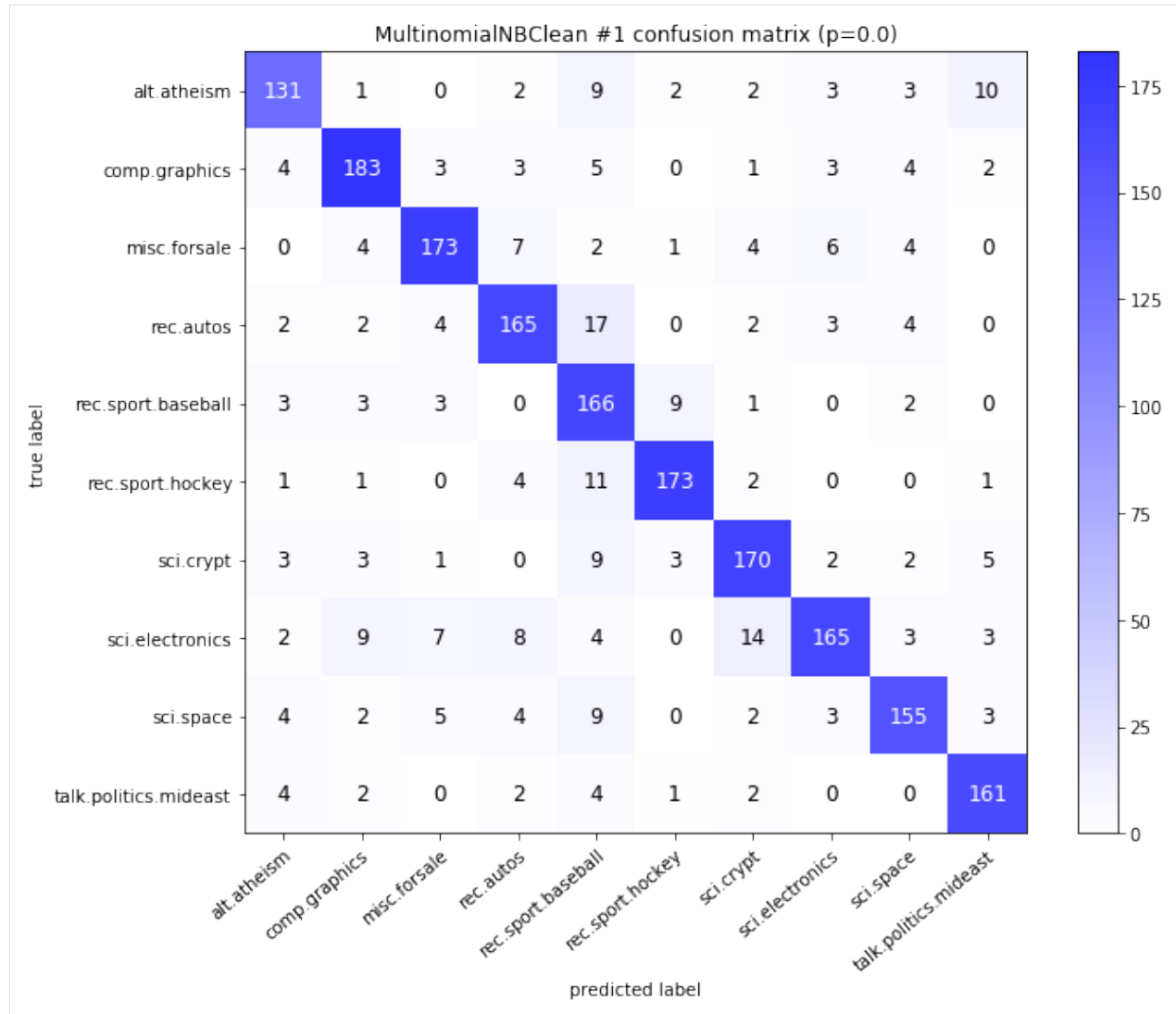


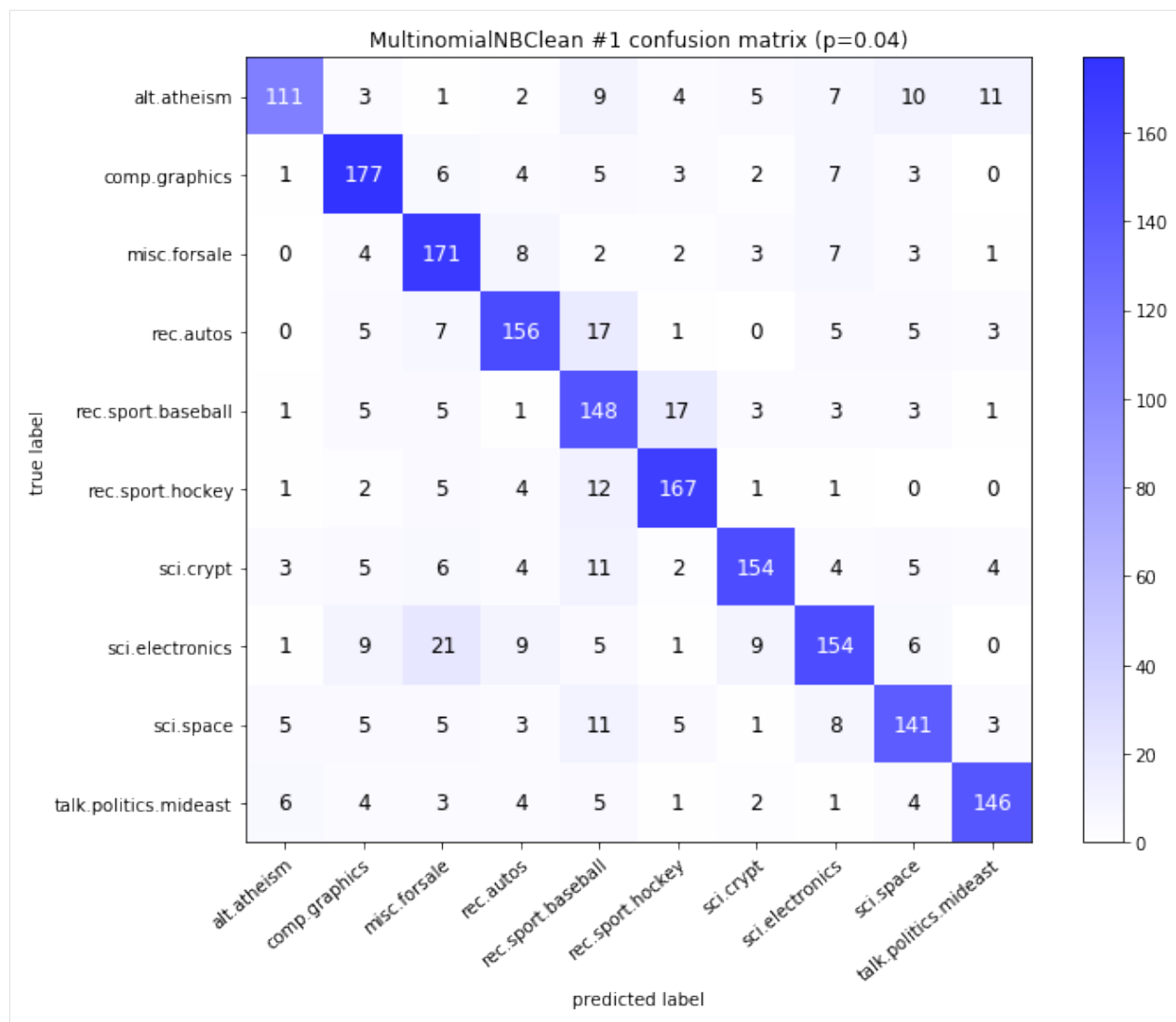


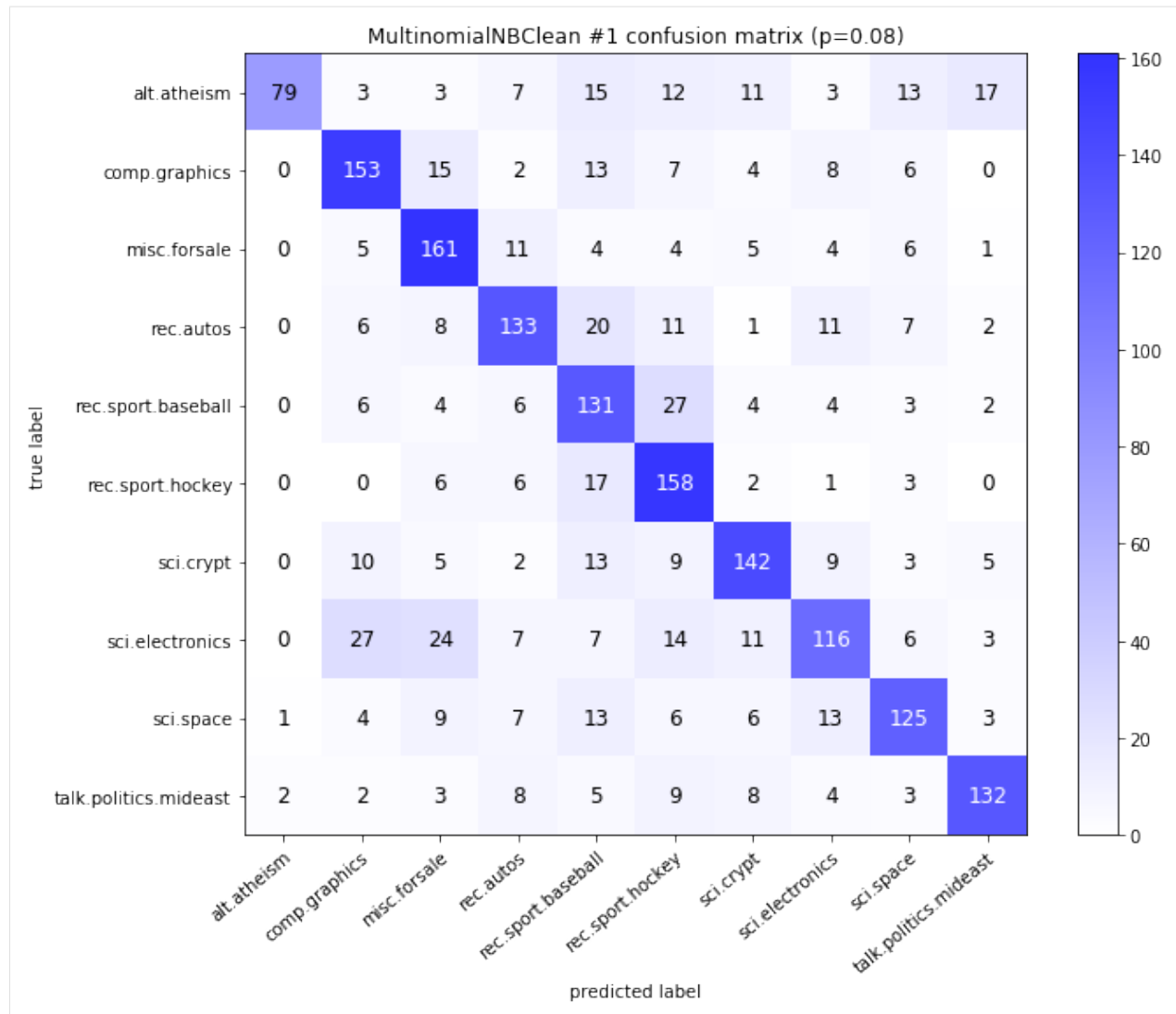


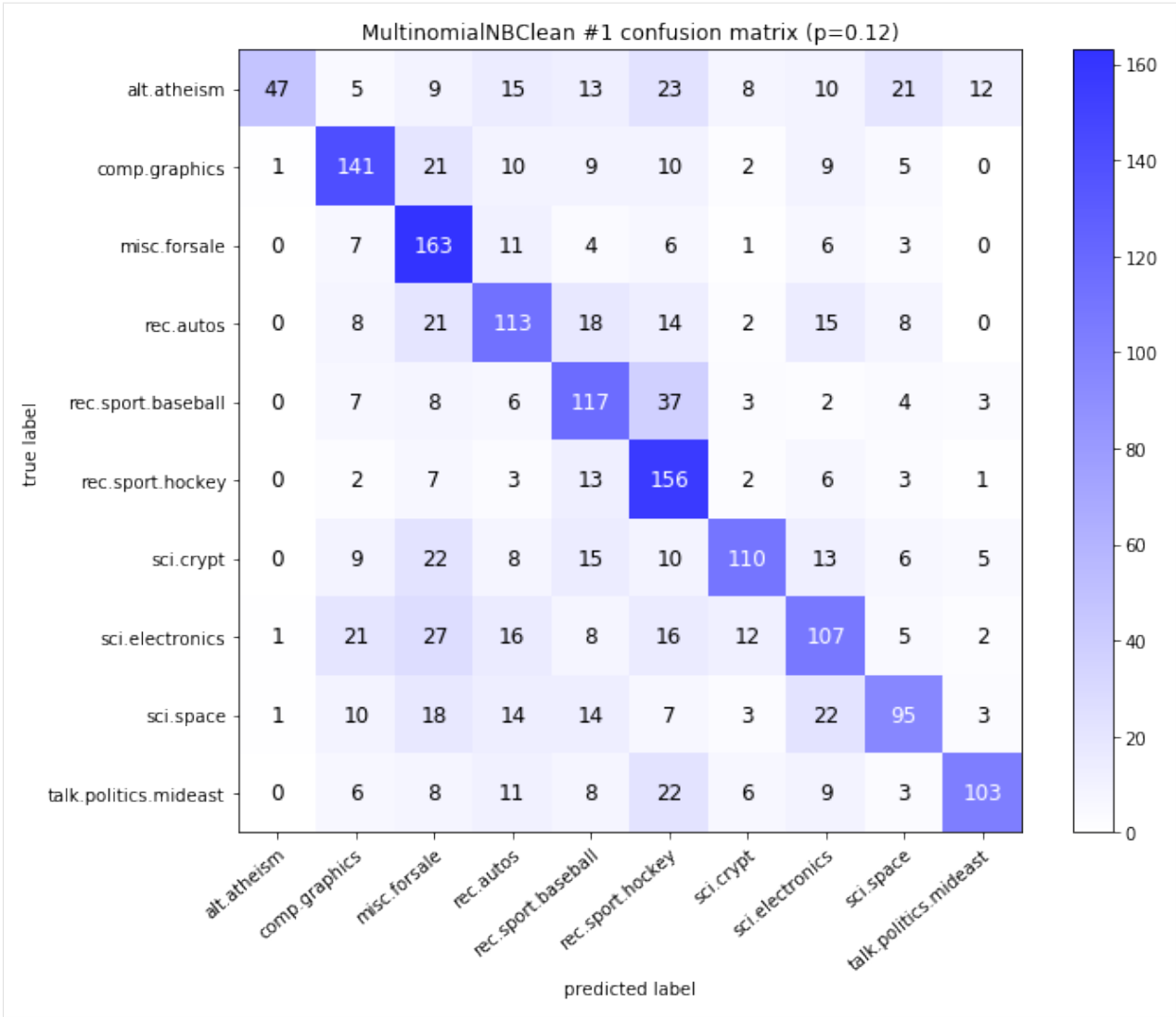


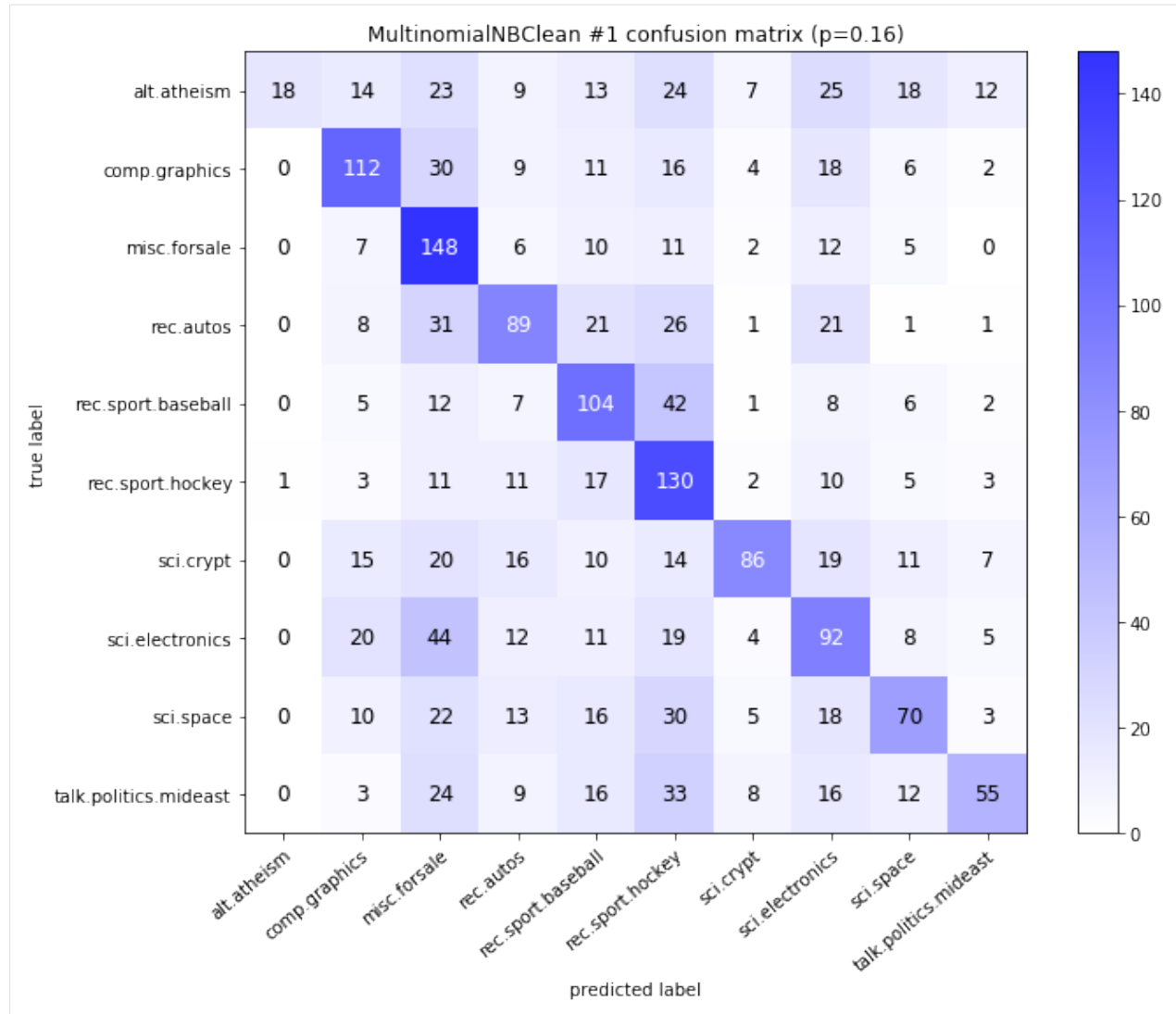


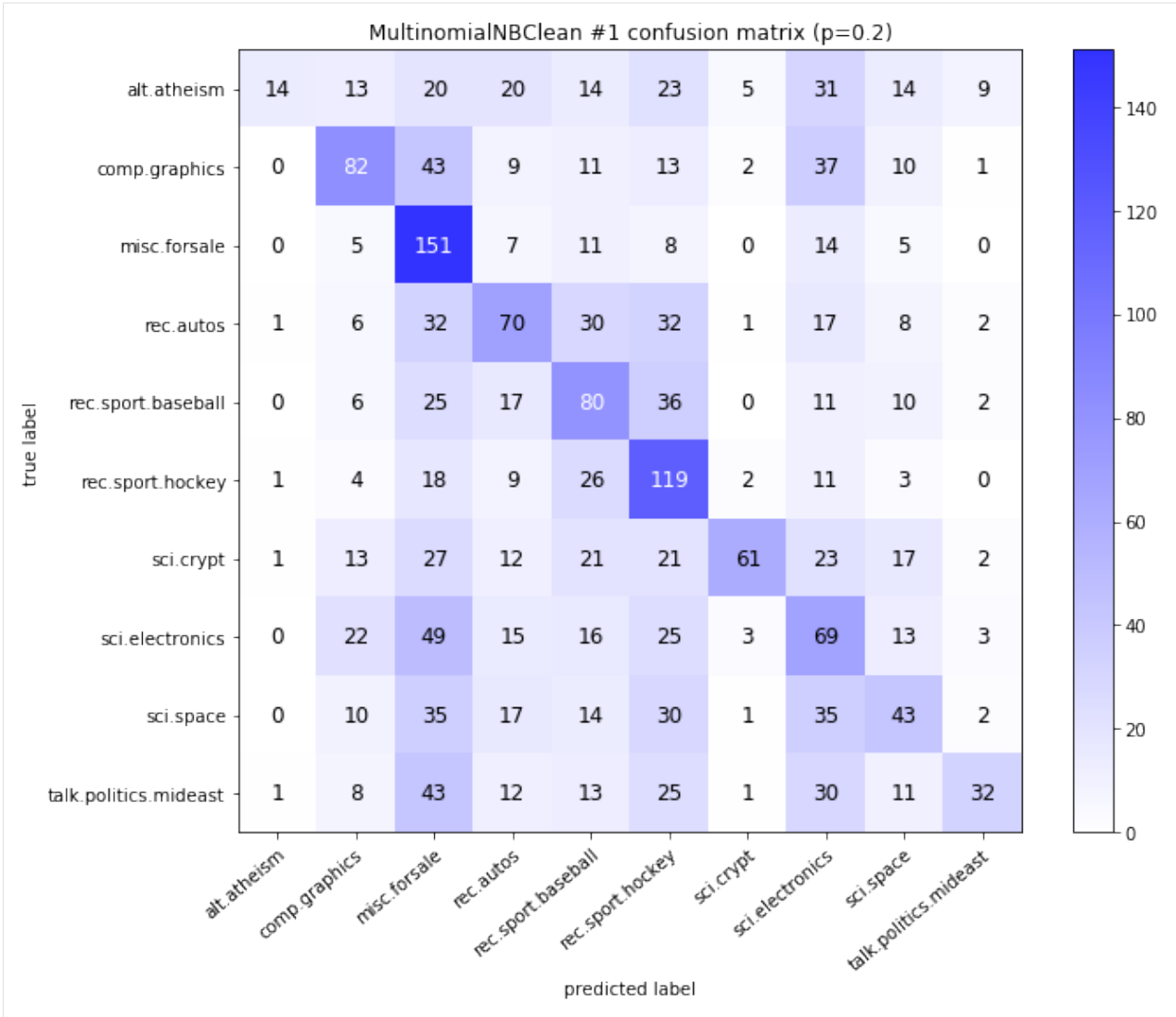


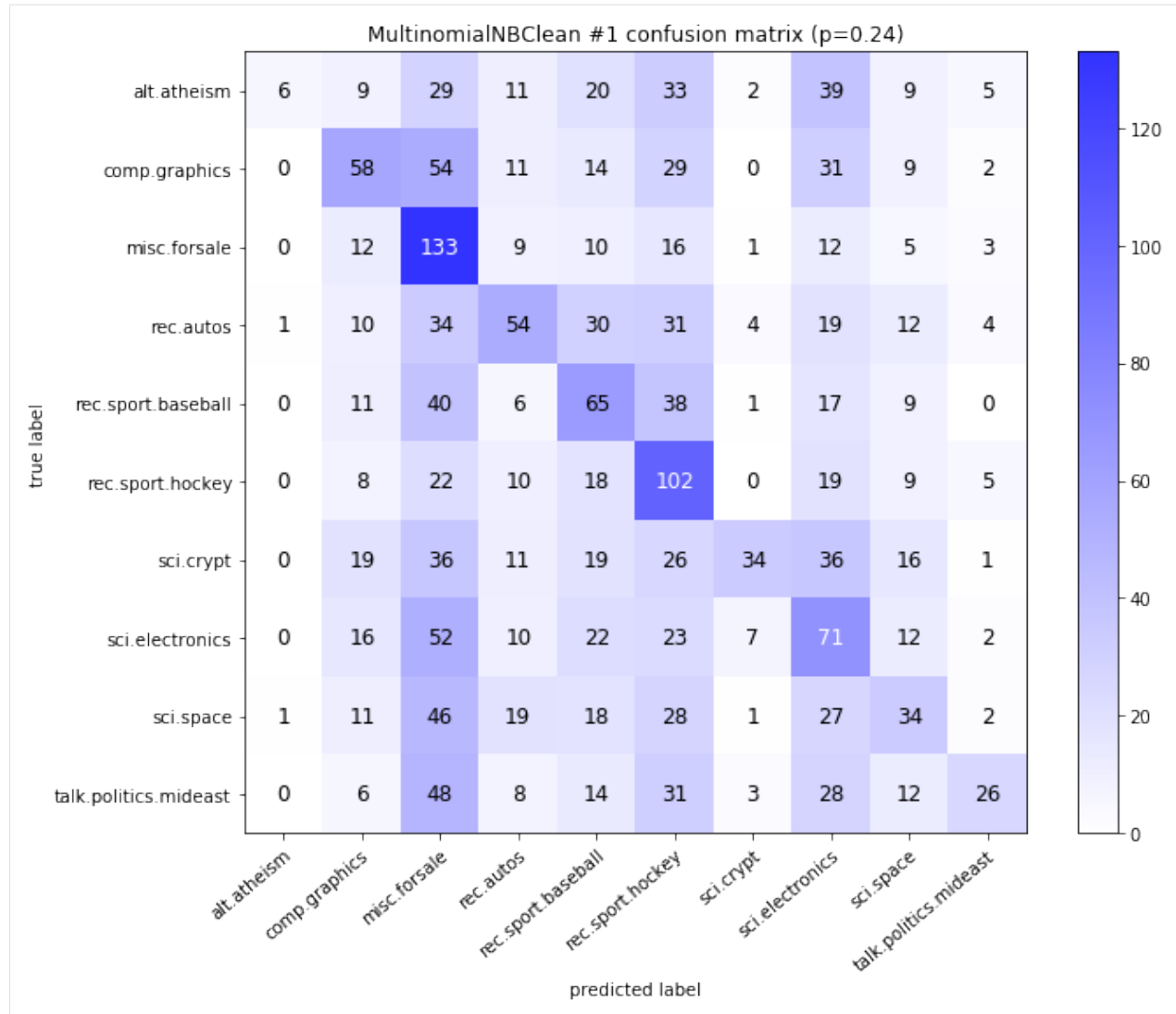


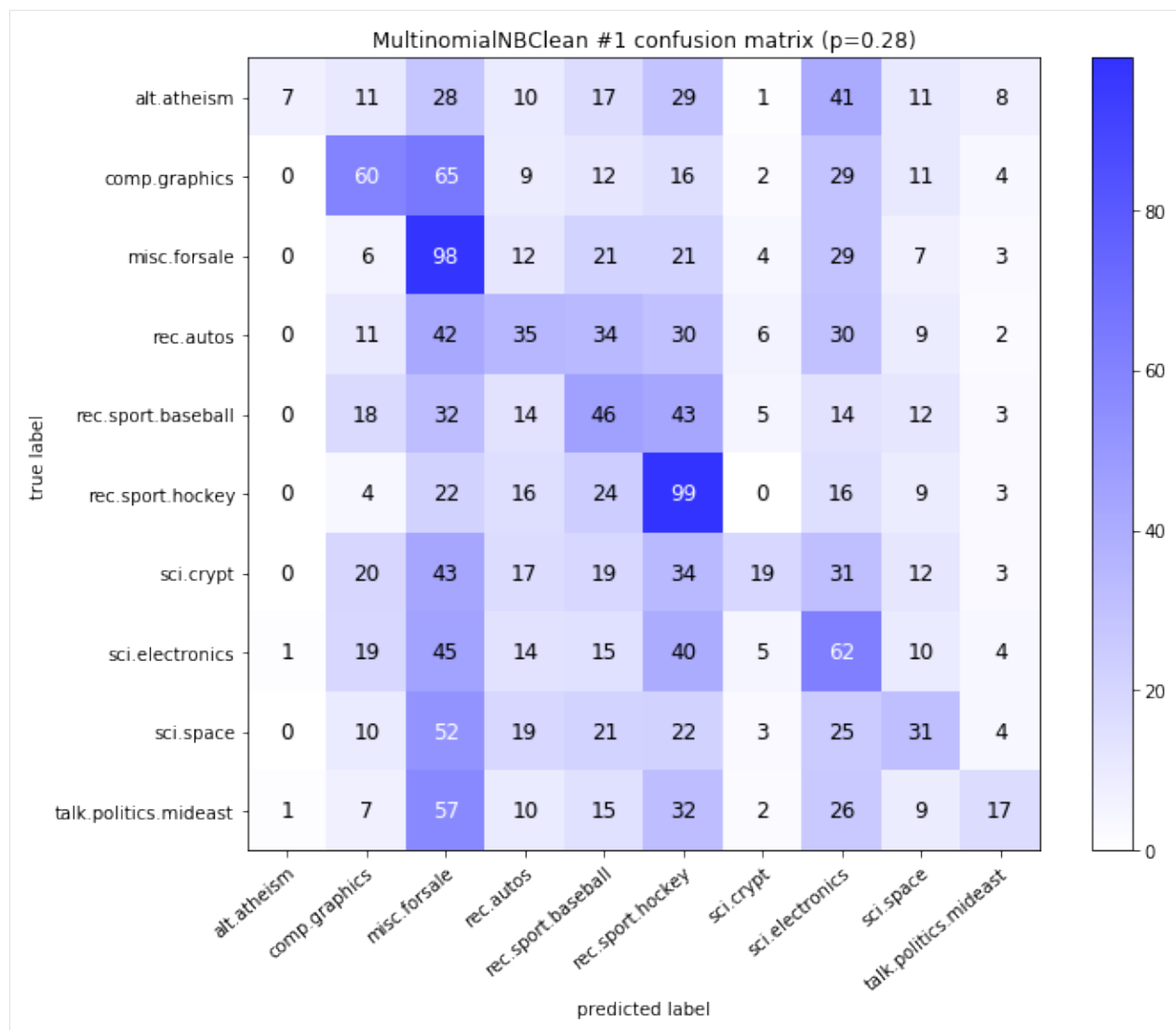












The notebook for this case study can be found [here](#).

4.4 Text classification: OCR error

```
[1]: import warnings
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
import numpy as np
from numba import NumbaDeprecationWarning, NumbaWarning
from numpy.random import RandomState
from sklearn.exceptions import ConvergenceWarning
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.svm import LinearSVC
```

(continues on next page)

(continued from previous page)

```

from dpemu import runner
from dpemu.dataset_utils import load_newsgroups
from dpemu.filters.text import OCRError
from dpemu.ml_utils import reduce_dimensions_sparse
from dpemu.nodes.array import Array
from dpemu.pg_utils import load_ocr_error_params, normalize_ocr_error_params
from dpemu.plotting_utils import visualize_best_model_params, visualize_scores, \
    visualize_classes, \
    print_results_by_model, visualize_confusion_matrices
from dpemu.utils import get_project_root

warnings.simplefilter("ignore", category=ConvergenceWarning)
warnings.simplefilter("ignore", category=NumbaDeprecationWarning)
warnings.simplefilter("ignore", category=NumbaWarning)

```

```

[2]: def get_data():
    data, labels, label_names, dataset_name = load_newsgroups("all", 10)
    train_data, test_data, train_labels, test_labels = train_test_split(data, labels, \
    ↪test_size=.2,
                                                    random_
    ↪state=RandomState(42))
    return train_data, test_data, train_labels, test_labels, label_names, dataset_name

```

```

[3]: def get_err_root_node():
    err_root_node = Array()
    err_root_node.addfilter(OCRError("normalized_params", "p"))
    return err_root_node

```

```

[4]: def get_err_params_list():
    p_steps = np.linspace(0, .98, num=8)
    params = load_ocr_error_params(f"{get_project_root()}/data/example_ocr_error_
    ↪config.json")
    normalized_params = normalize_ocr_error_params(params)
    err_params_list = [{
        "p": p,
        "normalized_params": normalized_params
    } for p in p_steps]

    return err_params_list

```

```

[5]: class Preprocessor:
    def __init__(self):
        self.random_state = RandomState(0)

    def run(self, train_data, test_data, _):
        vectorizer = TfidfVectorizer(max_df=0.5, min_df=2, stop_words="english")
        vectorized_train_data = vectorizer.fit_transform(train_data)
        vectorized_test_data = vectorizer.transform(test_data)

        reduced_test_data = reduce_dimensions_sparse(vectorized_test_data, self.
    ↪random_state)

        return vectorized_train_data, vectorized_test_data, {"reduced_test_data": \
    ↪reduced_test_data}

```

(continues on next page)

(continued from previous page)

```
[6]: class AbstractModel(ABC):

    def __init__(self):
        self.random_state = RandomState(42)

    @abstractmethod
    def get_fitted_model(self, train_data, train_labels, params):
        pass

    def run(self, train_data, test_data, params):
        train_labels = params["train_labels"]
        test_labels = params["test_labels"]

        fitted_model = self.get_fitted_model(train_data, train_labels, params)

        predicted_test_labels = fitted_model.predict(test_data)
        cm = confusion_matrix(test_labels, predicted_test_labels)
        return {
            "confusion_matrix": cm,
            "predicted_test_labels": predicted_test_labels,
            "test_mean_accuracy": round(np.mean(predicted_test_labels == test_labels),
→ 3),
            "train_mean_accuracy": fitted_model.score(train_data, train_labels),
        }

class MultinomialNBModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, train_data, train_labels, params):
        return MultinomialNB(params["alpha"]).fit(train_data, train_labels)

class LinearSVCModel(AbstractModel):

    def __init__(self):
        super().__init__()

    def get_fitted_model(self, train_data, train_labels, params):
        return LinearSVC(C=params["C"], random_state=self.random_state).fit(train_
→data, train_labels)

[7]: def get_model_params_dict_list(train_labels, test_labels):
    alpha_steps = [10 ** i for i in range(-4, 1)]
    C_steps = [10 ** k for k in range(-3, 2)]
    model_params_base = {"train_labels": train_labels, "test_labels": test_labels}
    return [
        {
            "model": MultinomialNBModel,
            "params_list": [{"alpha": alpha, **model_params_base} for alpha in alpha_
→steps],
            "use_clean_train_data": False
```

(continues on next page)

(continued from previous page)

```

    },
    {
        "model": MultinomialNBModel,
        "params_list": [{"alpha": alpha, **model_params_base} for alpha in alpha_
→ steps],
        "use_clean_train_data": True
    },
    {
        "model": LinearSVCModel,
        "params_list": [{"C": C, **model_params_base} for C in C_steps],
        "use_clean_train_data": False
    },
    {
        "model": LinearSVCModel,
        "params_list": [{"C": C, **model_params_base} for C in C_steps],
        "use_clean_train_data": True
    },
    },
    ]

```

```

[8]: def visualize(df, dataset_name, label_names, test_data):
    visualize_scores(
        df,
        score_names=["test_mean_accuracy", "train_mean_accuracy"],
        is_higher_score_better=[True, True],
        err_param_name="p",
        title=f"{dataset_name} classification scores with added error"
    )
    visualize_best_model_params(
        df,
        "MultinomialNB",
        model_params=["alpha"],
        score_names=["test_mean_accuracy"],
        is_higher_score_better=[True],
        err_param_name="p",
        title=f"Best parameters for {dataset_name} classification",
        y_log=True
    )
    visualize_best_model_params(
        df,
        "LinearSVC",
        model_params=["C"],
        score_names=["test_mean_accuracy"],
        is_higher_score_better=[True],
        err_param_name="p",
        title=f"Best parameters for {dataset_name} classification",
        y_log=True
    )
    visualize_classes(
        df,
        label_names,
        err_param_name="p",
        reduced_data_column="reduced_test_data",
        labels_column="test_labels",
        cmap="tab20",
        title=f"{dataset_name} test set (n={len(test_data)}) true classes with added_
→ error"

```

(continues on next page)

(continued from previous page)

```

)
visualize_confusion_matrices(
    df,
    label_names,
    score_name="test_mean_accuracy",
    is_higher_score_better=True,
    err_param_name="p",
    labels_col="test_labels",
    predictions_col="predicted_test_labels",
)
plt.show()

```

```

[9]: def main():
    train_data, test_data, train_labels, test_labels, label_names, dataset_name = get_
    ↪data()

    df = runner.run(
        train_data=train_data,
        test_data=test_data,
        preproc=Preprocessor,
        preproc_params=None,
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(train_labels, test_labels),
    )

    print_results_by_model(df, dropped_columns=[
        "train_labels", "test_labels", "reduced_test_data", "confusion_matrix",
    ↪"predicted_test_labels",
        "normalized_params"
    ])
    visualize(df, dataset_name, label_names, test_data)

```

Models LinearSVCClean and MultinomialNBClean have been trained with clean data and LinearSVC and MultinomialNB with erroneous data.

```

[10]: main()
100%|| 8/8 [09:23<00:00, 84.65s/it]
LinearSVC #1

```

	test_mean_accuracy	train_mean_accuracy	p	C	time_err	time_pre	time_
↪mod							
0	0.720	0.792343	0.00	0.001	14.081	20.823	0.
↪199							
1	0.809	0.884620	0.00	0.010	14.081	20.823	0.
↪251							
2	0.842	0.946658	0.00	0.100	14.081	20.823	0.
↪289							
3	0.846	0.973005	0.00	1.000	14.081	20.823	0.
↪461							
4	0.835	0.974043	0.00	10.000	14.081	20.823	2.
↪246							
5	0.645	0.760286	0.14	0.001	77.115	21.137	0.
↪266							
6	0.751	0.879948	0.14	0.010	77.115	21.137	0.
↪299							

(continues on next page)

(continued from previous page)

7	0.799	0.951979	0.14	0.100	77.115	21.137	0.
↪ 370							
8	0.798	0.973783	0.14	1.000	77.115	21.137	0.
↪ 625							
9	0.789	0.974173	0.14	10.000	77.115	21.137	3.
↪ 082							
10	0.591	0.719533	0.28	0.001	140.397	20.408	0.
↪ 264							
11	0.693	0.869825	0.28	0.010	140.397	20.408	0.
↪ 295							
12	0.737	0.956132	0.28	0.100	140.397	20.408	0.
↪ 376							
13	0.751	0.973783	0.28	1.000	140.397	20.408	0.
↪ 662							
14	0.742	0.974173	0.28	10.000	140.397	20.408	3.
↪ 285							
15	0.510	0.675146	0.42	0.001	220.260	21.367	0.
↪ 296							
16	0.630	0.848799	0.42	0.010	220.260	21.367	0.
↪ 387							
17	0.692	0.956652	0.42	0.100	220.260	21.367	0.
↪ 425							
18	0.699	0.973394	0.42	1.000	220.260	21.367	0.
↪ 748							
19	0.691	0.974043	0.42	10.000	220.260	21.367	3.
↪ 867							
20	0.443	0.627255	0.56	0.001	260.029	20.559	0.
↪ 285							
21	0.557	0.815185	0.56	0.010	260.029	20.559	0.
↪ 342							
22	0.642	0.957949	0.56	0.100	260.029	20.559	0.
↪ 417							
23	0.657	0.974043	0.56	1.000	260.029	20.559	0.
↪ 760							
24	0.650	0.974692	0.56	10.000	260.029	20.559	4.
↪ 079							
25	0.396	0.586892	0.70	0.001	345.086	20.691	0.
↪ 283							
26	0.504	0.788709	0.70	0.010	345.086	20.691	0.
↪ 353							
27	0.624	0.956003	0.70	0.100	345.086	20.691	0.
↪ 417							
28	0.639	0.974043	0.70	1.000	345.086	20.691	0.
↪ 764							
29	0.627	0.974562	0.70	10.000	345.086	20.691	4.
↪ 118							
30	0.371	0.575730	0.84	0.001	379.122	21.887	0.
↪ 282							
31	0.466	0.776639	0.84	0.010	379.122	21.887	0.
↪ 349							
32	0.581	0.958339	0.84	0.100	379.122	21.887	0.
↪ 415							
33	0.604	0.974043	0.84	1.000	379.122	21.887	0.
↪ 762							
34	0.603	0.974951	0.84	10.000	379.122	21.887	4.
↪ 236							
35	0.357	0.554705	0.98	0.001	532.418	21.617	0.
↪ 286							

(continues on next page)

(continued from previous page)

36	0.453	0.757430	0.98	0.010	532.418	21.617	0.
↪360							
37	0.580	0.955354	0.98	0.100	532.418	21.617	0.
↪419							
38	0.604	0.973653	0.98	1.000	532.418	21.617	0.
↪774							
39	0.592	0.974432	0.98	10.000	532.418	21.617	4.
↪261							
LinearSVCClean #1							
	test_mean_accuracy	train_mean_accuracy	p	C	time_err	time_pre	time_
↪mod							
0	0.720	0.792343	0.00	0.001	14.081	20.823	0.
↪197							
1	0.809	0.884620	0.00	0.010	14.081	20.823	0.
↪250							
2	0.842	0.946658	0.00	0.100	14.081	20.823	0.
↪288							
3	0.846	0.973005	0.00	1.000	14.081	20.823	0.
↪461							
4	0.835	0.974043	0.00	10.000	14.081	20.823	2.
↪240							
5	0.678	0.792343	0.14	0.001	77.115	21.137	0.
↪206							
6	0.771	0.884620	0.14	0.010	77.115	21.137	0.
↪259							
7	0.812	0.946658	0.14	0.100	77.115	21.137	0.
↪300							
8	0.801	0.973005	0.14	1.000	77.115	21.137	0.
↪482							
9	0.778	0.974043	0.14	10.000	77.115	21.137	2.
↪384							
10	0.627	0.792343	0.28	0.001	140.397	20.408	0.
↪184							
11	0.722	0.884620	0.28	0.010	140.397	20.408	0.
↪234							
12	0.755	0.946658	0.28	0.100	140.397	20.408	0.
↪268							
13	0.733	0.973005	0.28	1.000	140.397	20.408	0.
↪428							
14	0.718	0.974043	0.28	10.000	140.397	20.408	2.
↪086							
15	0.558	0.792343	0.42	0.001	220.260	21.367	0.
↪211							
16	0.650	0.884620	0.42	0.010	220.260	21.367	0.
↪246							
17	0.679	0.946658	0.42	0.100	220.260	21.367	0.
↪281							
18	0.661	0.973005	0.42	1.000	220.260	21.367	0.
↪450							
19	0.636	0.974043	0.42	10.000	220.260	21.367	2.
↪200							
20	0.457	0.792343	0.56	0.001	260.029	20.559	0.
↪183							
21	0.564	0.884620	0.56	0.010	260.029	20.559	0.
↪233							
22	0.592	0.946658	0.56	0.100	260.029	20.559	0.
↪270							

(continues on next page)

(continued from previous page)

23	0.551	0.973005	0.56	1.000	260.029	20.559	0.
↪433							
24	0.523	0.974043	0.56	10.000	260.029	20.559	2.
↪115							
25	0.366	0.792343	0.70	0.001	345.086	20.691	0.
↪177							
26	0.474	0.884620	0.70	0.010	345.086	20.691	0.
↪224							
27	0.508	0.946658	0.70	0.100	345.086	20.691	0.
↪258							
28	0.461	0.973005	0.70	1.000	345.086	20.691	0.
↪418							
29	0.418	0.974043	0.70	10.000	345.086	20.691	2.
↪038							
30	0.299	0.792343	0.84	0.001	379.122	21.887	0.
↪178							
31	0.382	0.884620	0.84	0.010	379.122	21.887	0.
↪226							
32	0.406	0.946658	0.84	0.100	379.122	21.887	0.
↪260							
33	0.371	0.973005	0.84	1.000	379.122	21.887	0.
↪417							
34	0.343	0.974043	0.84	10.000	379.122	21.887	2.
↪042							
35	0.239	0.792343	0.98	0.001	532.418	21.617	0.
↪177							
36	0.311	0.884620	0.98	0.010	532.418	21.617	0.
↪224							
37	0.317	0.946658	0.98	0.100	532.418	21.617	0.
↪259							
38	0.281	0.973005	0.98	1.000	532.418	21.617	0.
↪415							
39	0.257	0.974043	0.98	10.000	532.418	21.617	2.
↪037							
MultinomialNB #1							
	test_mean_accuracy	train_mean_accuracy	p	alpha	time_err	time_pre	time_
↪mod							
0	0.834	0.961583	0.00	0.0001	14.081	20.823	0.
↪038							
1	0.843	0.960675	0.00	0.0010	14.081	20.823	0.
↪032							
2	0.851	0.958209	0.00	0.0100	14.081	20.823	0.
↪032							
3	0.852	0.951979	0.00	0.1000	14.081	20.823	0.
↪032							
4	0.832	0.925892	0.00	1.0000	14.081	20.823	0.
↪032							
5	0.763	0.968981	0.14	0.0001	77.115	21.137	0.
↪054							
6	0.786	0.968981	0.14	0.0010	77.115	21.137	0.
↪044							
7	0.805	0.968332	0.14	0.0100	77.115	21.137	0.
↪044							
8	0.810	0.966126	0.14	0.1000	77.115	21.137	0.
↪044							
9	0.794	0.940818	0.14	1.0000	77.115	21.137	0.
↪057							

(continues on next page)

(continued from previous page)

10	0.728	0.970539	0.28	0.0001	140.397	20.408	0.
↪048							
11	0.747	0.970409	0.28	0.0010	140.397	20.408	0.
↪046							
12	0.770	0.970019	0.28	0.0100	140.397	20.408	0.
↪046							
13	0.782	0.968073	0.28	0.1000	140.397	20.408	0.
↪046							
14	0.736	0.946009	0.28	1.0000	140.397	20.408	0.
↪046							
15	0.691	0.971317	0.42	0.0001	220.260	21.367	0.
↪073							
16	0.706	0.971317	0.42	0.0010	220.260	21.367	0.
↪055							
17	0.736	0.971317	0.42	0.0100	220.260	21.367	0.
↪055							
18	0.742	0.970409	0.42	0.1000	220.260	21.367	0.
↪055							
19	0.665	0.946398	0.42	1.0000	220.260	21.367	0.
↪055							
20	0.665	0.971317	0.56	0.0001	260.029	20.559	0.
↪053							
21	0.691	0.971317	0.56	0.0010	260.029	20.559	0.
↪050							
22	0.705	0.970928	0.56	0.0100	260.029	20.559	0.
↪050							
23	0.708	0.970279	0.56	0.1000	260.029	20.559	0.
↪050							
24	0.601	0.941077	0.56	1.0000	260.029	20.559	0.
↪050							
25	0.625	0.971836	0.70	0.0001	345.086	20.691	0.
↪056							
26	0.649	0.971836	0.70	0.0010	345.086	20.691	0.
↪050							
27	0.678	0.971836	0.70	0.0100	345.086	20.691	0.
↪050							
28	0.687	0.971058	0.70	0.1000	345.086	20.691	0.
↪050							
29	0.551	0.928358	0.70	1.0000	345.086	20.691	0.
↪050							
30	0.615	0.972356	0.84	0.0001	379.122	21.887	0.
↪053							
31	0.637	0.972356	0.84	0.0010	379.122	21.887	0.
↪050							
32	0.655	0.972226	0.84	0.0100	379.122	21.887	0.
↪050							
33	0.649	0.972226	0.84	0.1000	379.122	21.887	0.
↪050							
34	0.508	0.923945	0.84	1.0000	379.122	21.887	0.
↪050							
35	0.595	0.971707	0.98	0.0001	532.418	21.617	0.
↪055							
36	0.618	0.971707	0.98	0.0010	532.418	21.617	0.
↪051							
37	0.651	0.971707	0.98	0.0100	532.418	21.617	0.
↪051							
38	0.650	0.970668	0.98	0.1000	532.418	21.617	0.
↪051							

(continues on next page)

(continued from previous page)

39	0.480	0.912914	0.98	1.0000	532.418	21.617	0.
↪051							
MultinomialNBClean #1							
	test_mean_accuracy	train_mean_accuracy	p	alpha	time_err	time_pre	time_
↪mod							
0	0.834	0.961583	0.00	0.0001	14.081	20.823	0.
↪035							
1	0.843	0.960675	0.00	0.0010	14.081	20.823	0.
↪032							
2	0.851	0.958209	0.00	0.0100	14.081	20.823	0.
↪032							
3	0.852	0.951979	0.00	0.1000	14.081	20.823	0.
↪032							
4	0.832	0.925892	0.00	1.0000	14.081	20.823	0.
↪032							
5	0.778	0.961583	0.14	0.0001	77.115	21.137	0.
↪044							
6	0.801	0.960675	0.14	0.0010	77.115	21.137	0.
↪031							
7	0.818	0.958209	0.14	0.0100	77.115	21.137	0.
↪031							
8	0.825	0.951979	0.14	0.1000	77.115	21.137	0.
↪031							
9	0.808	0.925892	0.14	1.0000	77.115	21.137	0.
↪031							
10	0.708	0.961583	0.28	0.0001	140.397	20.408	0.
↪033							
11	0.729	0.960675	0.28	0.0010	140.397	20.408	0.
↪029							
12	0.751	0.958209	0.28	0.0100	140.397	20.408	0.
↪029							
13	0.765	0.951979	0.28	0.1000	140.397	20.408	0.
↪029							
14	0.754	0.925892	0.28	1.0000	140.397	20.408	0.
↪029							
15	0.575	0.961583	0.42	0.0001	220.260	21.367	0.
↪034							
16	0.605	0.960675	0.42	0.0010	220.260	21.367	0.
↪029							
17	0.643	0.958209	0.42	0.0100	220.260	21.367	0.
↪033							
18	0.683	0.951979	0.42	0.1000	220.260	21.367	0.
↪029							
19	0.680	0.925892	0.42	1.0000	220.260	21.367	0.
↪029							
20	0.451	0.961583	0.56	0.0001	260.029	20.559	0.
↪031							
21	0.477	0.960675	0.56	0.0010	260.029	20.559	0.
↪027							
22	0.517	0.958209	0.56	0.0100	260.029	20.559	0.
↪027							
23	0.563	0.951979	0.56	0.1000	260.029	20.559	0.
↪027							
24	0.582	0.925892	0.56	1.0000	260.029	20.559	0.
↪027							
25	0.337	0.961583	0.70	0.0001	345.086	20.691	0.
↪030							

(continues on next page)

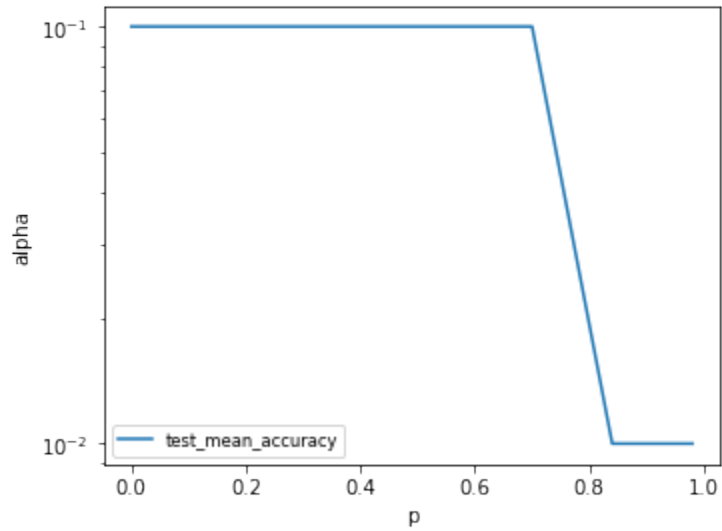
(continued from previous page)

26	0.363	0.960675	0.70	0.0010	345.086	20.691	0.
↪026							
27	0.387	0.958209	0.70	0.0100	345.086	20.691	0.
↪026							
28	0.445	0.951979	0.70	0.1000	345.086	20.691	0.
↪026							
29	0.451	0.925892	0.70	1.0000	345.086	20.691	0.
↪026							
30	0.261	0.961583	0.84	0.0001	379.122	21.887	0.
↪031							
31	0.272	0.960675	0.84	0.0010	379.122	21.887	0.
↪027							
32	0.290	0.958209	0.84	0.0100	379.122	21.887	0.
↪027							
33	0.320	0.951979	0.84	0.1000	379.122	21.887	0.
↪027							
34	0.353	0.925892	0.84	1.0000	379.122	21.887	0.
↪027							
35	0.216	0.961583	0.98	0.0001	532.418	21.617	0.
↪030							
36	0.218	0.960675	0.98	0.0010	532.418	21.617	0.
↪027							
37	0.233	0.958209	0.98	0.0100	532.418	21.617	0.
↪026							
38	0.261	0.951979	0.98	0.1000	532.418	21.617	0.
↪026							
39	0.276	0.925892	0.98	1.0000	532.418	21.617	0.
↪026							

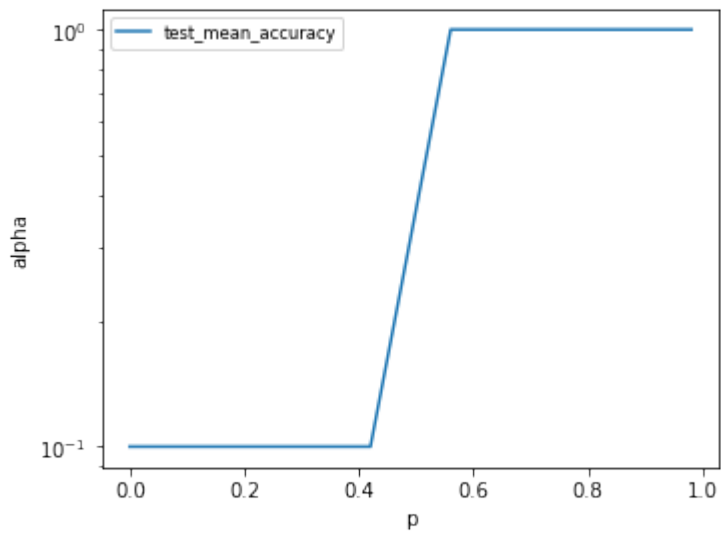
/wrk/users/thalvari/dpEmu/dpemu/plotting_utils.py:299: RuntimeWarning: More than 20_
↪figures have been opened. Figures created through the pyplot interface (`matplotlib.
↪pyplot.figure`) are retained until explicitly closed and may consume too much_
↪memory. (To control this warning, see the rcParam `figure.max_open_warning`).
fig, ax = plt.subplots(figsize=(10, 8))

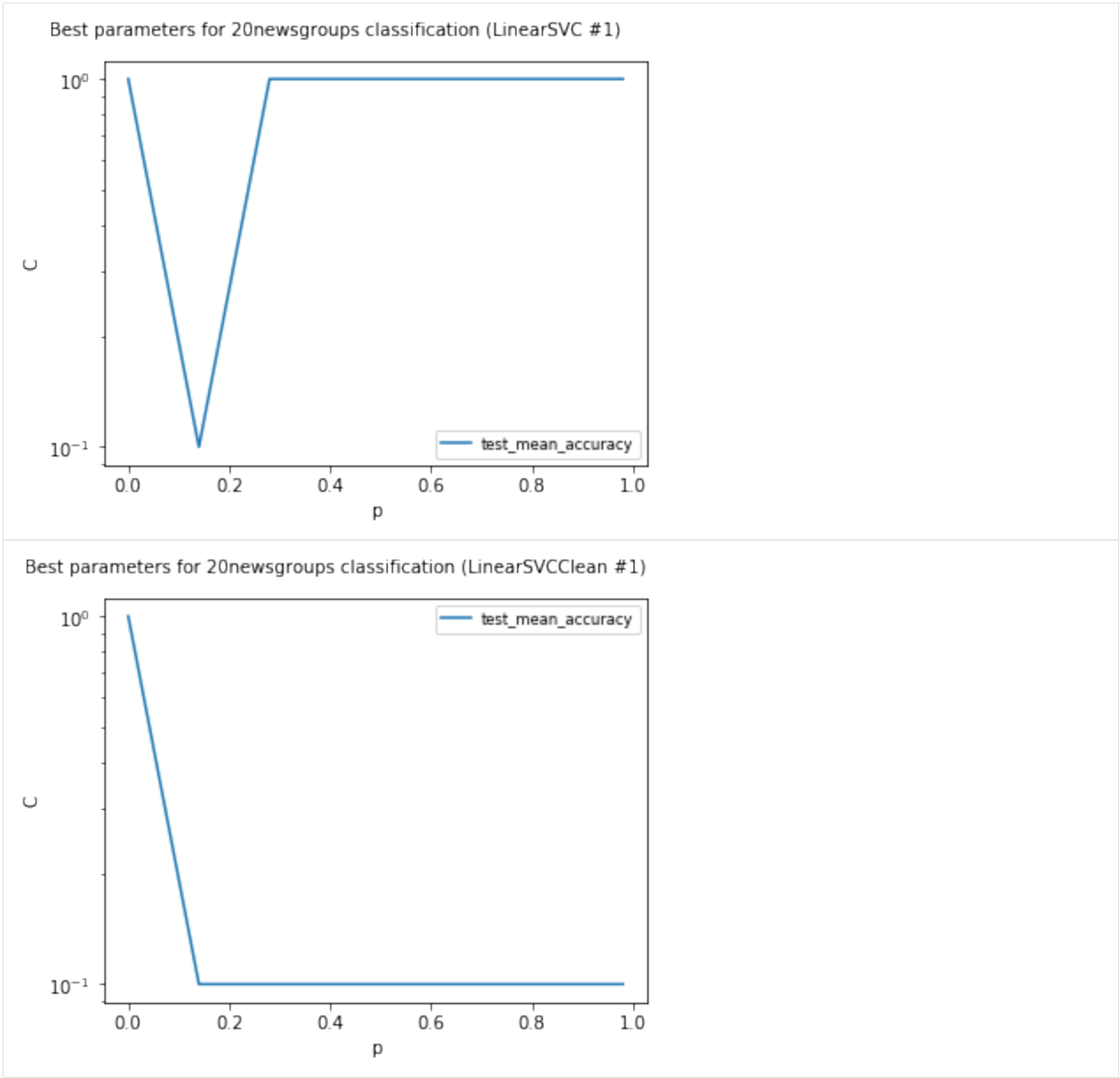


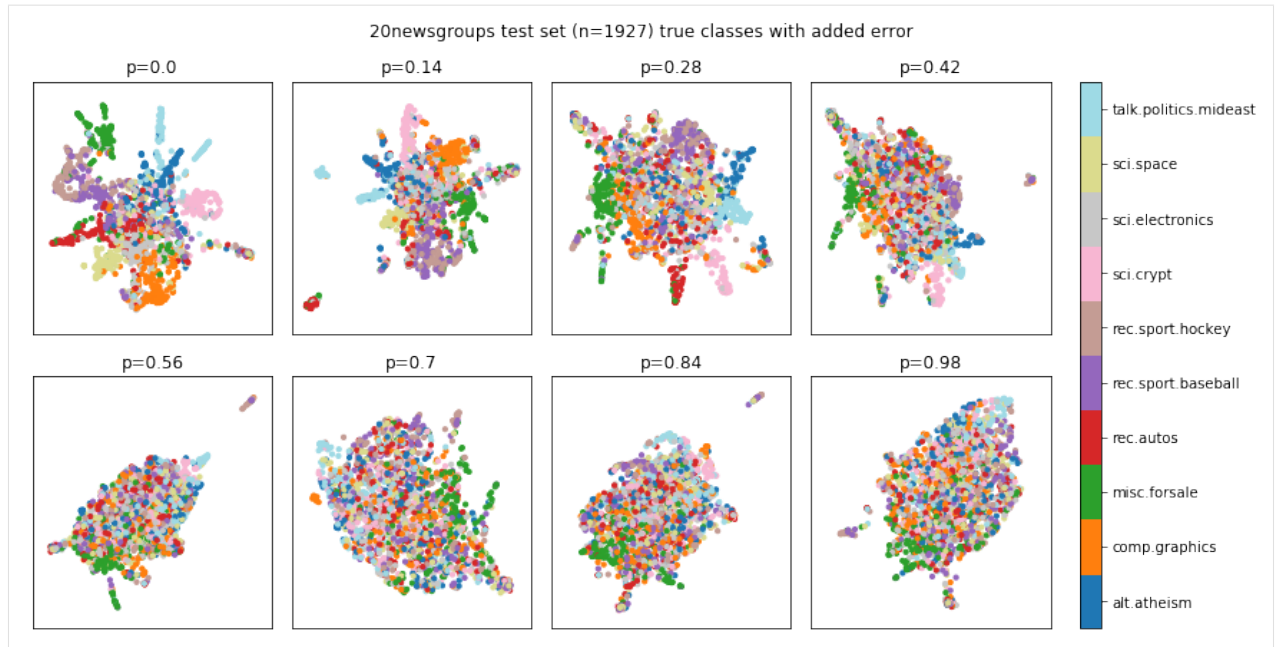
Best parameters for 20newsgroups classification (MultinomialNB #1)

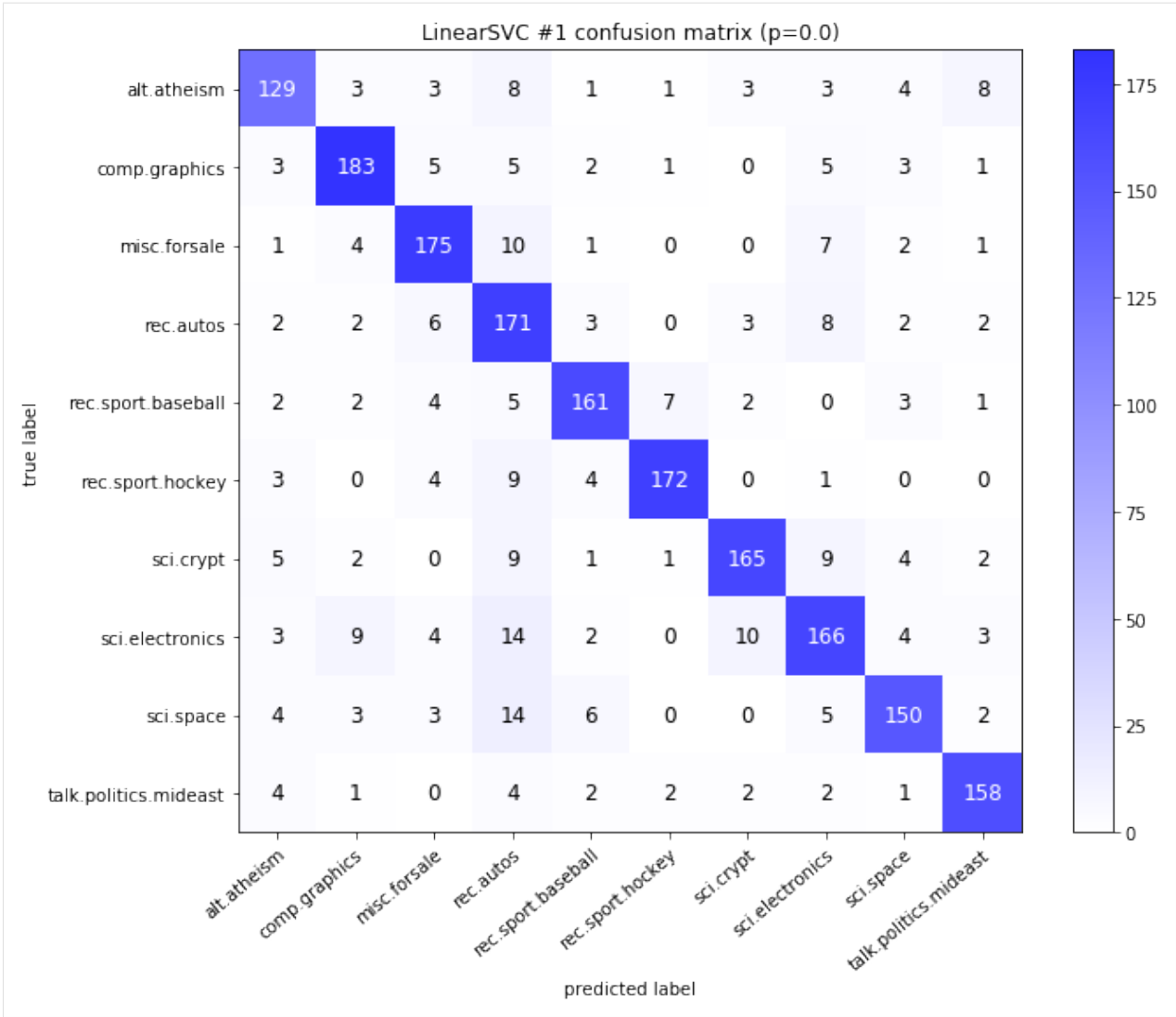


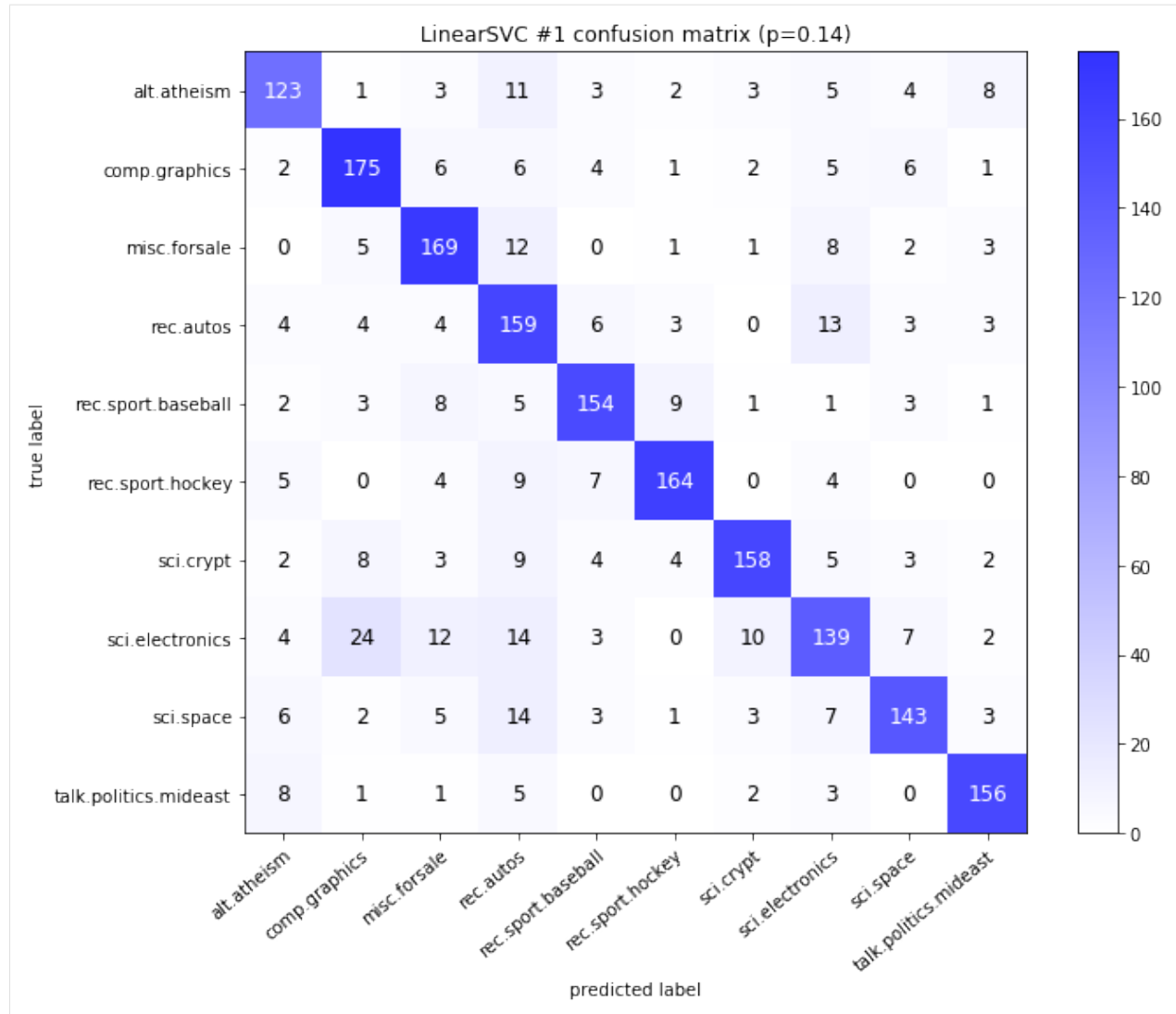
Best parameters for 20newsgroups classification (MultinomialNBClean #1)

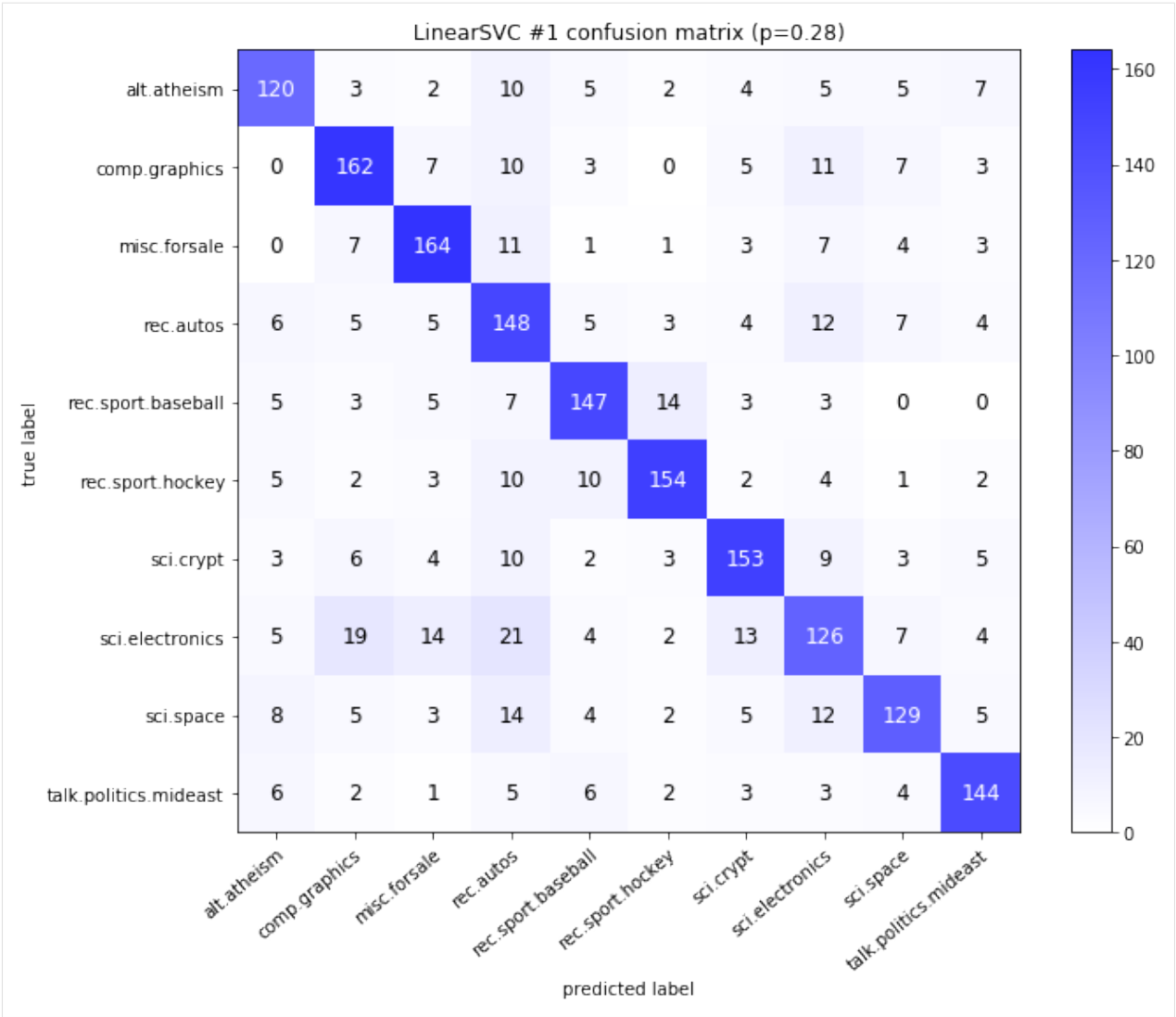


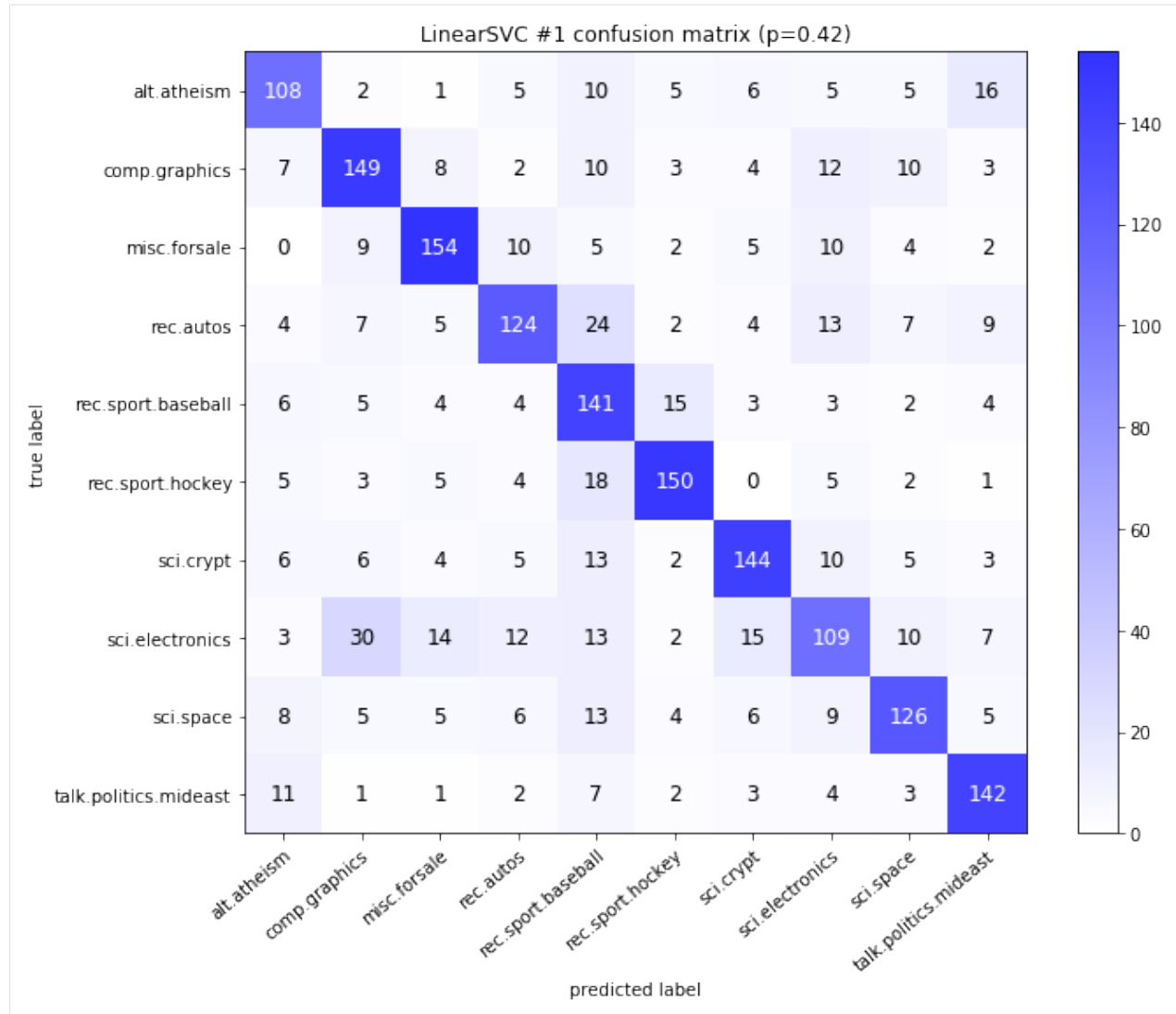


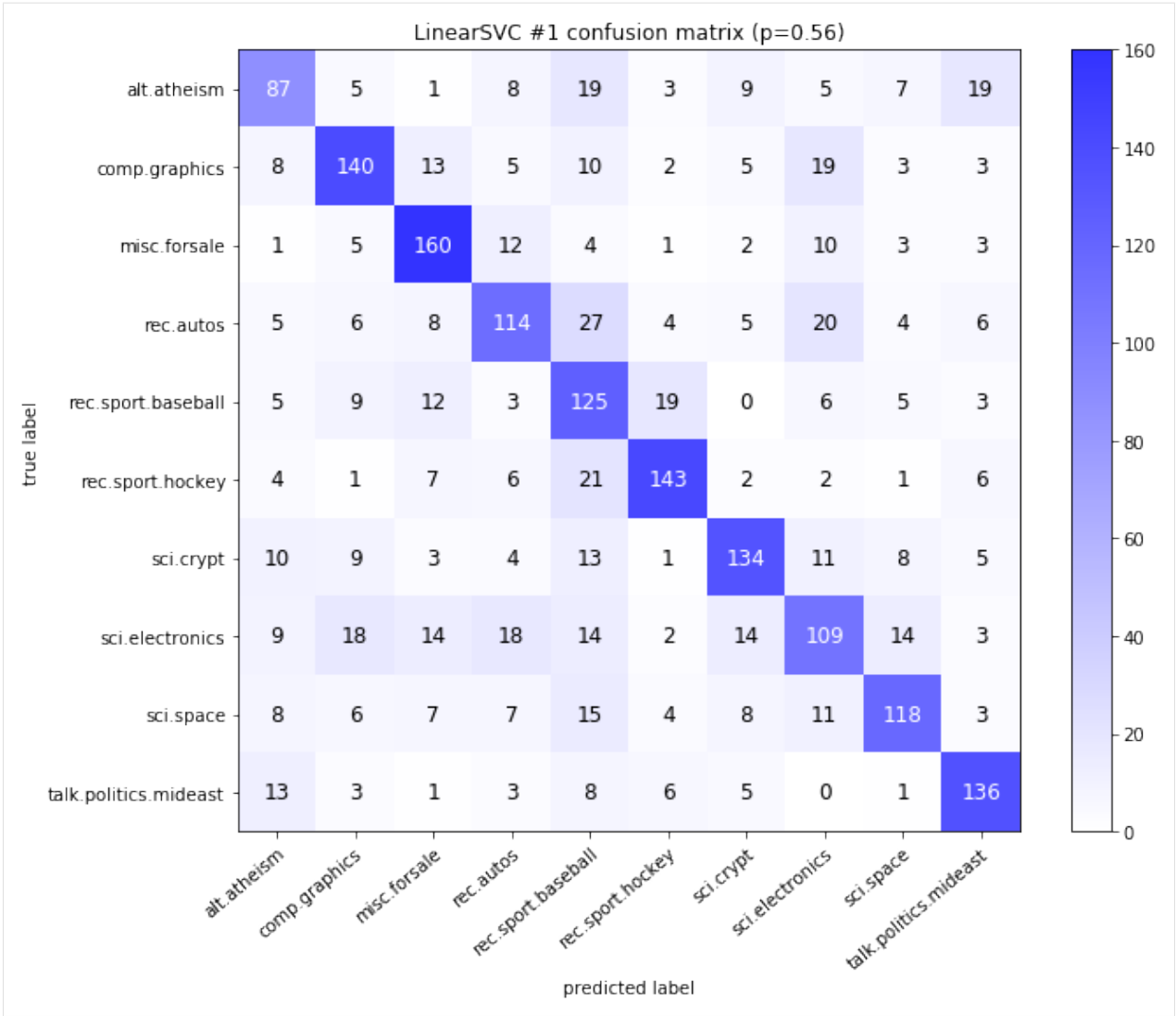


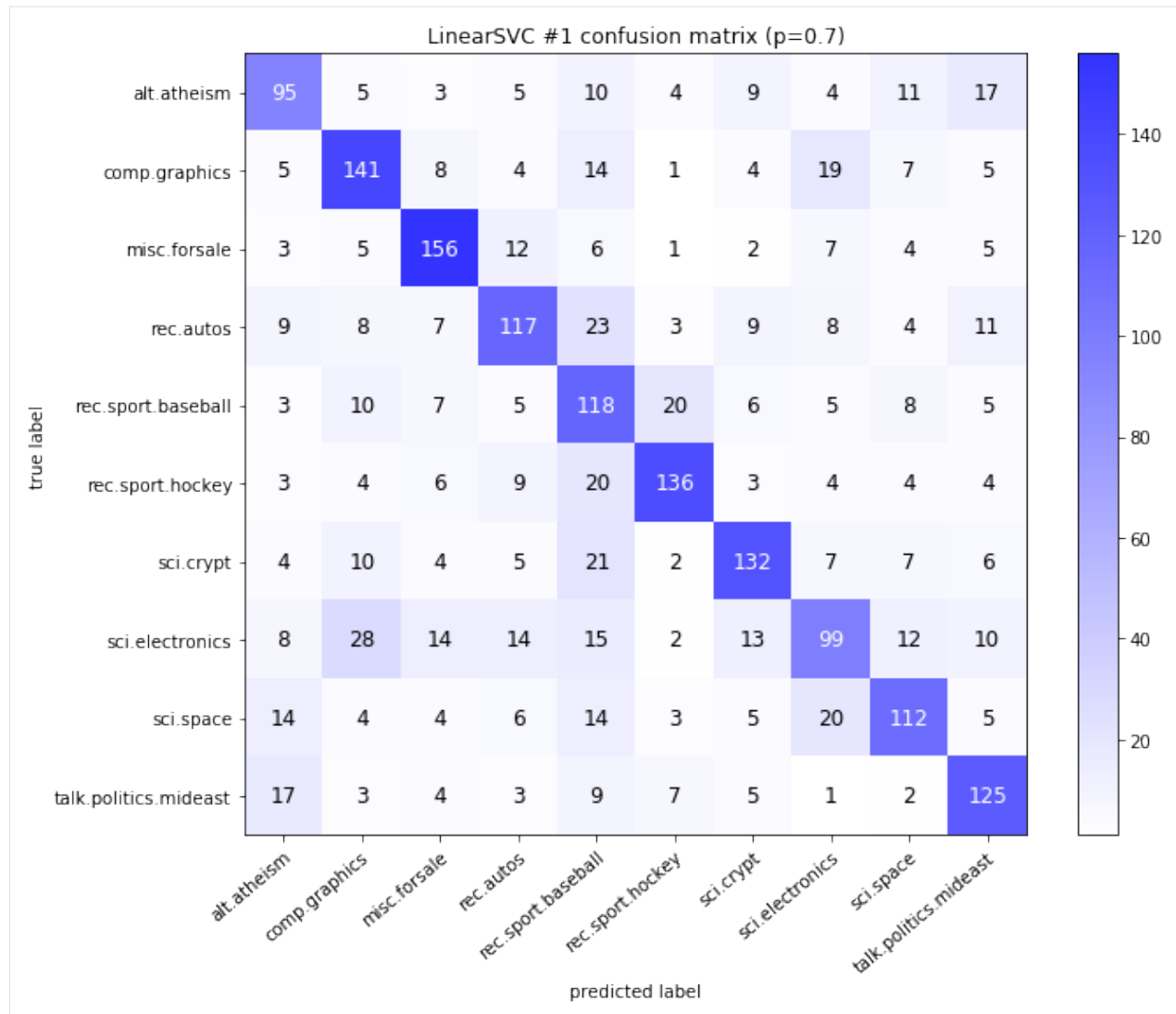


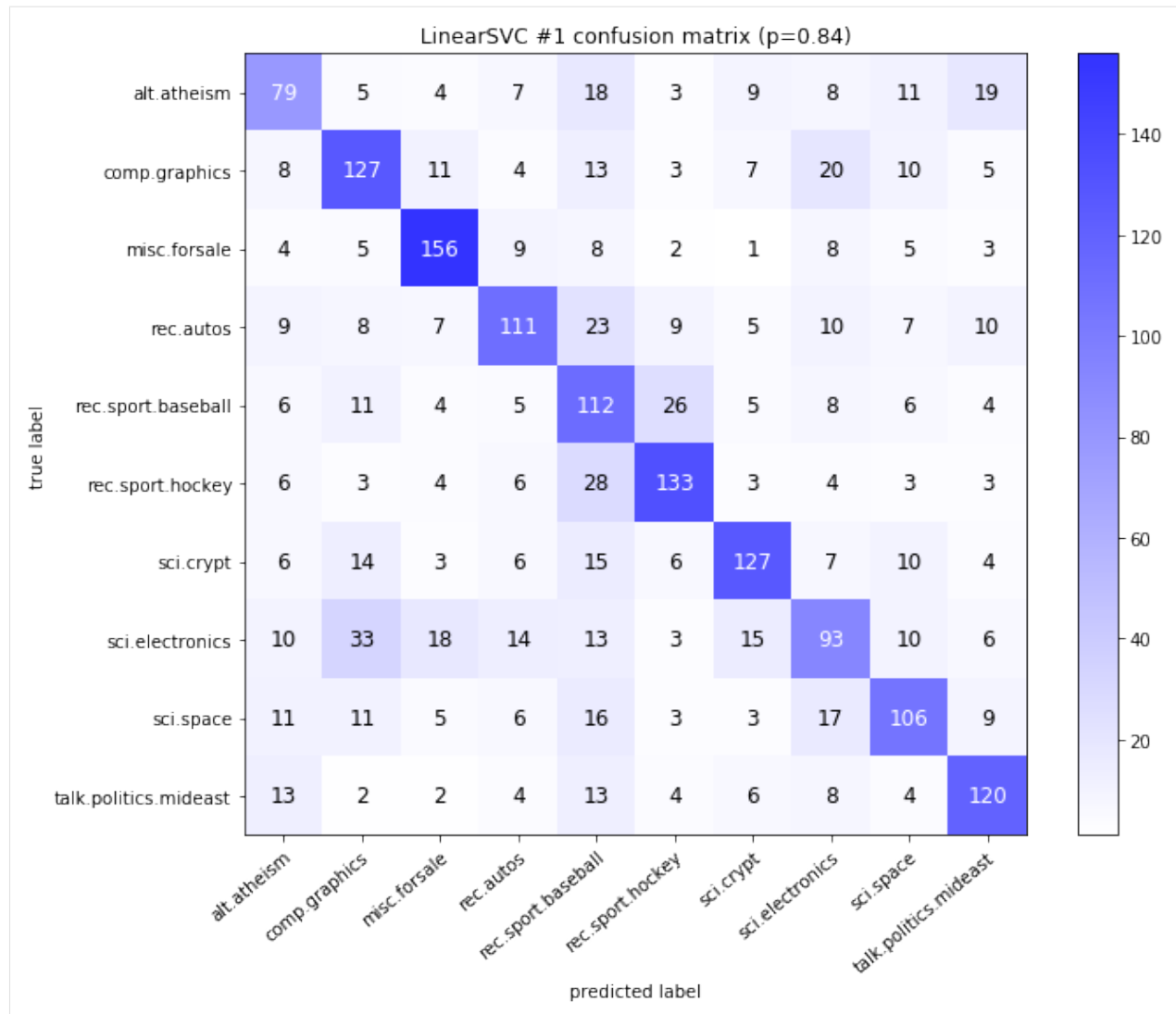


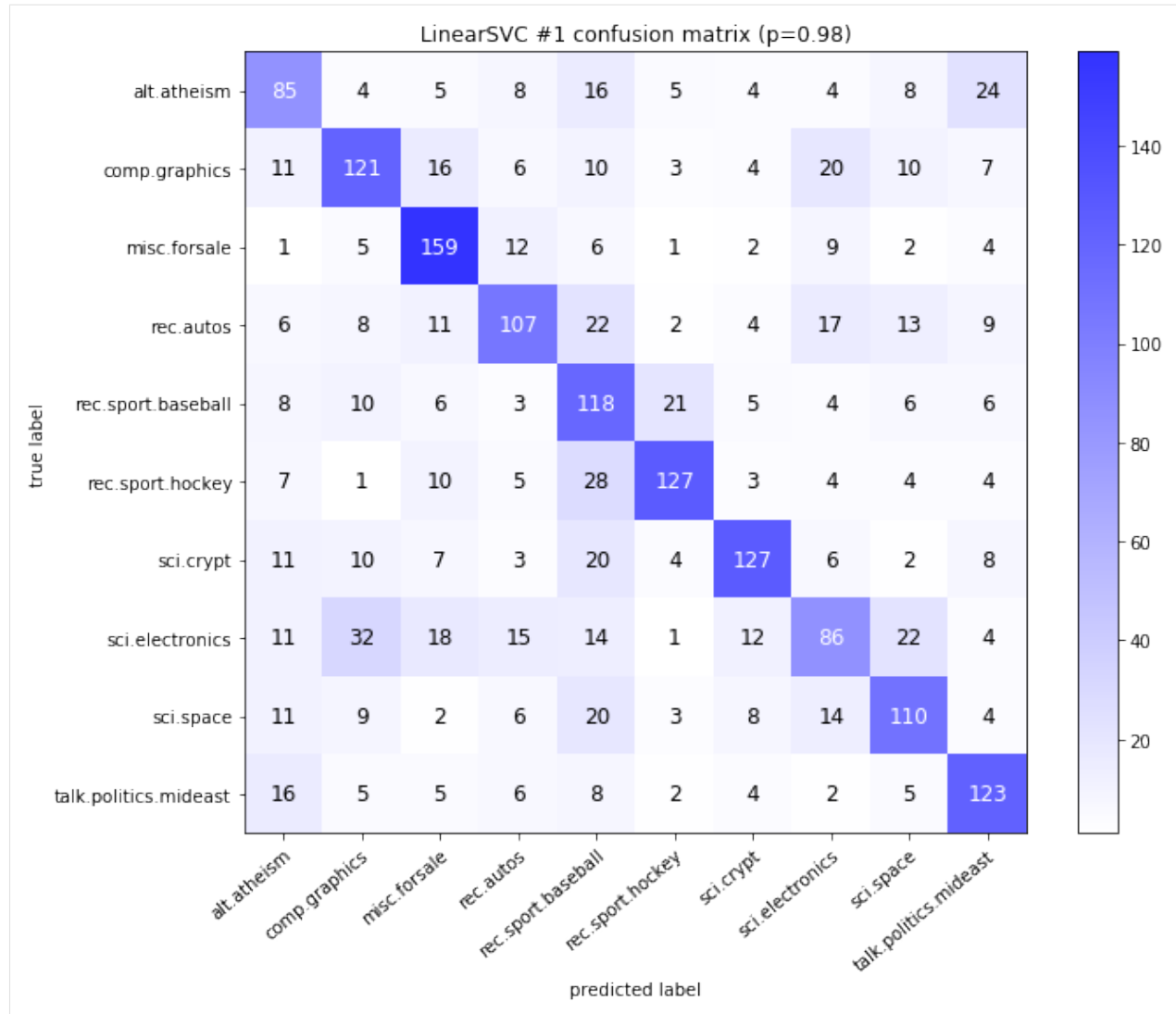


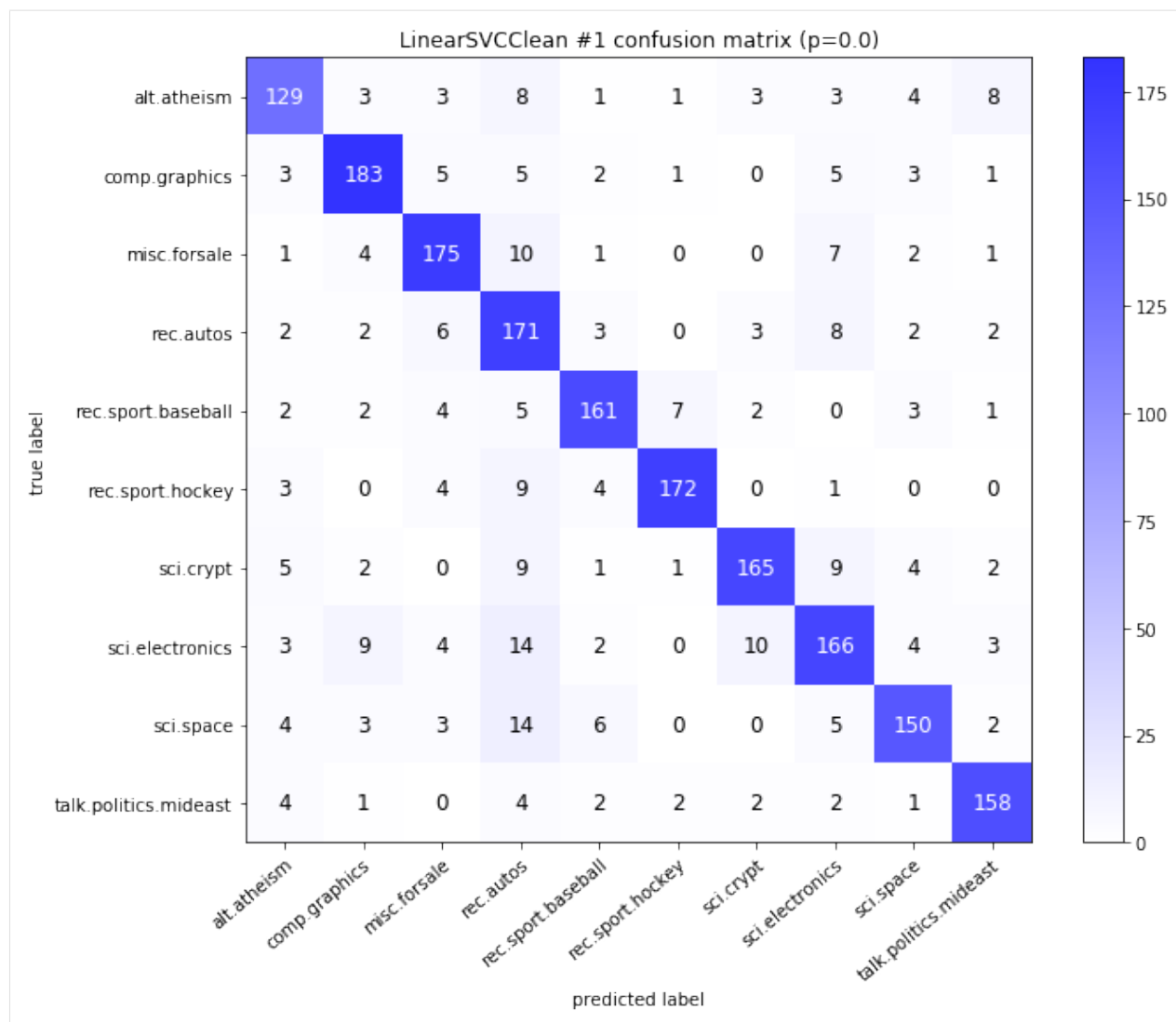


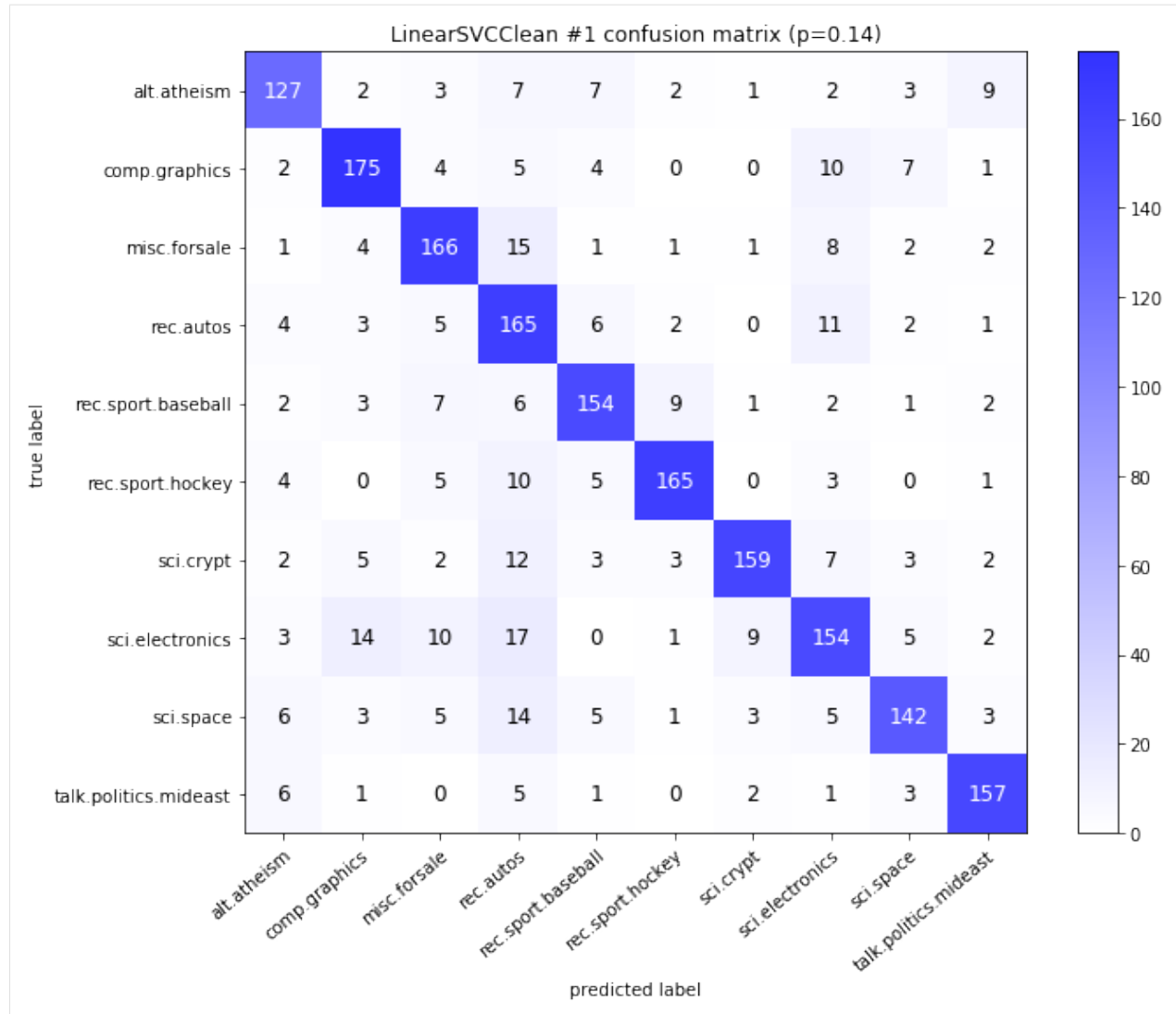


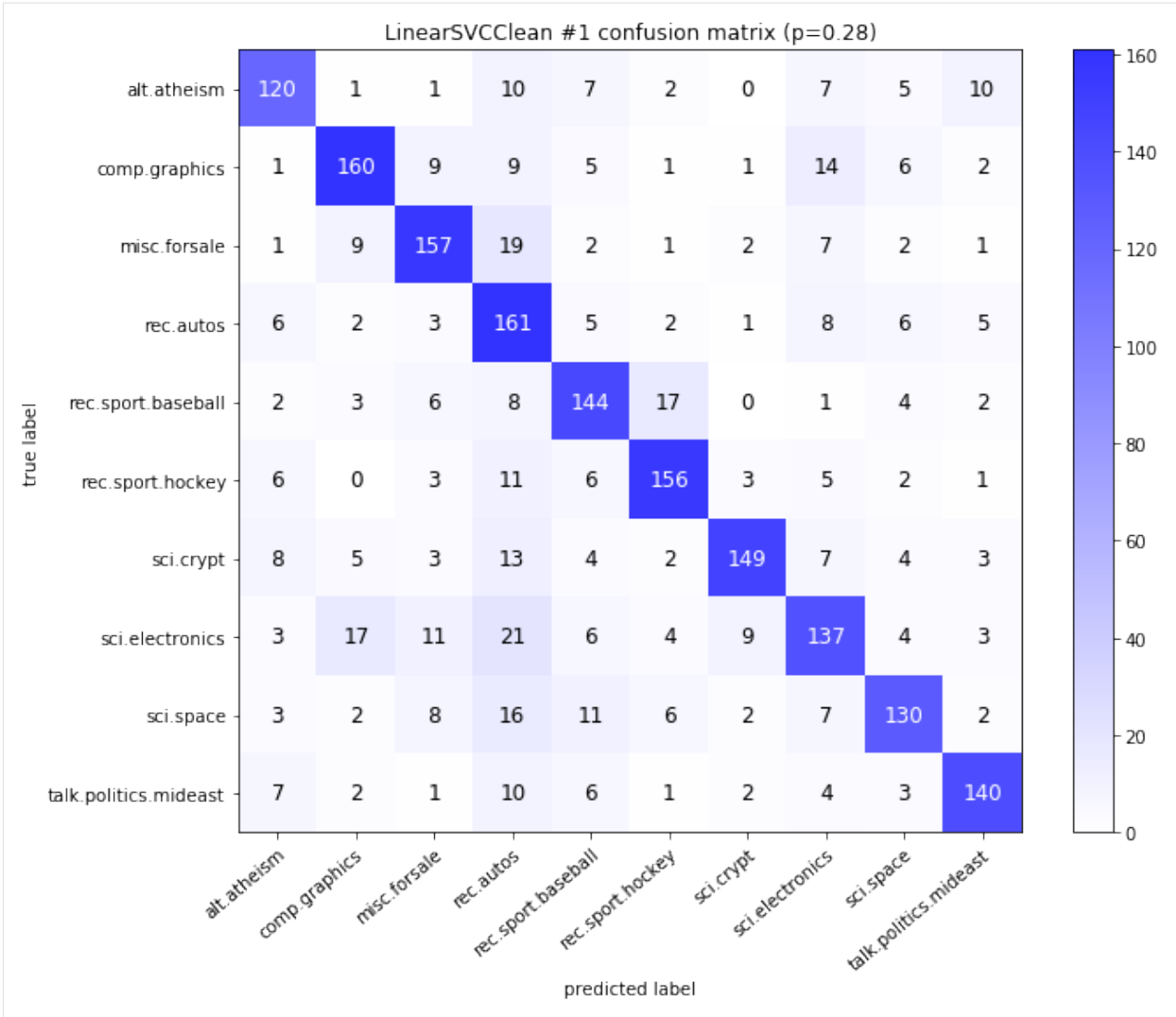


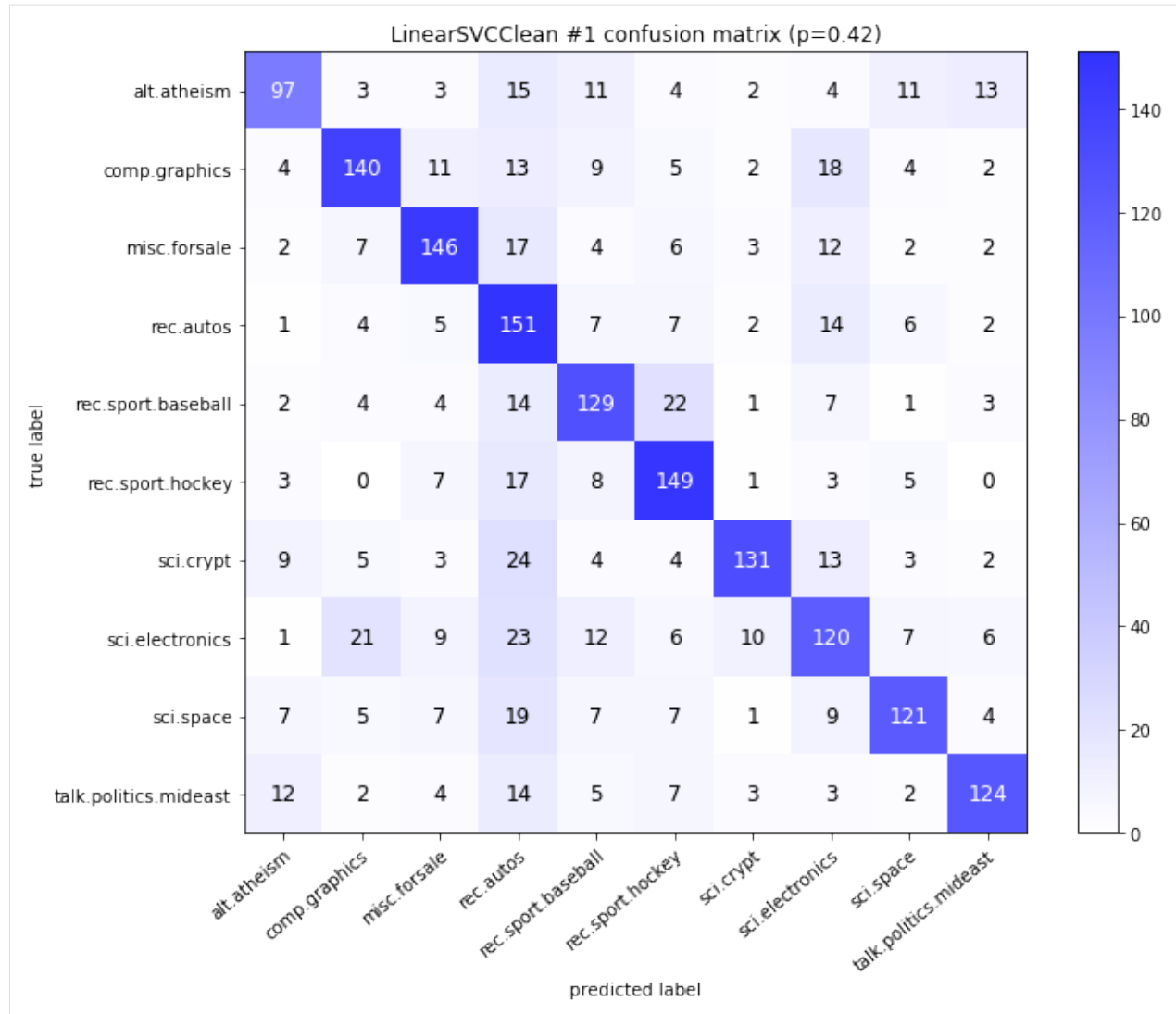


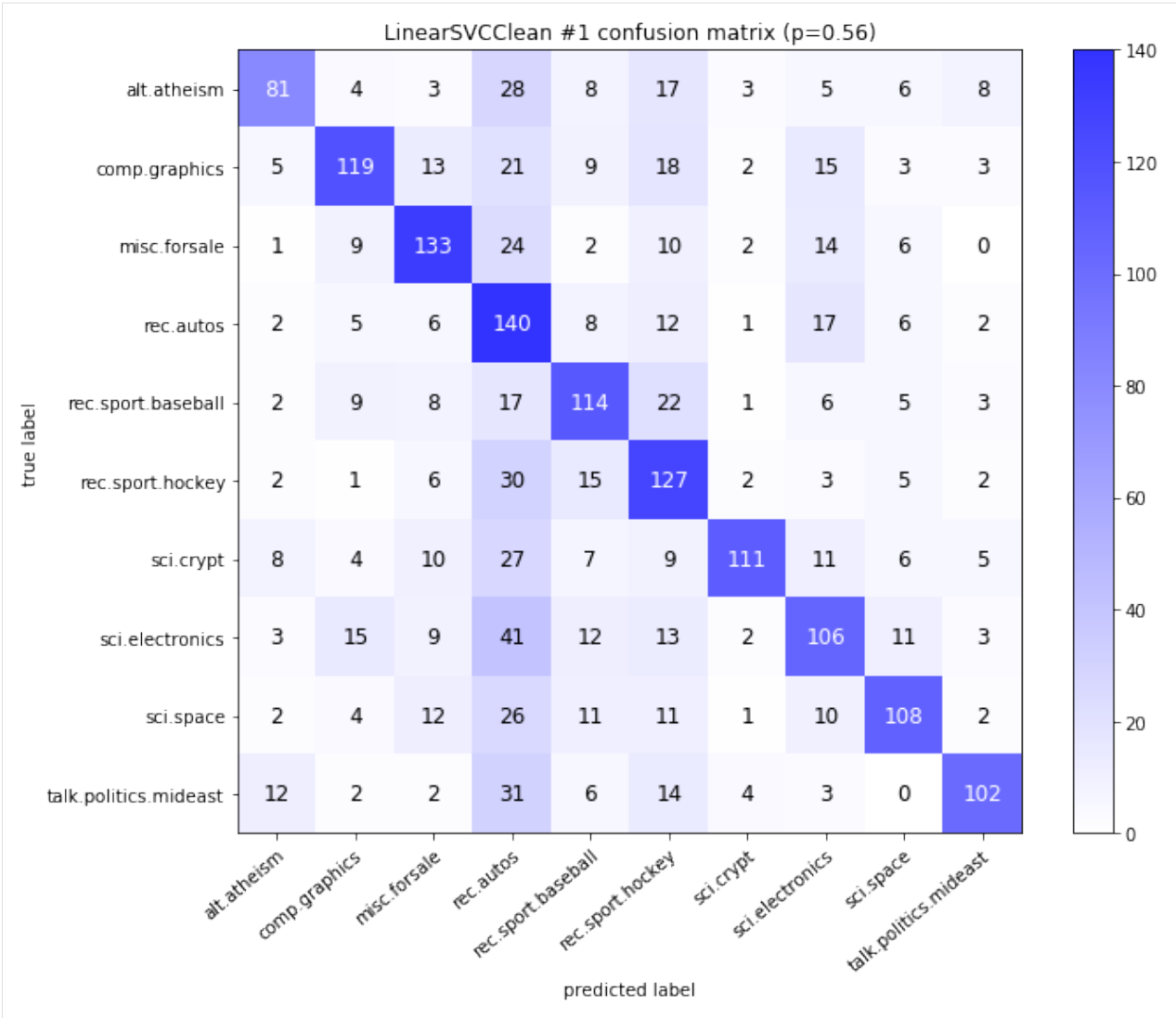


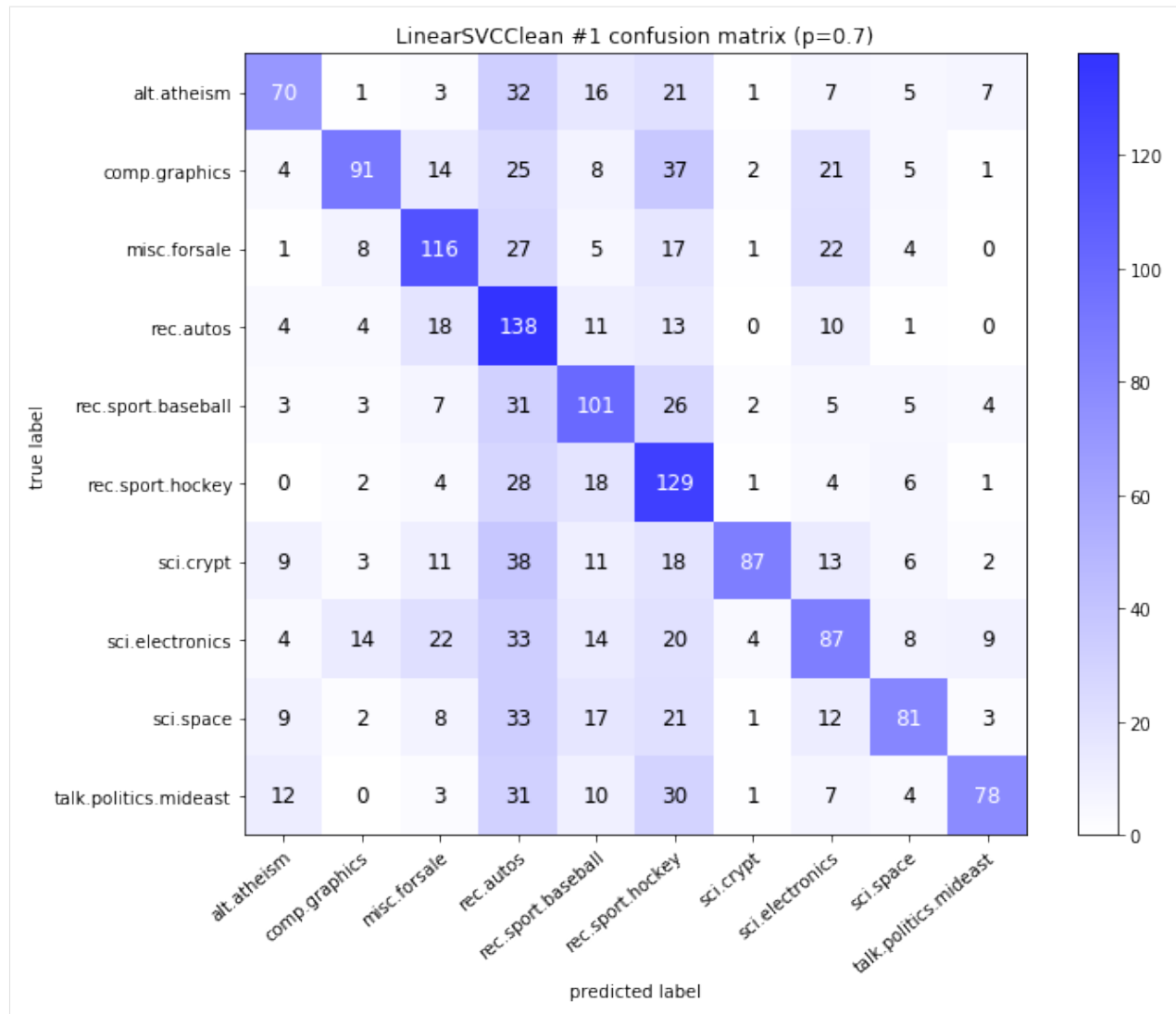


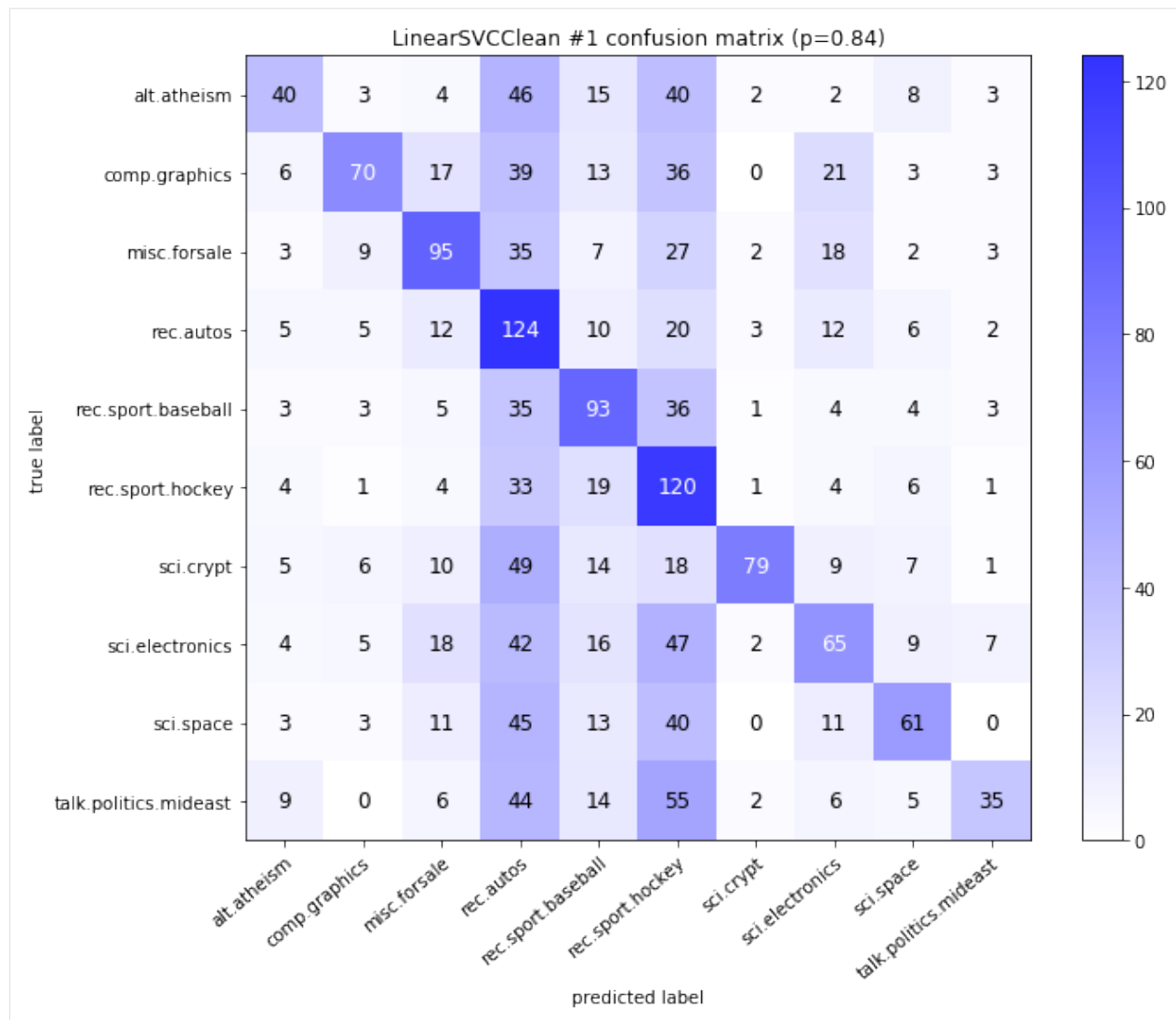


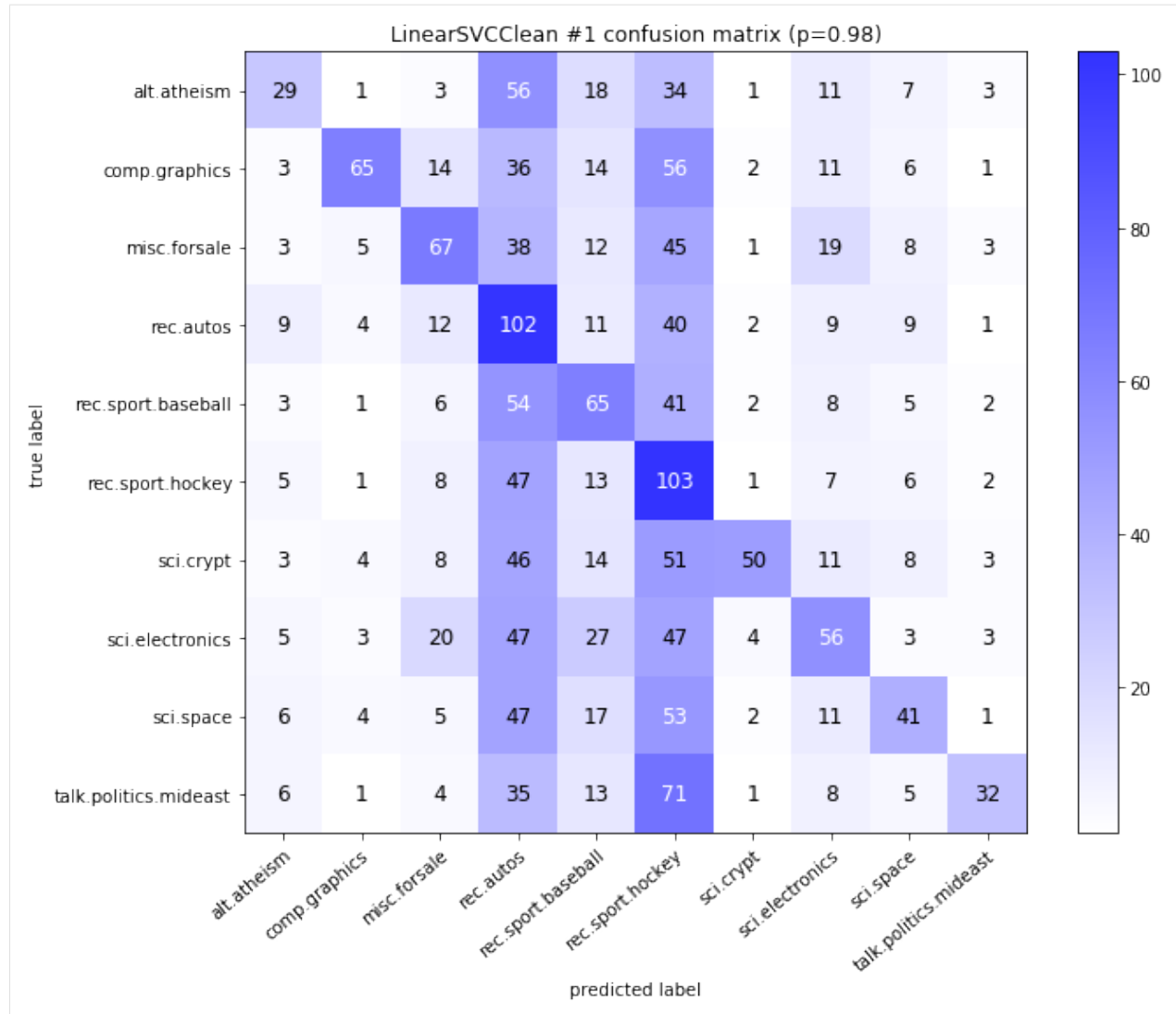


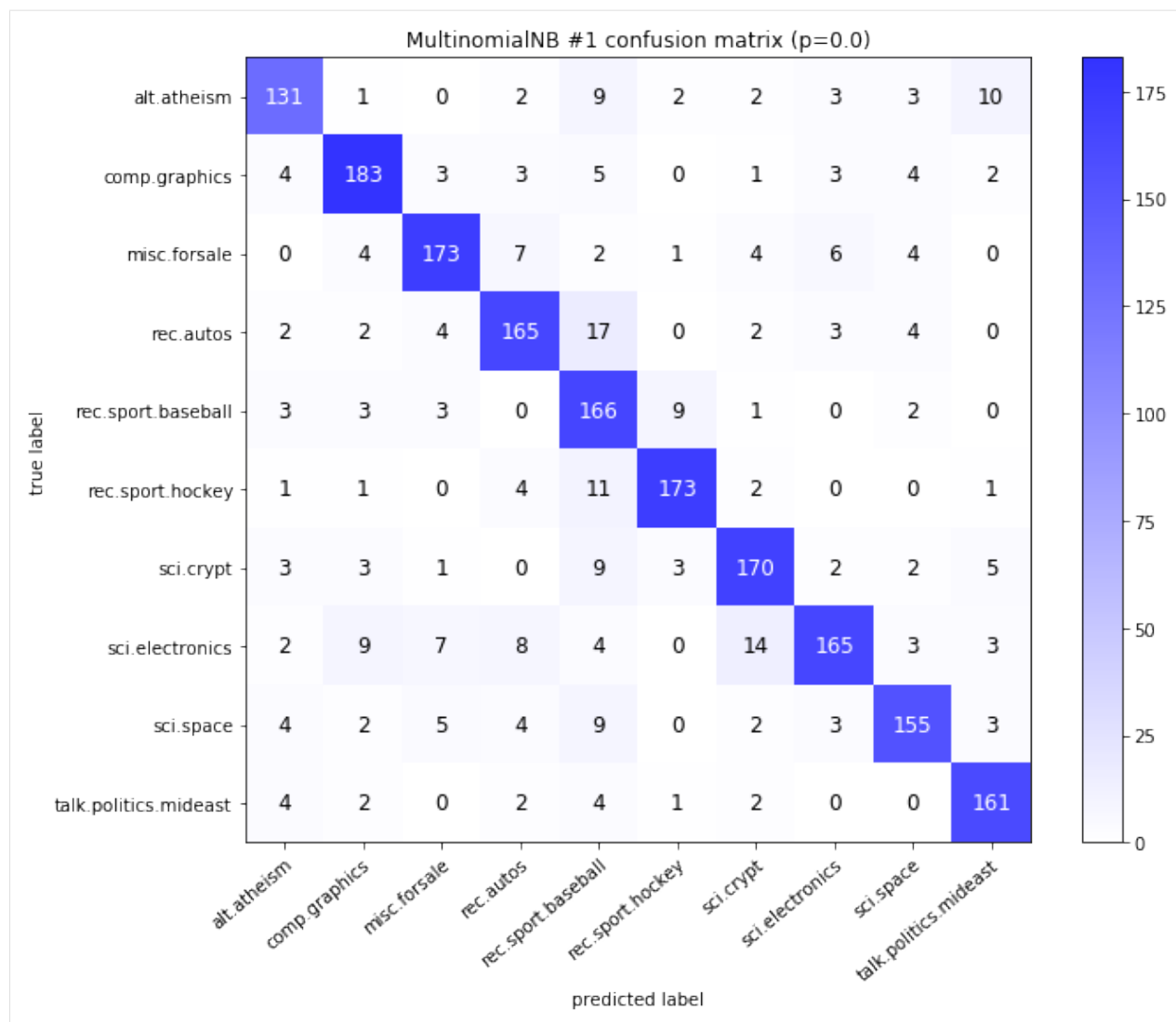


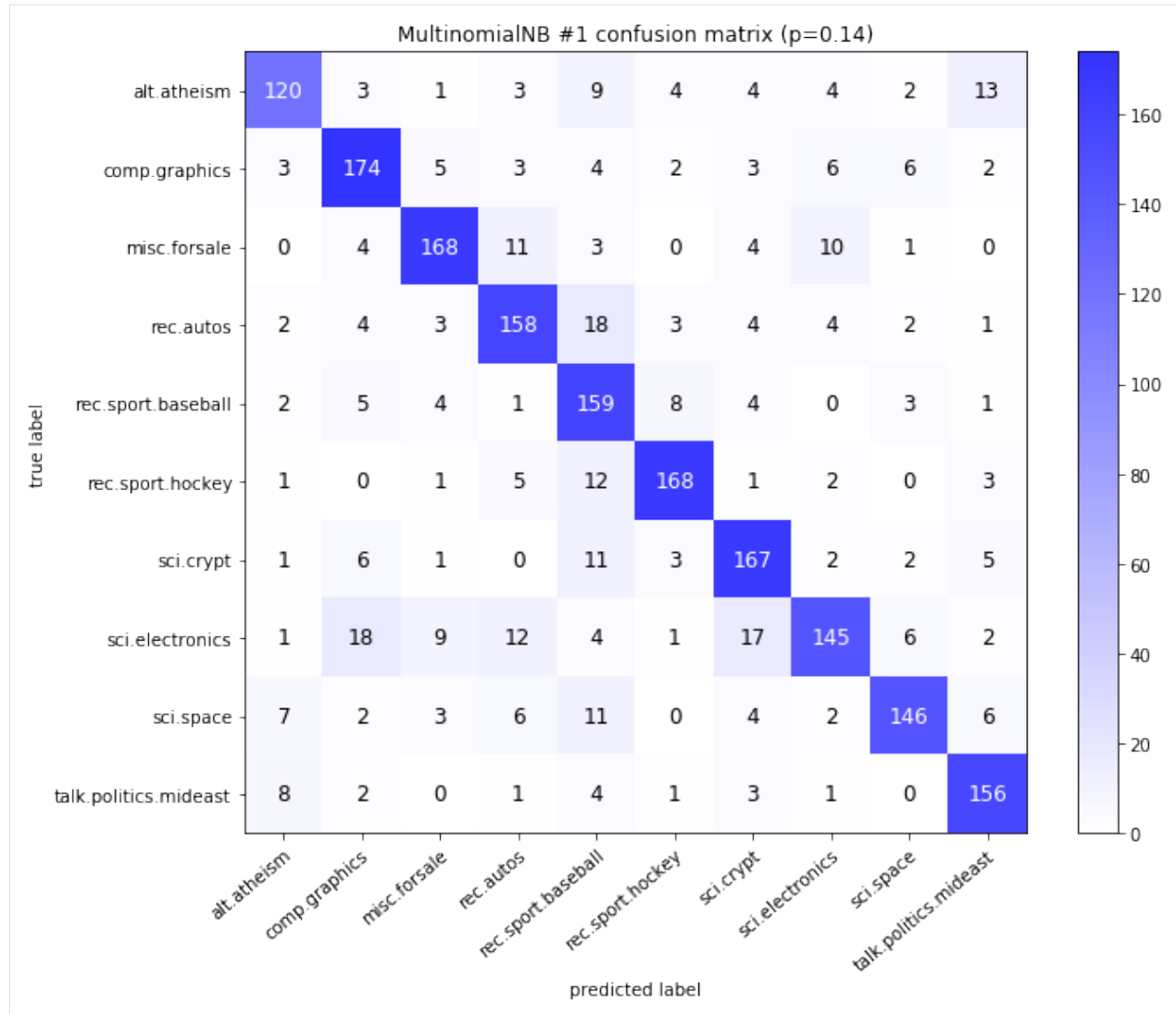


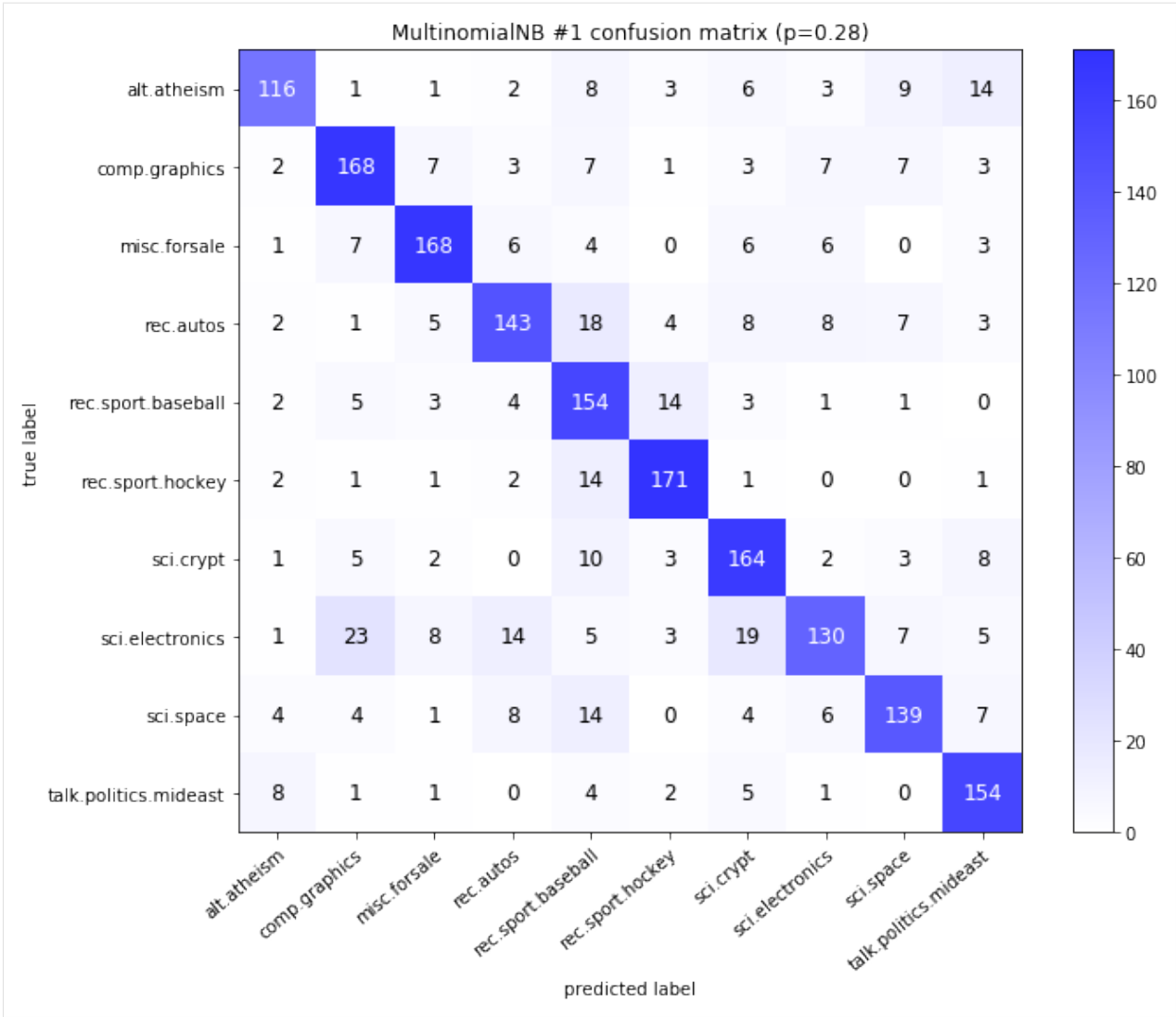


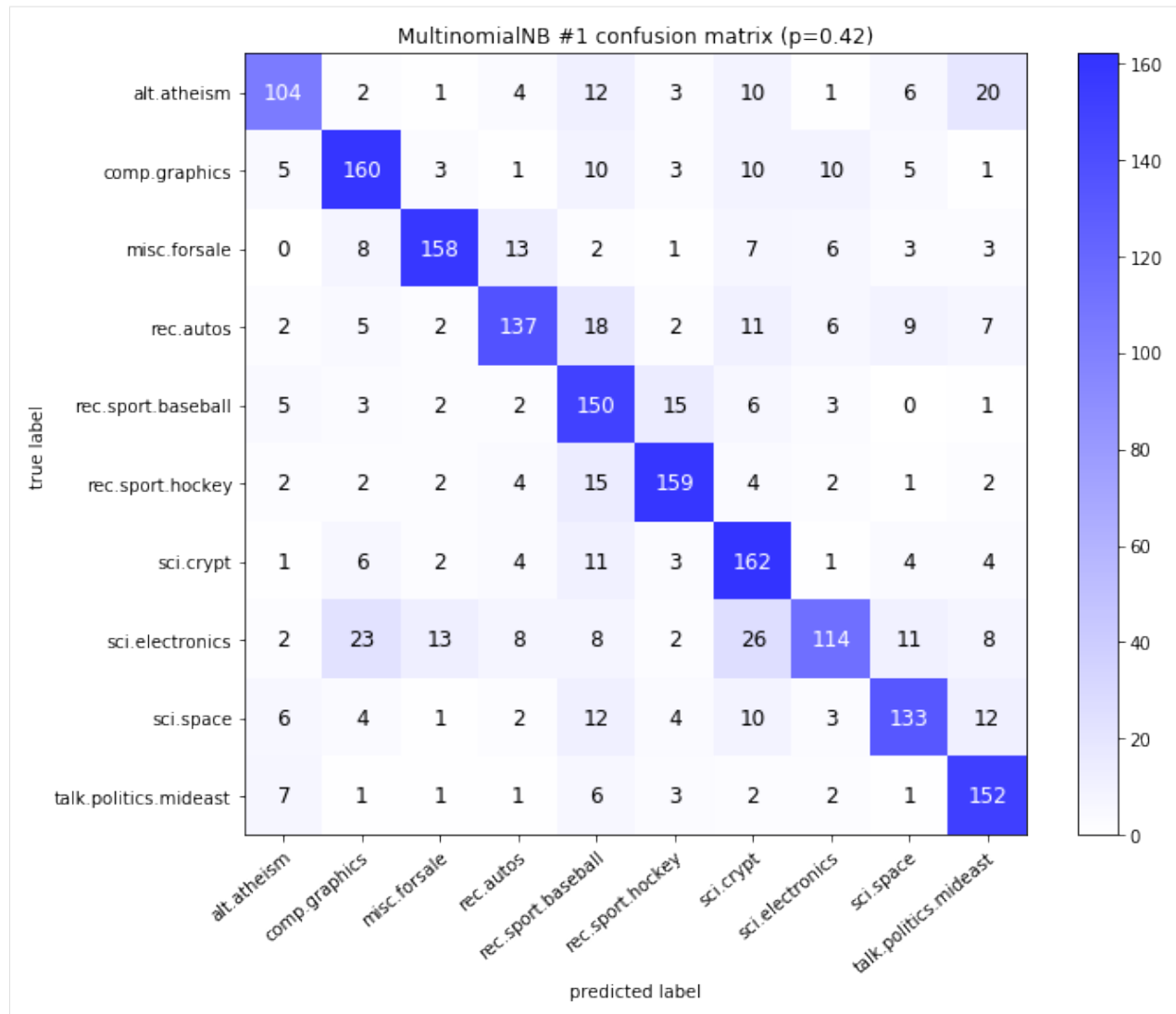


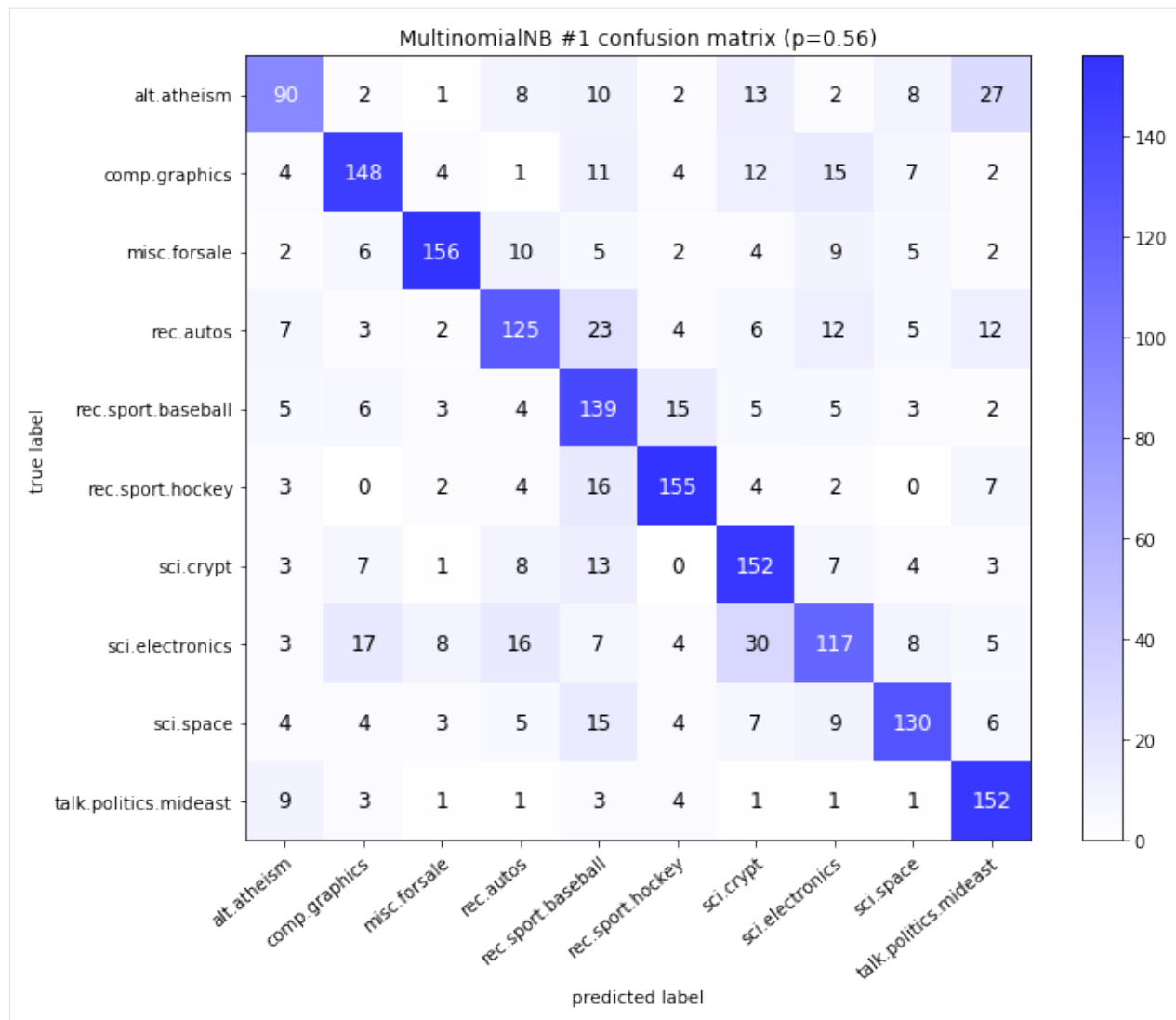


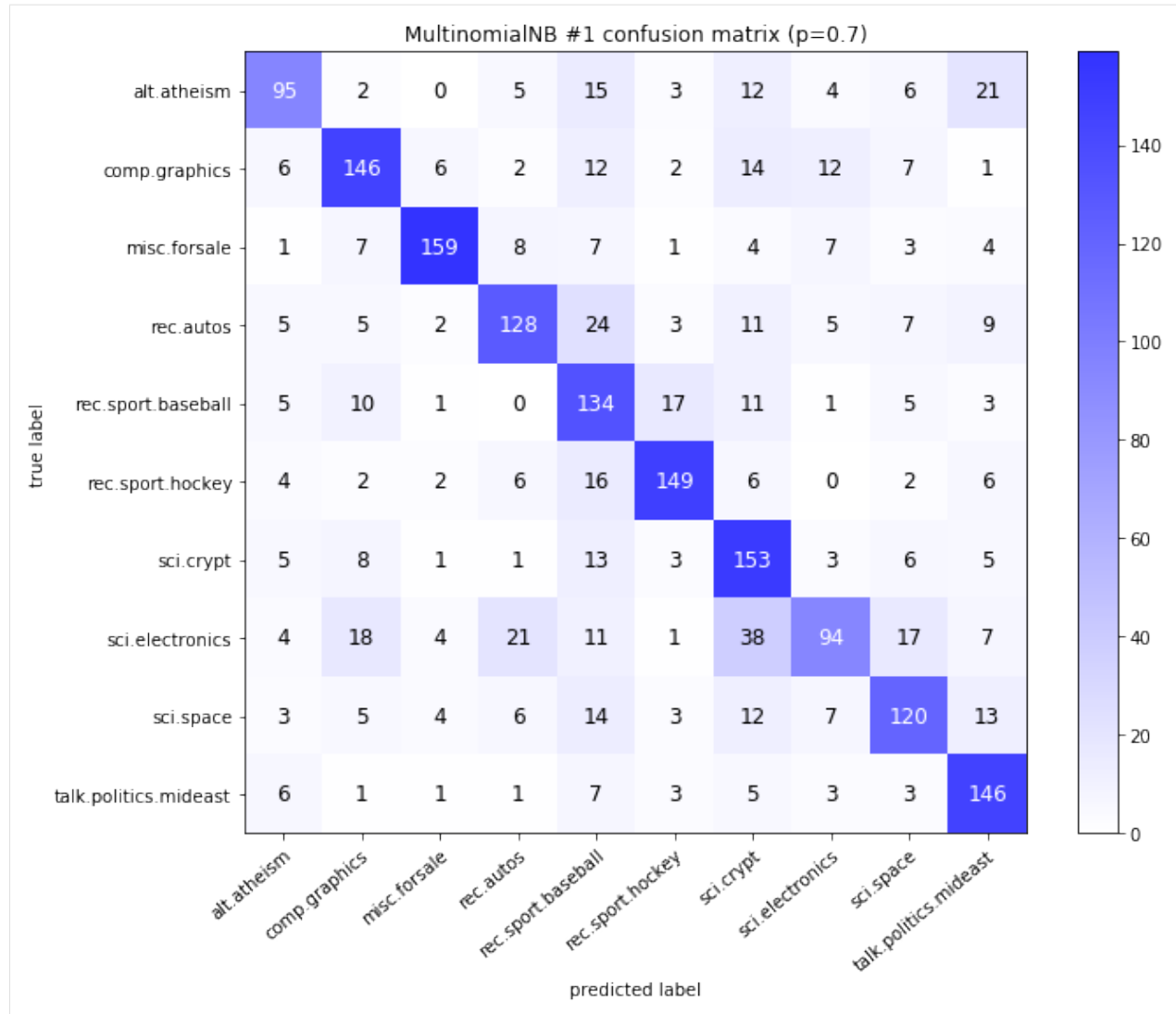


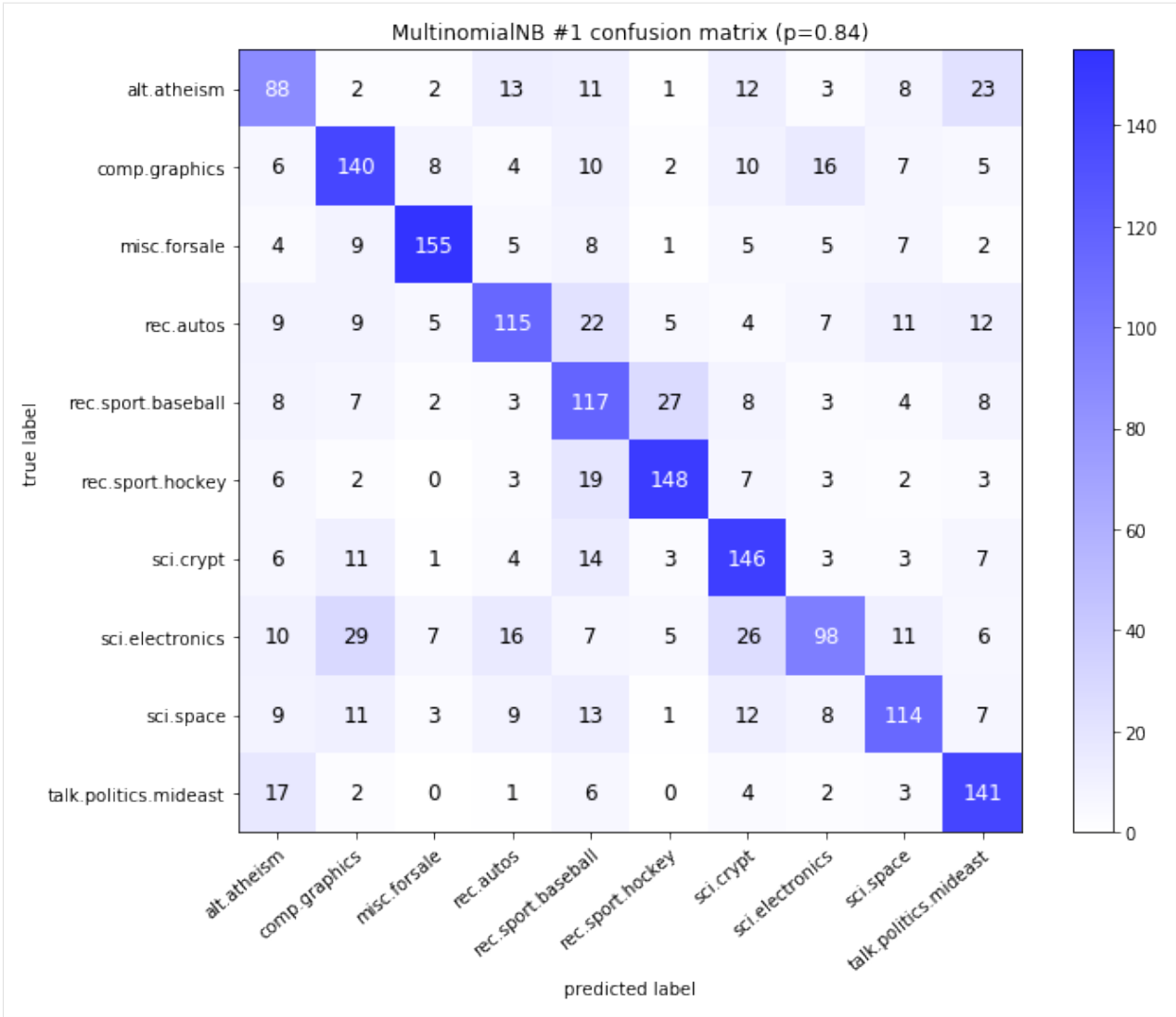


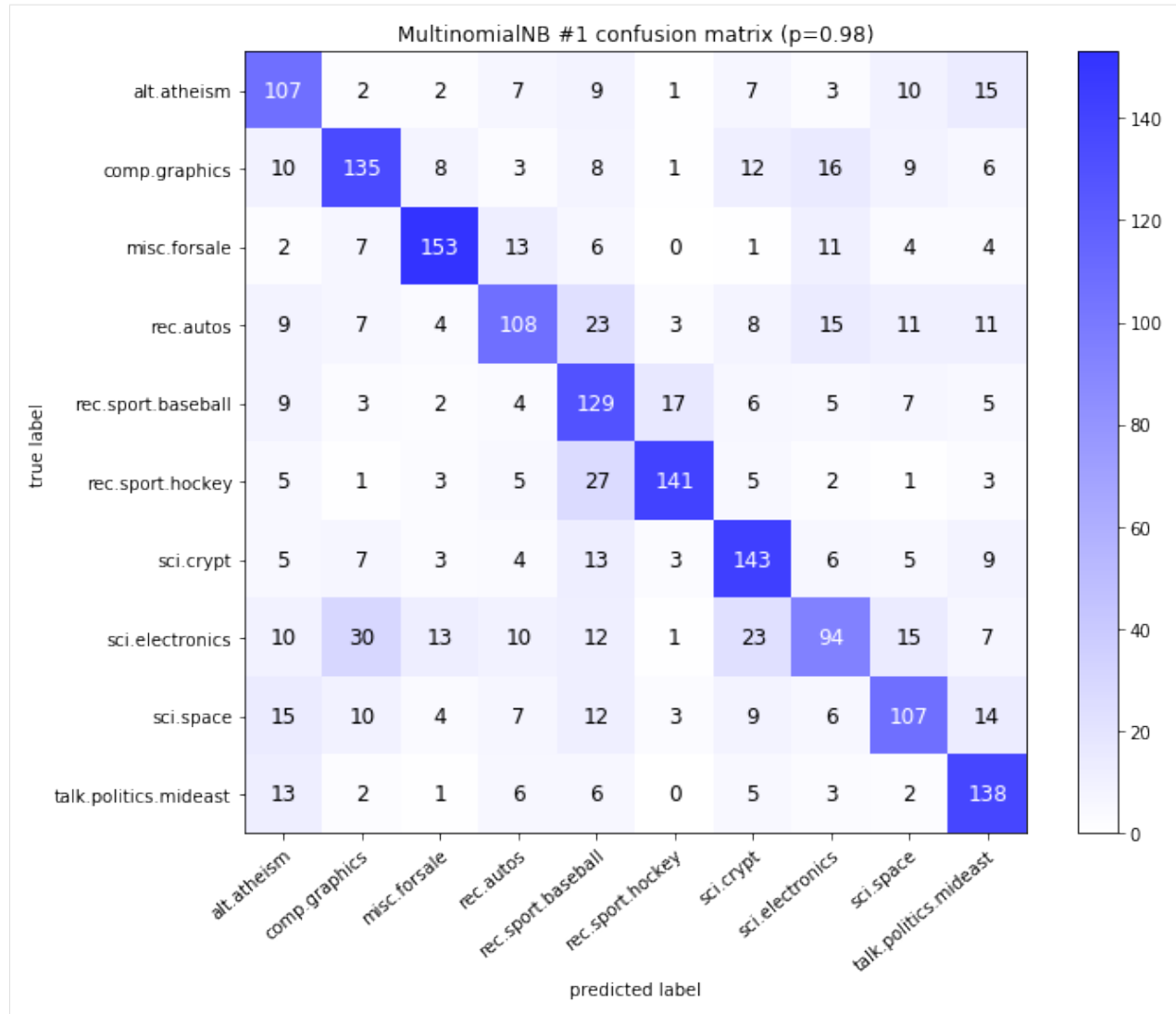


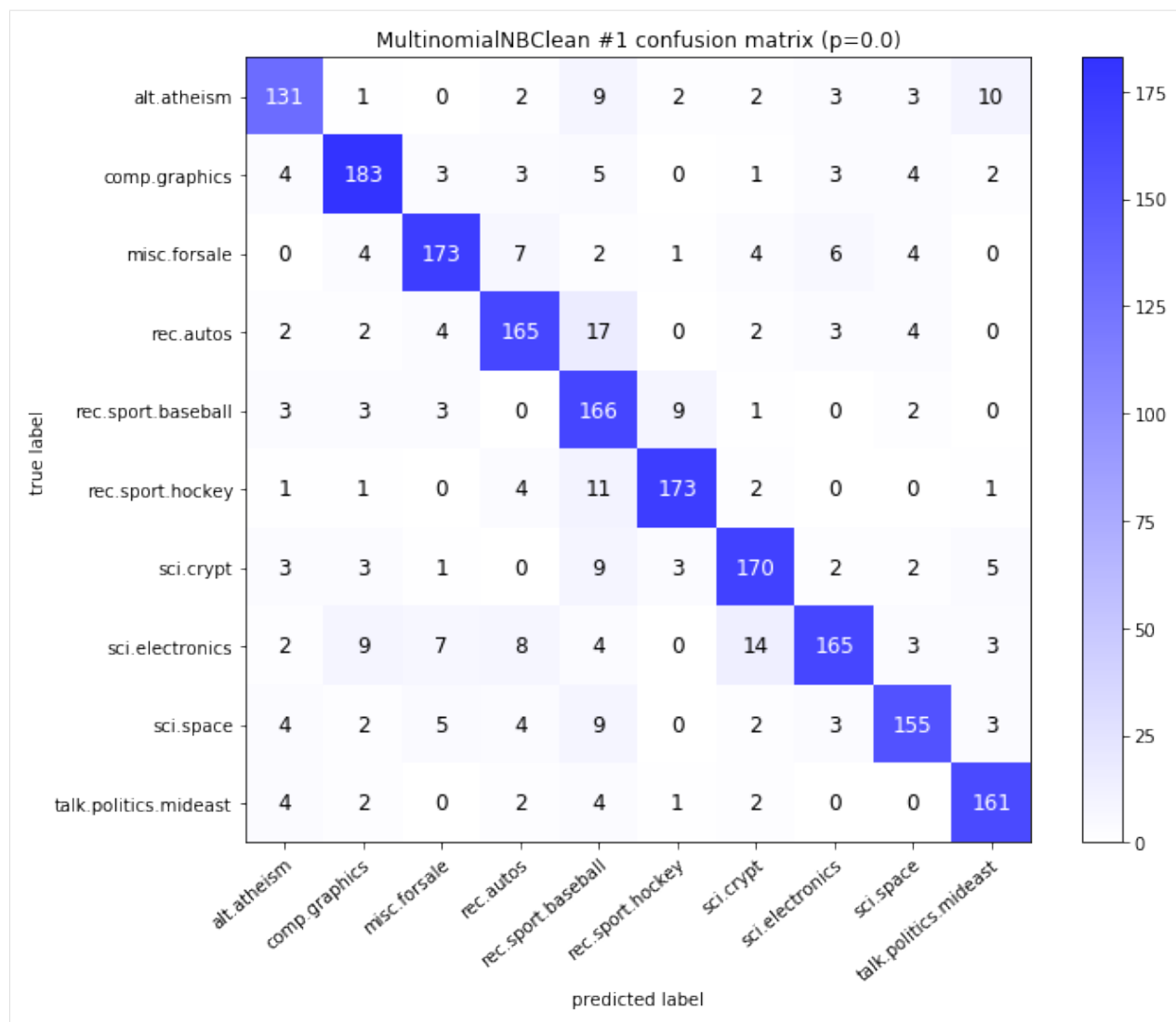


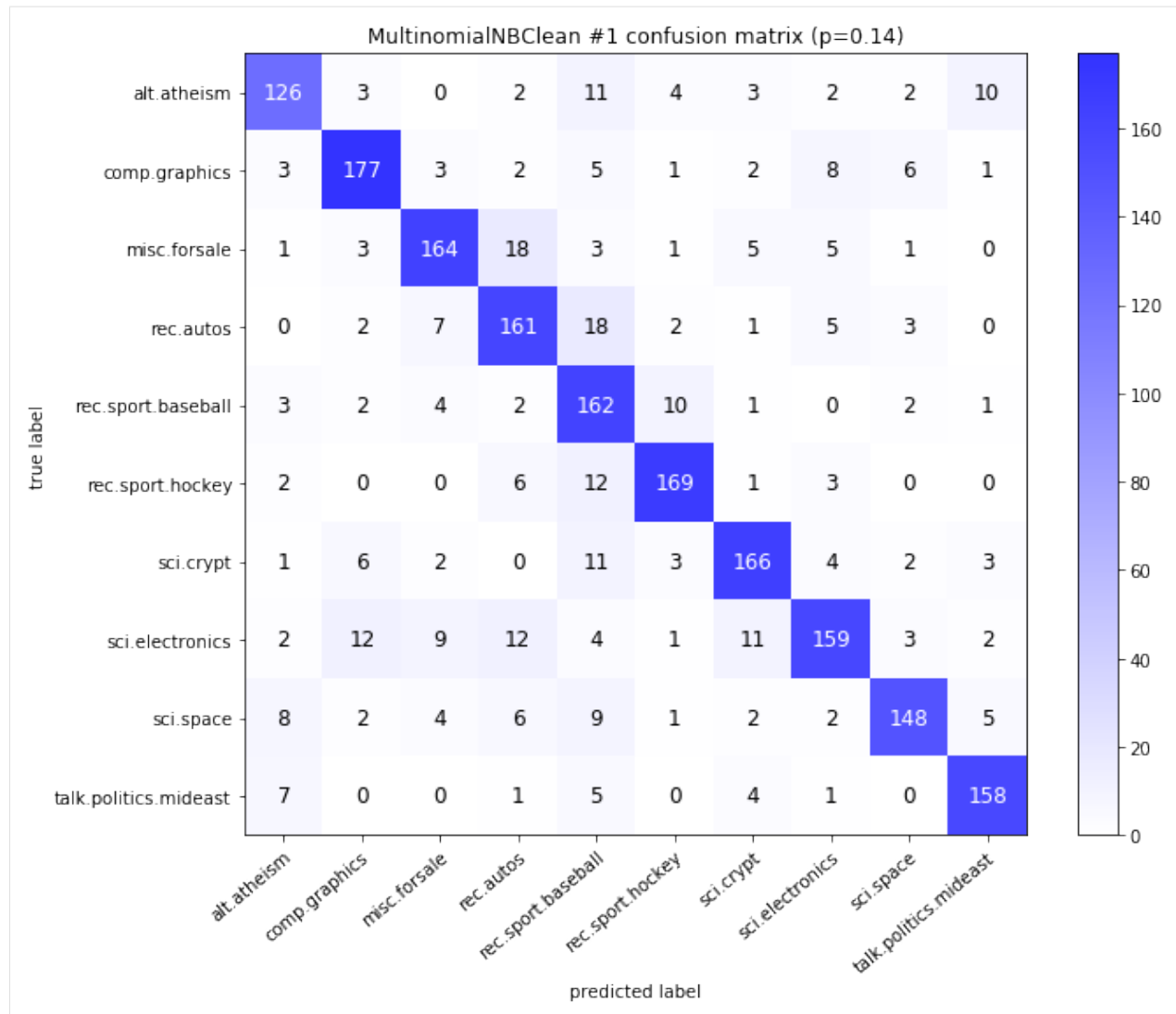


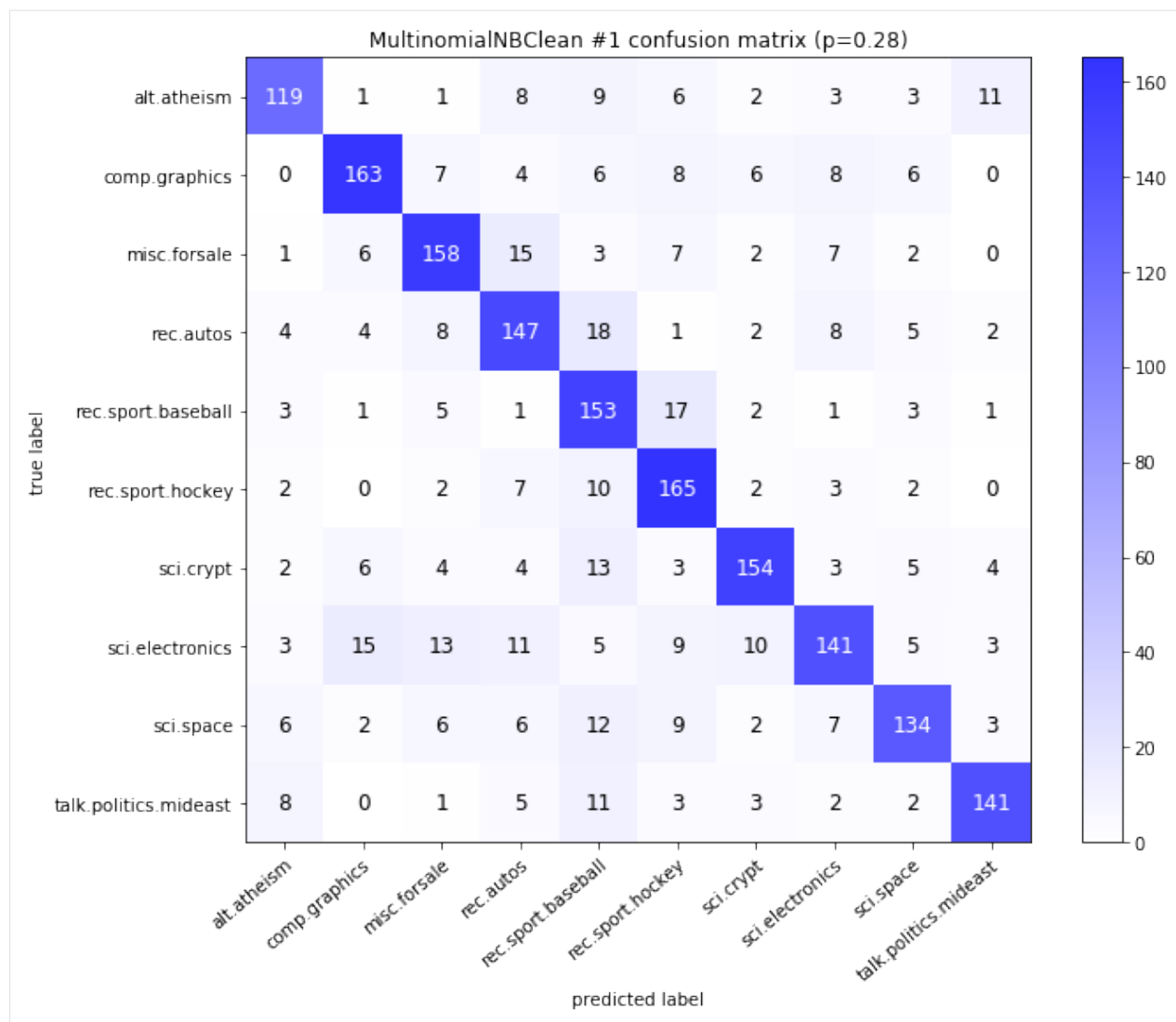


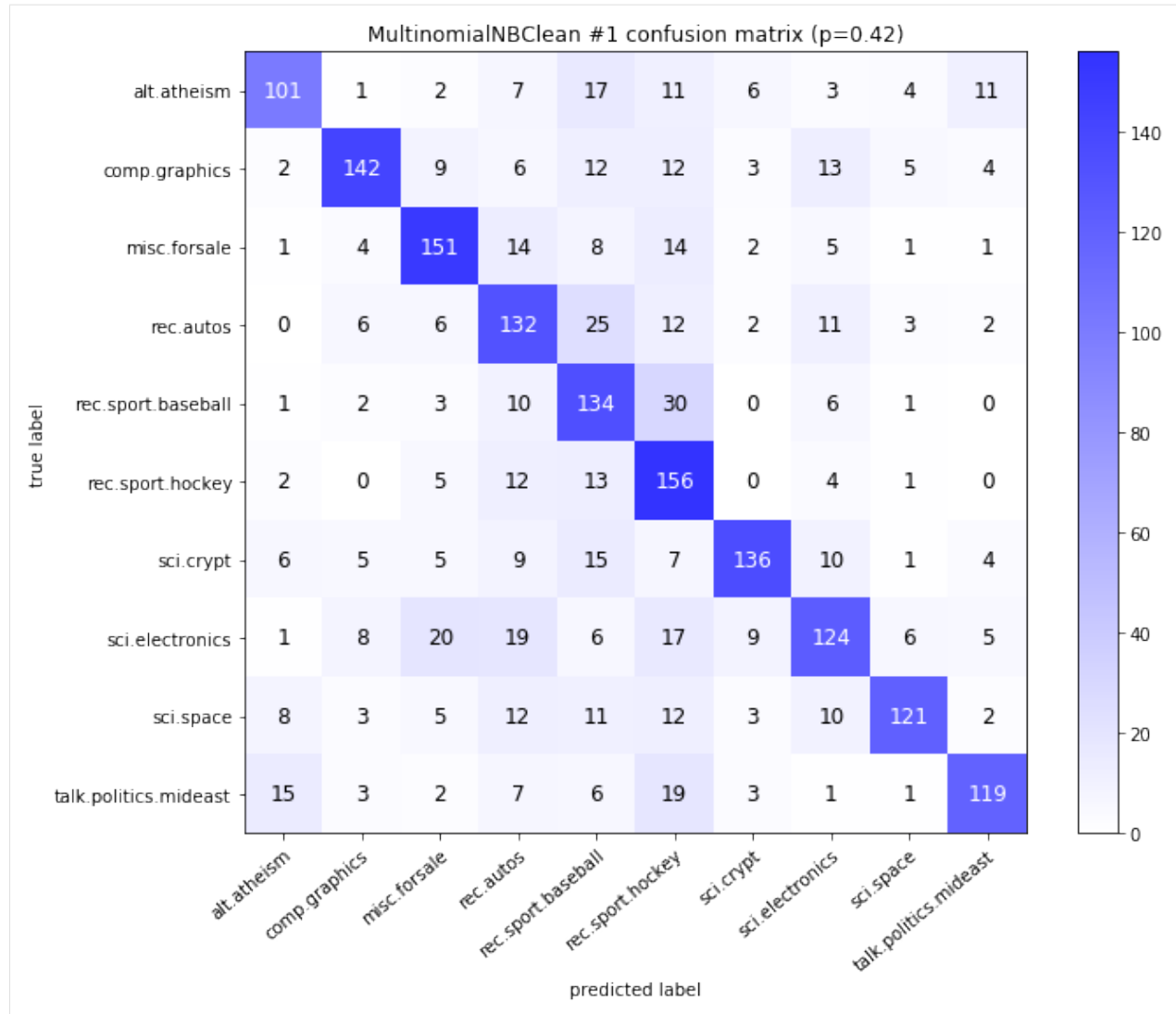


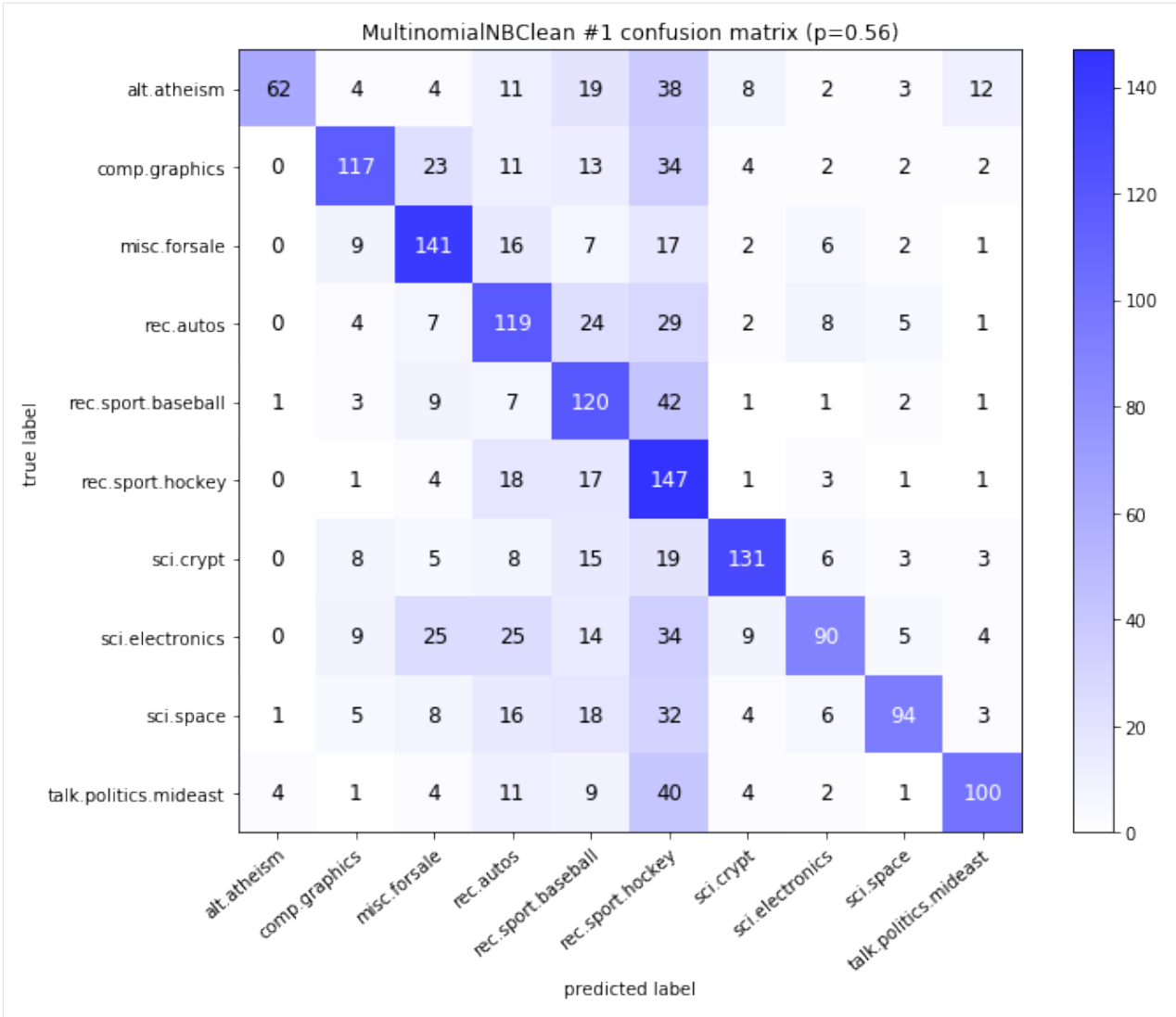


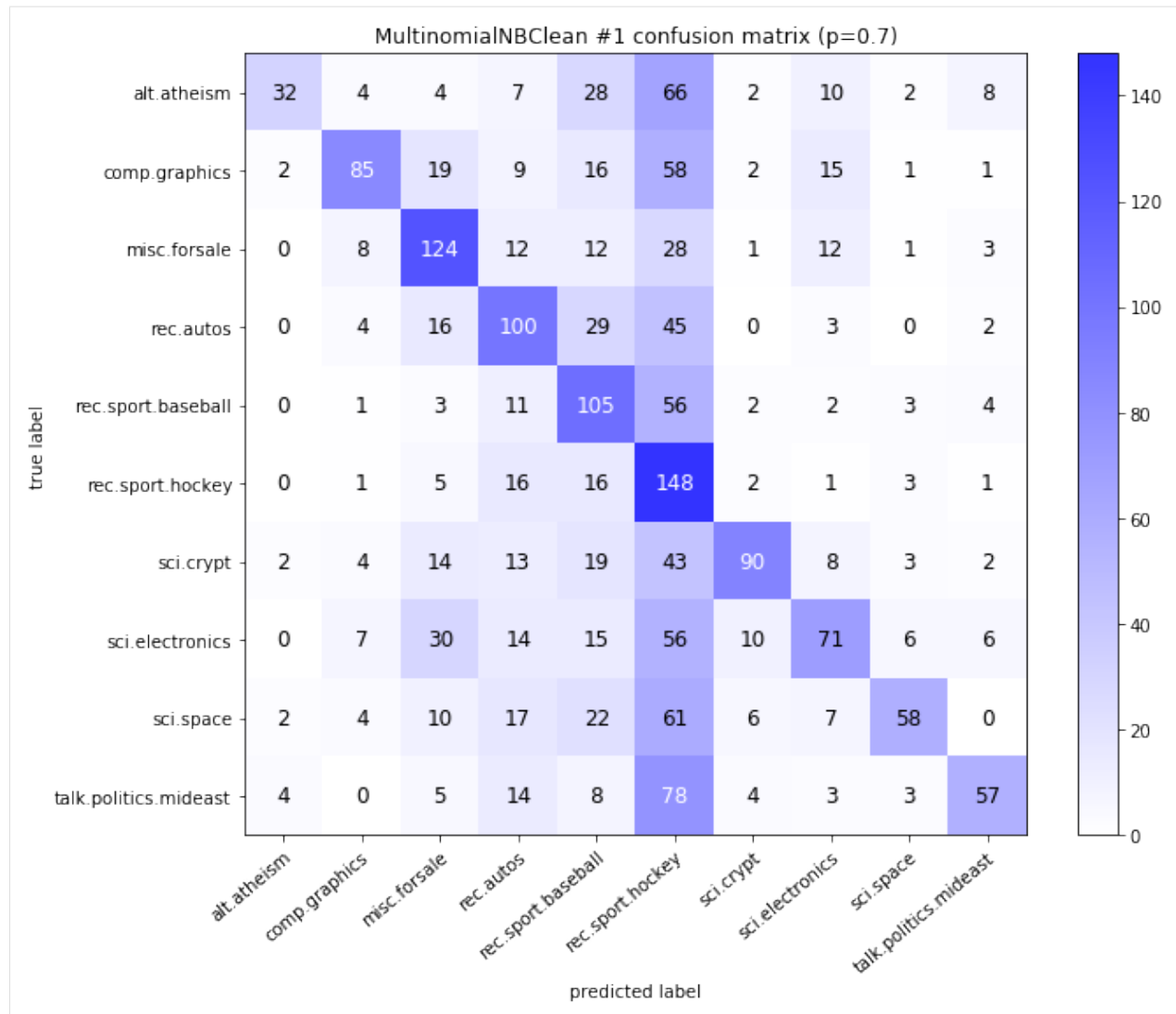


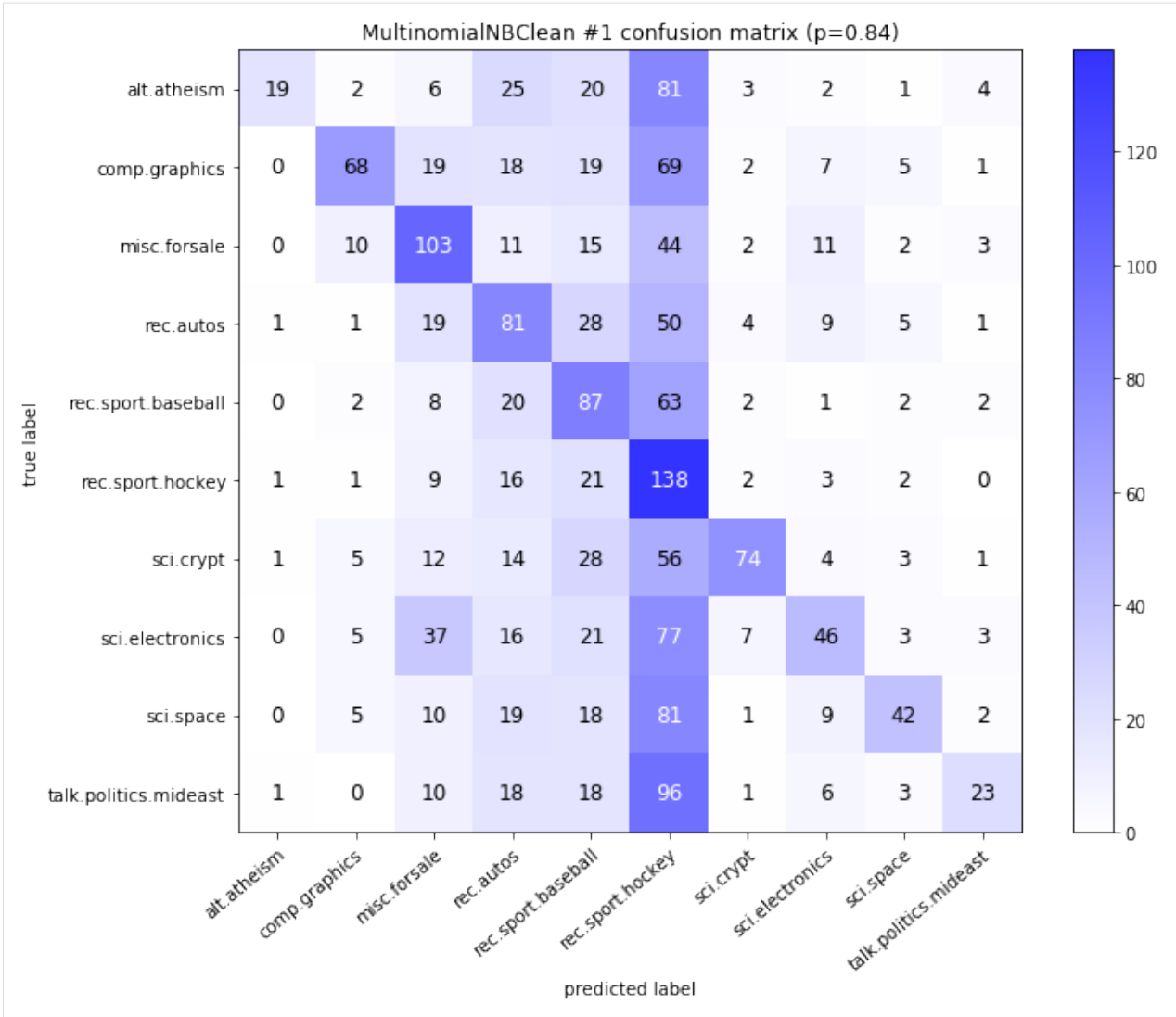


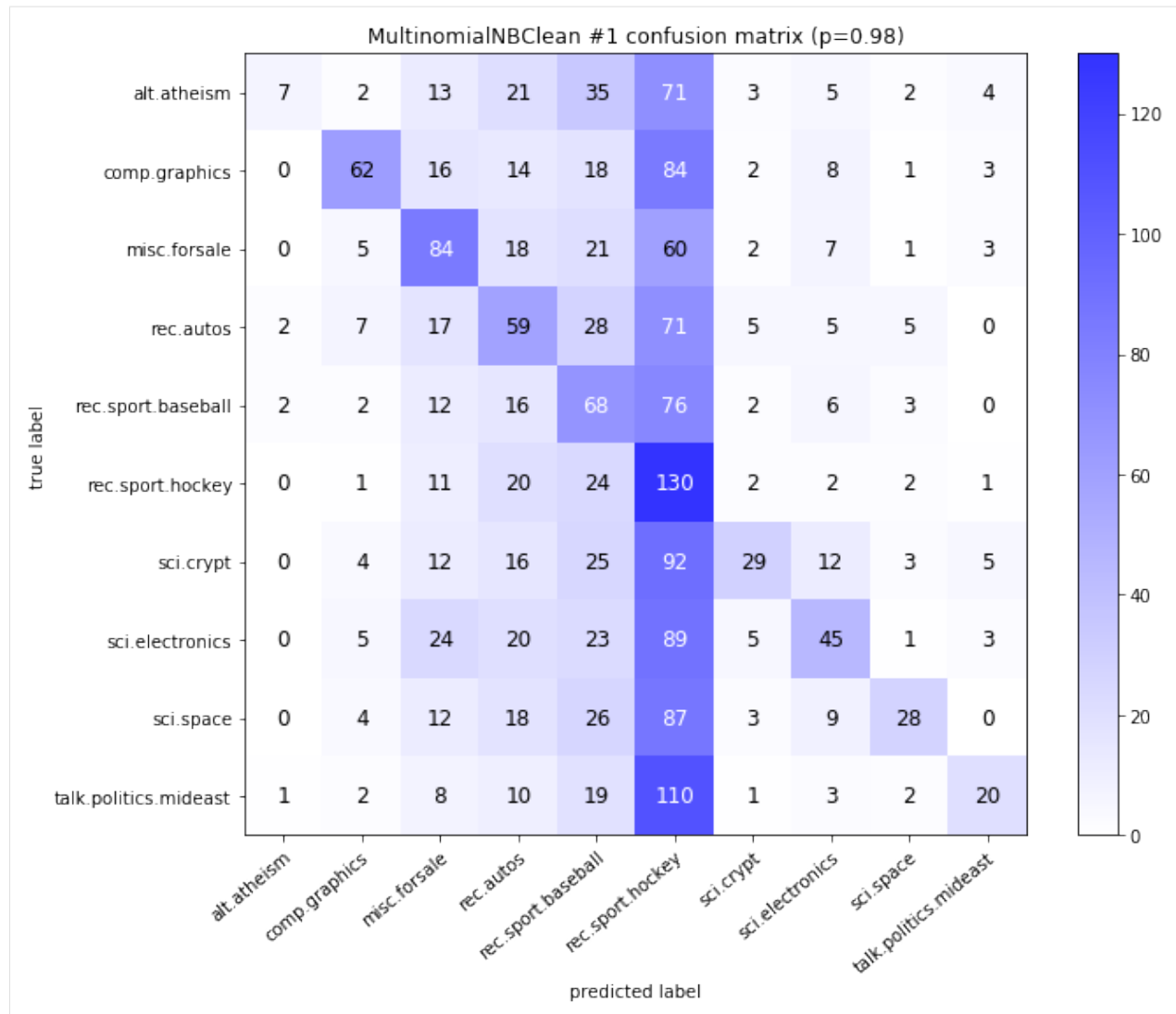












The notebook for this case study can be found [here](#).

4.5 Object detection: Added brightness

Note: *Object detection requirements.*

Warning: Runtimes can be several hours even on clusters.

We compared the performance of models from FaceBook's Detectron project and YOLOv3 model from Joseph Redmon, when different error sources were added. The models from FaceBook's Detectron project were FasterRCNN, MaskRCNN and RetinaNet.

```
[1]: import re
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
import numpy as np
from PIL import Image

from dpemu import runner
from dpemu.dataset_utils import load_coco_val_2017
from dpemu.filters.image import Brightness
from dpemu.ml_utils import run_ml_module_using_cli, load_yolov3
from dpemu.nodes import Array, Series
from dpemu.plotting_utils import print_results_by_model, visualize_scores
from dpemu.utils import get_project_root
```

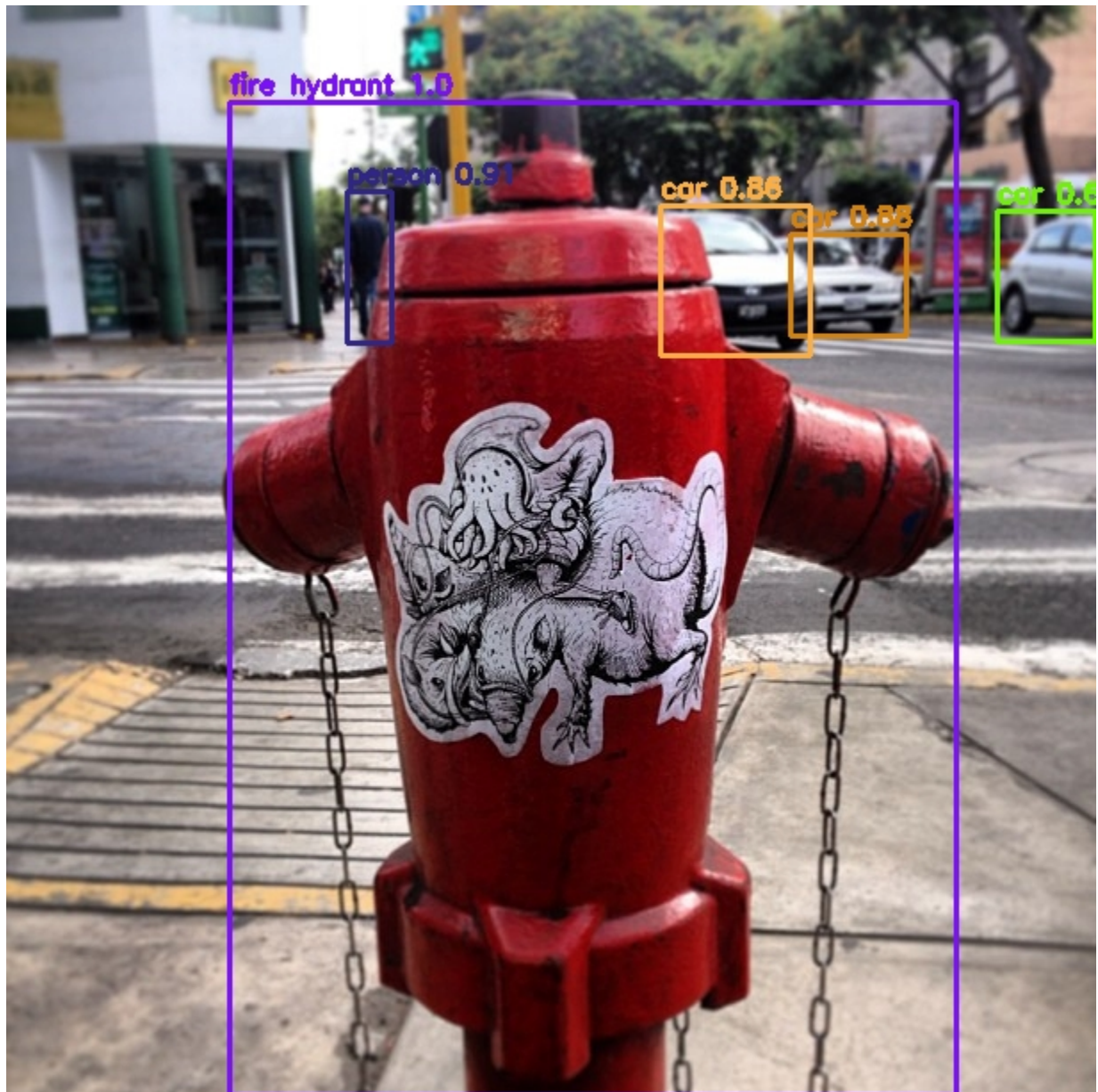
We used 118 287 jpg images (COCO train2017) as the train set and 5000 images (COCO val2017) as the test set to calculate the mAP-50 scores.

```
[ ]: def get_data():
    imgs, _, _, img_filenames = load_coco_val_2017()
    return imgs, img_filenames
```

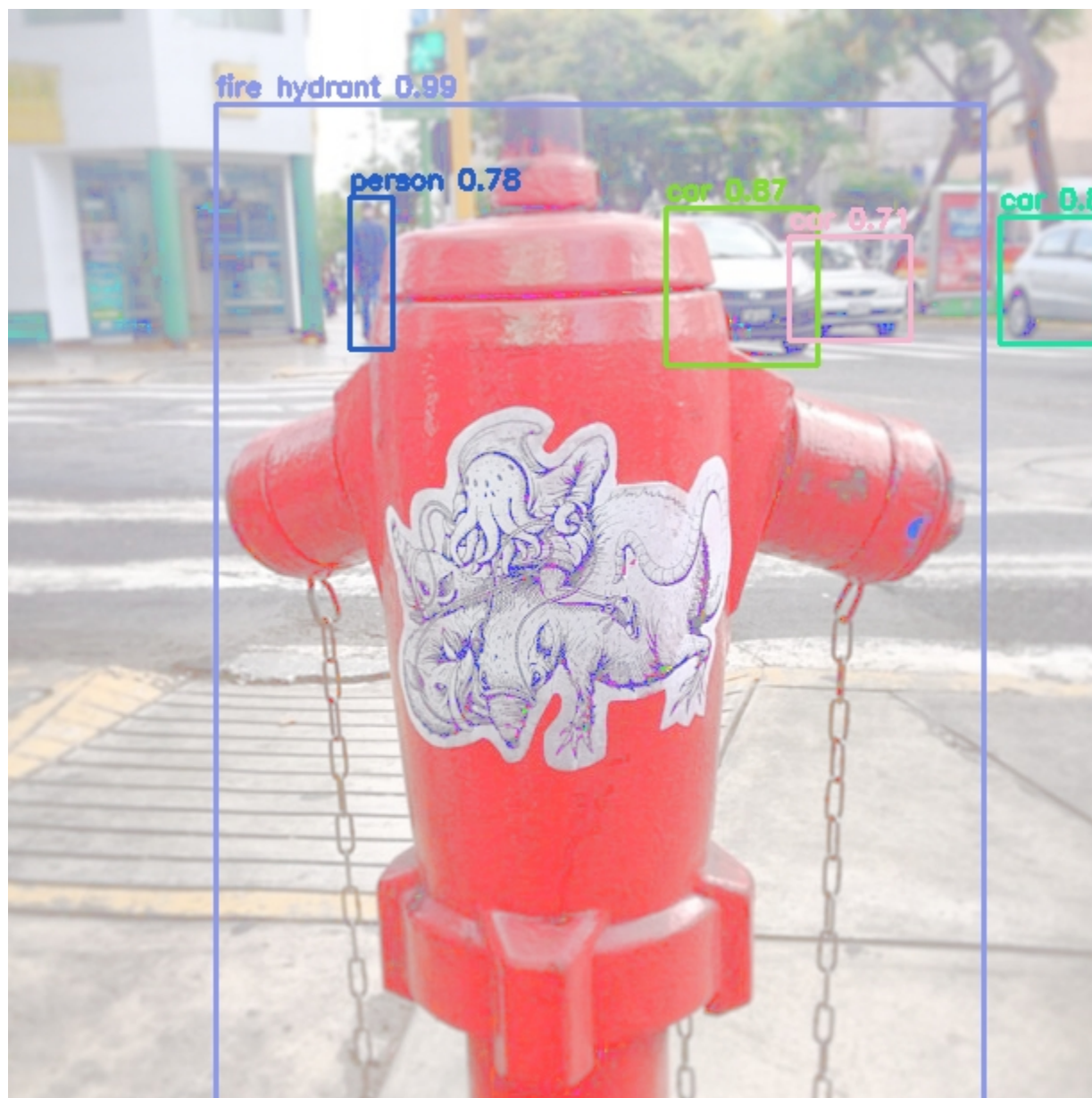
```
[2]: def get_err_root_node():
    err_node = Array()
    err_root_node = Series(err_node)
    err_node.addfilter(Brightness("tar", "rat", "range"))
    return err_root_node
```

4.5.1 Examples from run_yolo_example.py

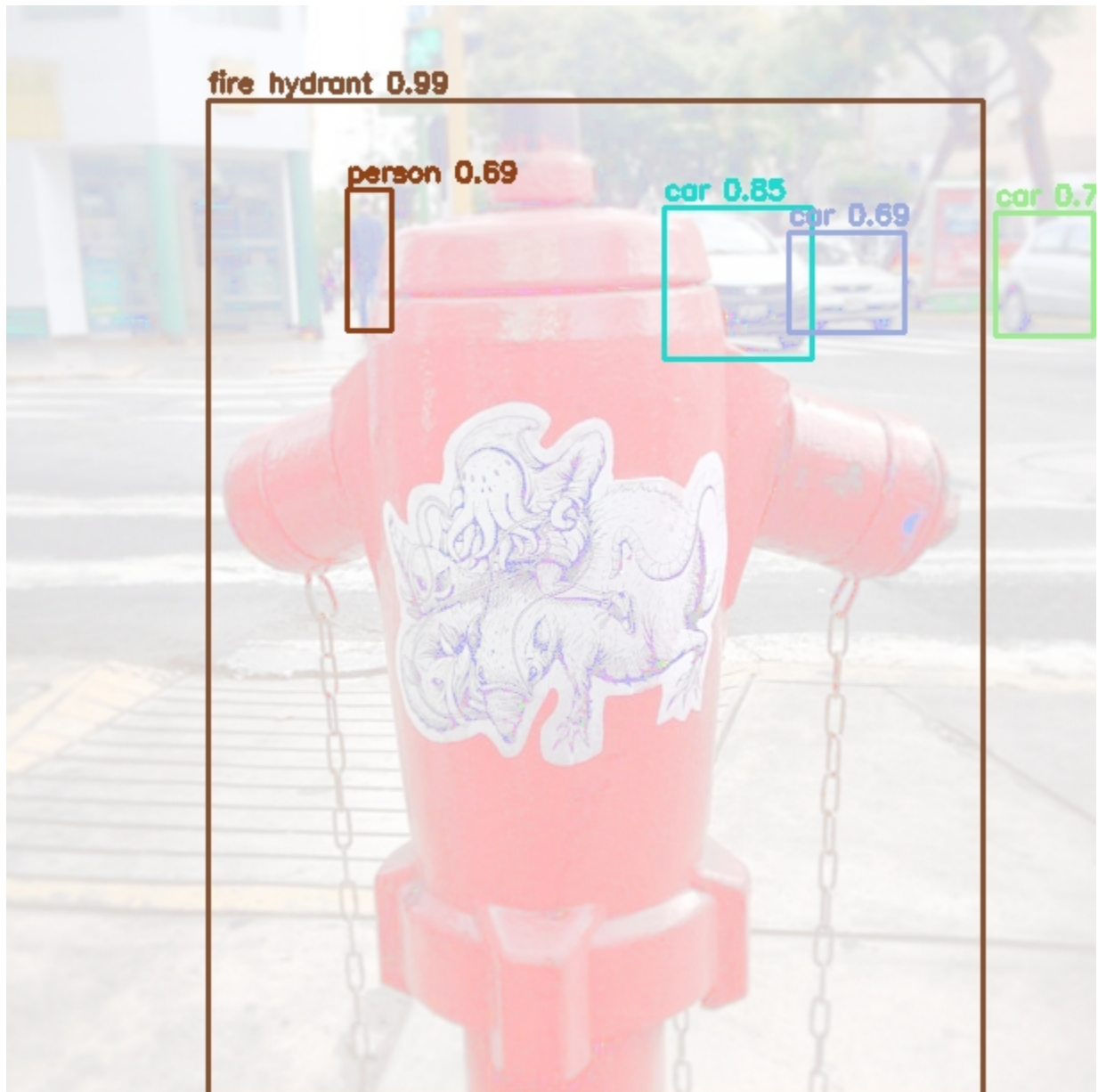
rat: 0.0



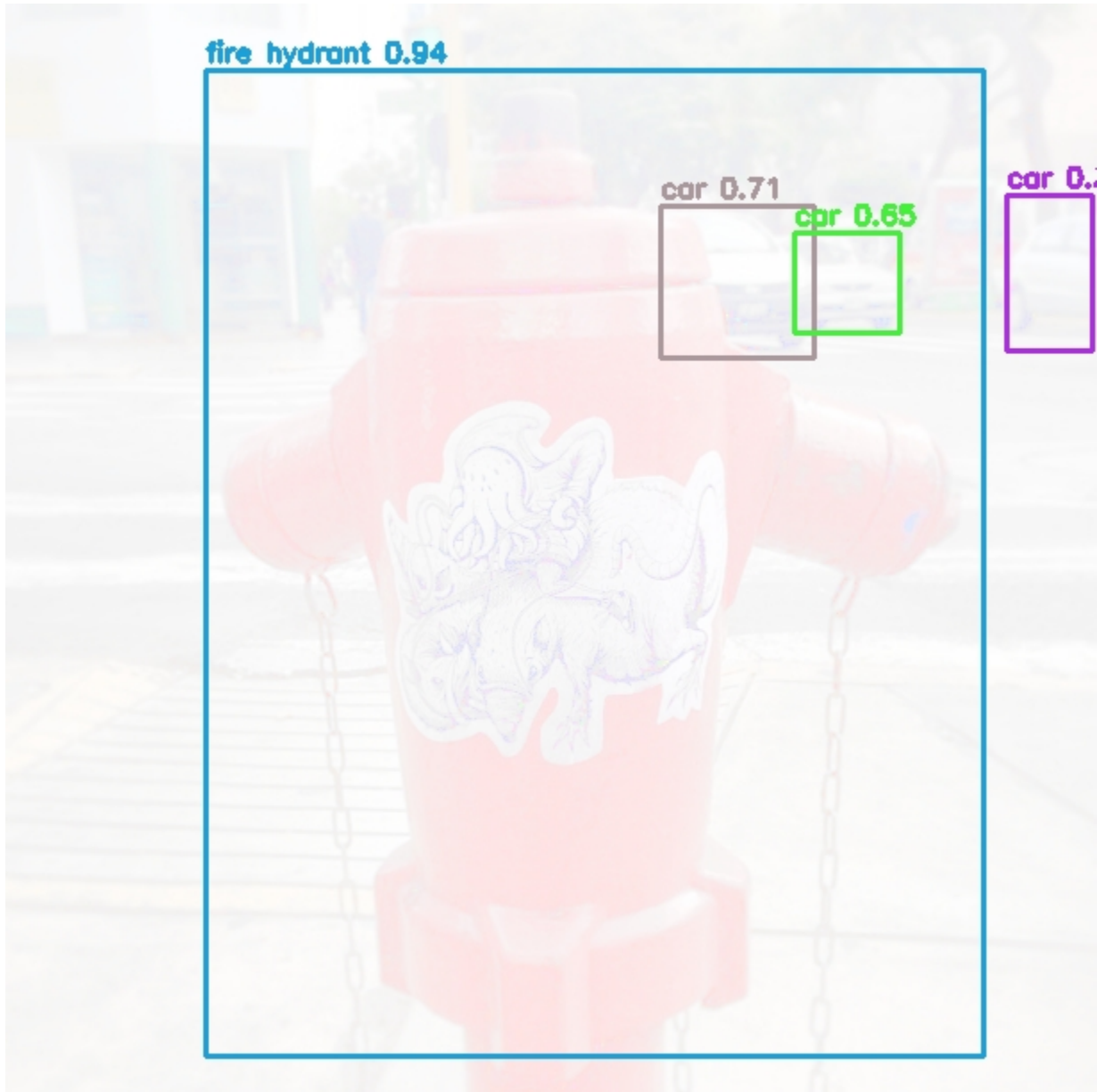
rat: 0.47



rat: 0.93



rat: 1.4



```
[3]: def get_err_params_list():
    rat_steps = np.linspace(0, 1.4, num=8)
    return [{"tar": 1, "rat": rat, "range": 255} for rat in rat_steps]
```

```
[4]: class Preprocessor:

    def run(self, _, imgs, params):
        img_filenames = params["img_filenames"]

        for i, img_arr in enumerate(imgs):
            img = Image.fromarray(img_arr)
            path_to_img = f"{get_project_root()}/tmp/val2017/" + img_filenames[i]
```

(continues on next page)

(continued from previous page)

```

        img.save(path_to_img, "jpeg", quality=100)

    return None, imgs, {}

```

Detectron's model zoo had pretrained weights for FasterRCNN, MaskRCNN and RetinaNet. YOLOv3's weights were trained by us, using the Kale cluster of University of Helsinki. The training took approximately five days when two NVIDIA Tesla V100 GPUs were used.

```

[5]: class YOLOv3Model:

    def run(self, _, imgs, params):
        path_to_yolov3_weights, path_to_yolov3_cfg = load_yolov3()

        cline = f"{get_project_root()}/libs/darknet/darknet detector map {get_project_
→root()}/data/coco.data \
                {path_to_yolov3_cfg} {path_to_yolov3_weights}"
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"\(mAP@0.50\) = (\d+\.\d+)", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

class AbstractDetectronModel(ABC):

    def run(self, _, imgs, params):
        path_to_cfg = self.get_path_to_cfg()
        url_to_weights = self.get_url_to_weights()

        cline = f"""{get_project_root()}/libs/Detectron/tools/test_net.py \
                --cfg {path_to_cfg} \
                TEST.WEIGHTS {url_to_weights} \
                NUM_GPUS 1 \
                TEST.DATASETS '("coco_2017_val",)' \
                MODEL.MASK_ON False \
                OUTPUT_DIR {get_project_root()}/tmp \
                DOWNLOAD_CACHE {get_project_root()}/tmp"""
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"IoU=0.50      \|| area=    all \|| maxDets=100 ] = (\d+\.\d+)
→", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

    @abstractmethod
    def get_path_to_cfg(self):
        pass

    @abstractmethod
    def get_url_to_weights(self):
        pass

class FasterRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
→faster_rcnn_X-101-64x4d-FPN_1x.yaml"

```

(continues on next page)

(continued from previous page)

```

def get_url_to_weights(self):
    return (
        "https://dl.fbaipublicfiles.com/detectron/35858015/12_2017_baselines/"
        "e2e_faster_rcnn_X-101-64x4d-FPN_1x.yaml.01_40_54.1xc565DE/output/train/"
        "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
    )

class MaskRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
↪mask_rcnn_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36494496/12_2017_baselines/"
            "e2e_mask_rcnn_X-101-64x4d-FPN_1x.yaml.07_50_11.fkwVtEvg/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
        )

class RetinaNetModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/
↪retinanet_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36768875/12_2017_baselines/"
            "retinanet_X-101-64x4d-FPN_1x.yaml.08_34_37.FSXgMpzP/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/retinanet/model_final.pkl"
        )

```

```

[6]: def get_model_params_dict_list():
    return [
        {"model": FasterRCNNModel, "params_list": [{}]},
        {"model": MaskRCNNModel, "params_list": [{}]},
        {"model": RetinaNetModel, "params_list": [{}]},
        {"model": YOLOv3Model, "params_list": [{}]},
    ]

```

```

[7]: def visualize(df):
    visualize_scores(
        df,
        score_names=["mAP-50"],
        is_higher_score_better=[True],
        err_param_name="rat",
        title="Object detection with added brightness"
    )
    plt.show()

```

```
[8]: def main():
    imgs, img_filenames = get_data()

    df = runner.run(
        train_data=None,
        test_data=imgs,
        preproc=Preprocessor,
        preproc_params={"img_filenames": img_filenames},
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(),
        n_processes=1
    )

    print_results_by_model(df, dropped_columns=["tar", "range"])
    visualize(df)
```

```
[9]: main()

loading annotations into memory...
Done (t=0.51s)
creating index...
index created!

0%|          | 0/8 [00:00<?, ?it/s]

12%|         | 1/8 [47:55<5:35:27, 2875.41s/it]

25%|        | 2/8 [1:35:32<4:46:59, 2869.98s/it]

38%|       | 3/8 [2:22:59<3:58:35, 2863.07s/it]

50%|      | 4/8 [3:10:11<3:10:15, 2853.80s/it]

62%|     | 5/8 [3:57:10<2:22:10, 2843.37s/it]
```

(continues on next page)

(continued from previous page)

75%| | 6/8 [4:43:53<1:34:22, 2831.07s/it]

88%| | 7/8 [5:30:14<46:56, 2816.10s/it]

100%|| 8/8 [6:16:10<00:00, 2798.09s/it]

FasterRCNN #1

	mAP-50	rat	time_err	time_pre	time_mod
0	0.637	0.0	25.482	191.538	782.473
1	0.568	0.2	24.634	193.666	773.755
2	0.522	0.4	24.803	190.588	771.368
3	0.495	0.6	25.511	183.167	768.999
4	0.453	0.8	24.606	178.801	767.146
5	0.395	1.0	25.253	170.736	762.966
6	0.317	1.2	24.191	163.135	760.870
7	0.224	1.4	24.586	153.263	758.272

MaskRCNN #1

	mAP-50	rat	time_err	time_pre	time_mod
0	0.643	0.0	25.482	191.538	773.686
1	0.575	0.2	24.634	193.666	771.586
2	0.526	0.4	24.803	190.588	768.721
3	0.503	0.6	25.511	183.167	767.139
4	0.466	0.8	24.606	178.801	763.531
5	0.411	1.0	25.253	170.736	761.261
6	0.333	1.2	24.191	163.135	758.198
7	0.244	1.4	24.586	153.263	754.346

RetinaNet #1

	mAP-50	rat	time_err	time_pre	time_mod
0	0.594	0.0	25.482	191.538	1007.041
1	0.529	0.2	24.634	193.666	1006.339
2	0.488	0.4	24.803	190.588	1003.530
3	0.464	0.6	25.511	183.167	999.181
4	0.429	0.8	24.606	178.801	998.645
5	0.378	1.0	25.253	170.736	995.535
6	0.311	1.2	24.191	163.135	988.472
7	0.228	1.4	24.586	153.263	980.072

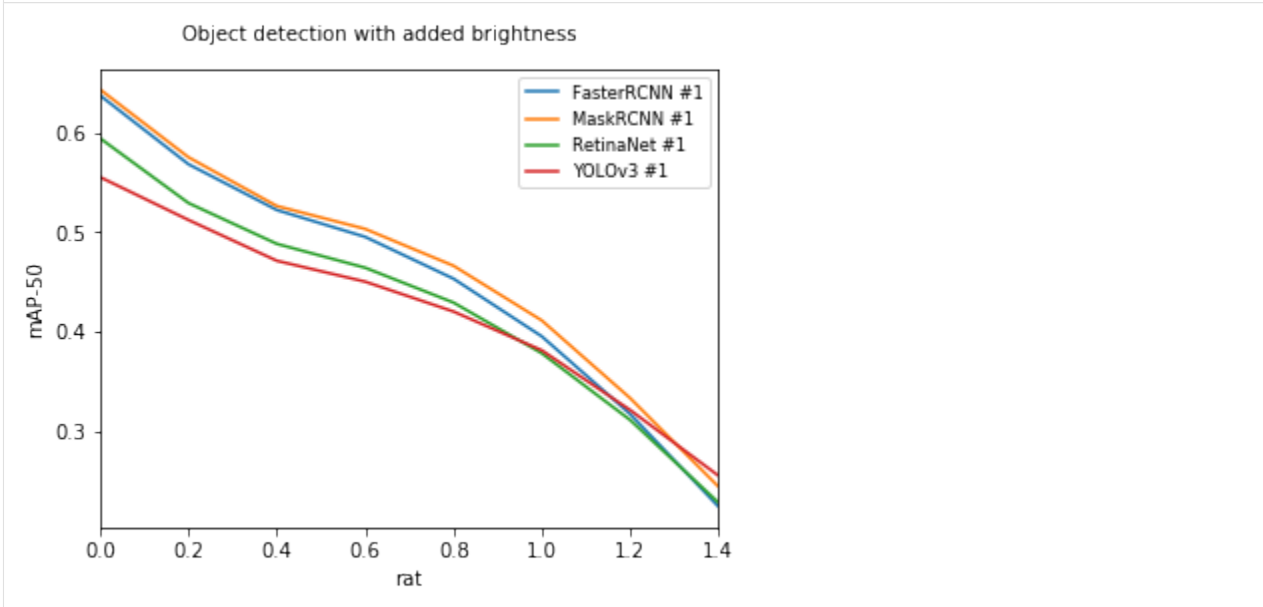
YOLOv3 #1

	mAP-50	rat	time_err	time_pre	time_mod
0	0.555	0.0	25.482	191.538	91.312
1	0.512	0.2	24.634	193.666	82.249
2	0.471	0.4	24.803	190.588	81.664
3	0.450	0.6	25.511	183.167	81.264

(continues on next page)

(continued from previous page)

4	0.420	0.8	24.606	178.801	81.255
5	0.381	1.0	25.253	170.736	81.076
6	0.321	1.2	24.191	163.135	80.890
7	0.255	1.4	24.586	153.263	79.988



The notebook for this case study can be found [here](#).

4.6 Object detection: Added snow

Note: *Object detection requirements.*

Warning: Runtimes can be several hours even on clusters.

We compared the performance of models from FaceBook's Detectron project and YOLOv3 model from Joseph Redmon, when different error sources were added. The models from FaceBook's Detectron project were FasterRCNN, MaskRCNN and RetinaNet.

```
[1]: import re
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
from PIL import Image

from dpemu import runner
from dpemu.dataset_utils import load_coco_val_2017
from dpemu.filters.image import Snow
from dpemu.ml_utils import run_ml_module_using_cli, load_yolov3
from dpemu.nodes import Array, Series
from dpemu.plotting_utils import print_results_by_model, visualize_scores
from dpemu.utils import get_project_root
```

We used 118 287 jpg images (COCO train2017) as the train set and 5000 images (COCO val2017) as the test set to calculate the mAP-50 scores.

```
[2]: def get_data():  
    imgs, _, _, img_filenames = load_coco_val_2017()  
    return imgs, img_filenames
```

```
[3]: def get_err_root_node():  
    err_node = Array()  
    err_root_node = Series(err_node)  
    err_node.addfilter(Snow("snowflake_probability", "snowflake_alpha", "snowstorm_  
↪alpha"))  
    return err_root_node
```

4.6.1 Examples from run_yolo_example.py

snowflake_probability: 0.0001



snowflake_probability: 0.001



snowflake_probability: 0.01



snowflake_probability: 0.1



```
[4]: def get_err_params_list():
      return [{"snowflake_probability": p, "snowflake_alpha": .4, "snowstorm_alpha": 0}
              for p in [10 ** i for i in range(-4, 0)]]
```

```
[5]: class Preprocessor:

      def run(self, _, imgs, params):
          img_filenames = params["img_filenames"]

          for i, img_arr in enumerate(imgs):
              img = Image.fromarray(img_arr)
              path_to_img = f"{get_project_root()}/tmp/val2017/" + img_filenames[i]
```

(continues on next page)

(continued from previous page)

```

img.save(path_to_img, "jpeg", quality=100)

return None, imgs, {}

```

Detectron's model zoo had pretrained weights for FasterRCNN, MaskRCNN and RetinaNet. YOLOv3's weights were trained by us, using the Kale cluster of University of Helsinki. The training took approximately five days when two NVIDIA Tesla V100 GPUs were used.

```

[6]: class YOLOv3Model:

    def run(self, _, imgs, params):
        path_to_yolov3_weights, path_to_yolov3_cfg = load_yolov3()

        cline = f"{get_project_root()}/libs/darknet/darknet detector map {get_project_
→root()}/data/coco.data \
                {path_to_yolov3_cfg} {path_to_yolov3_weights}"
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"\(mAP@0.50\) = (\d+\.\d+)", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

class AbstractDetectronModel(ABC):

    def run(self, _, imgs, params):
        path_to_cfg = self.get_path_to_cfg()
        url_to_weights = self.get_url_to_weights()

        cline = f"""{get_project_root()}/libs/Detectron/tools/test_net.py \
--cfg {path_to_cfg} \
TEST.WEIGHTS {url_to_weights} \
NUM_GPUS 1 \
TEST.DATASETS '("coco_2017_val",)' \
MODEL.MASK_ON False \
OUTPUT_DIR {get_project_root()}/tmp \
DOWNLOAD_CACHE {get_project_root()}/tmp"""
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"IoU=0.50      \|| area=    all \|| maxDets=100 ] = (\d+\.\d+)
→", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

    @abstractmethod
    def get_path_to_cfg(self):
        pass

    @abstractmethod
    def get_url_to_weights(self):
        pass

class FasterRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
→faster_rcnn_X-101-64x4d-FPN_1x.yaml"

```

(continues on next page)

(continued from previous page)

```

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/35858015/12_2017_baselines/"
            "e2e_faster_rcnn_X-101-64x4d-FPN_1x.yaml.01_40_54.1xc565DE/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
        )

class MaskRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
↪mask_rcnn_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36494496/12_2017_baselines/"
            "e2e_mask_rcnn_X-101-64x4d-FPN_1x.yaml.07_50_11.fkwVtEvg/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
        )

class RetinaNetModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/
↪retinanet_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36768875/12_2017_baselines/"
            "retinanet_X-101-64x4d-FPN_1x.yaml.08_34_37.FSXgMpzP/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/retinanet/model_final.pkl"
        )

```

```

[7]: def get_model_params_dict_list():
    return [
        {"model": FasterRCNNModel, "params_list": [{}]},
        {"model": MaskRCNNModel, "params_list": [{}]},
        {"model": RetinaNetModel, "params_list": [{}]},
        {"model": YOLOv3Model, "params_list": [{}]},
    ]

```

```

[8]: def visualize(df):
    visualize_scores(
        df,
        score_names=["mAP-50"],
        is_higher_score_better=[True],
        err_param_name="snowflake_probability",
        title="Object detection with added snow",
        x_log=True
    )
    plt.show()

```

```
[9]: def main():
    imgs, img_filenames = get_data()

    df = runner.run(
        train_data=None,
        test_data=imgs,
        preproc=Preprocessor,
        preproc_params={"img_filenames": img_filenames},
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(),
        n_processes=1
    )

    print_results_by_model(df, dropped_columns=["snowflake_alpha", "snowstorm_alpha"])
    visualize(df)
```

```
[10]: main()
```

```
loading annotations into memory...
```

```
Done (t=0.47s)
```

```
creating index...
```

```
index created!
```

```
0%|          | 0/4 [00:00<?, ?it/s]
```

```
25%|         | 1/4 [54:40<2:44:01, 3280.50s/it]
```

```
50%|        | 2/4 [1:49:43<1:49:34, 3287.14s/it]
```

```
75%|       | 3/4 [2:50:25<56:33, 3393.72s/it]
```

```
100%|| 4/4 [4:42:35<00:00, 4394.49s/it]
```

```
FasterRCNN #1
```

	mAP-50	snowflake_probability	time_err	time_pre	time_mod
0	0.633	0.0001	324.606	201.185	822.460
1	0.607	0.0010	354.531	200.952	817.257
2	0.472	0.0100	683.656	213.076	818.777
3	0.134	0.1000	3803.646	215.833	809.986

MaskRCNN #1					
	mAP-50	snowflake_probability	time_err	time_pre	time_mod
0	0.638	0.0001	324.606	201.185	815.829
1	0.615	0.0010	354.531	200.952	815.793
2	0.481	0.0100	683.656	213.076	815.570
3	0.139	0.1000	3803.646	215.833	803.844
RetinaNet #1					
	mAP-50	snowflake_probability	time_err	time_pre	time_mod
0	0.591	0.0001	324.606	201.185	1011.160
1	0.565	0.0010	354.531	200.952	1015.663
2	0.434	0.0100	683.656	213.076	1012.165
3	0.138	0.1000	3803.646	215.833	993.220
YOLOv3 #1					
	mAP-50	snowflake_probability	time_err	time_pre	time_mod
0	0.552	0.0001	324.606	201.185	101.363
1	0.519	0.0010	354.531	200.952	93.684
2	0.385	0.0100	683.656	213.076	92.153
3	0.088	0.1000	3803.646	215.833	96.144

snowflake_probability	FasterRCNN #1	MaskRCNN #1	RetinaNet #1	YOLOv3 #1
10 ⁻⁴	0.638	0.638	0.591	0.552
10 ⁻³	0.615	0.615	0.565	0.519
10 ⁻²	0.481	0.481	0.434	0.385
10 ⁻¹	0.139	0.138	0.138	0.088

The notebook for this case study can be found [here](#).

4.7 Object detection: JPEG compression

Note: *Object detection requirements.*

Warning: Runtimes can be several hours even on clusters.

We compared the performance of models from FaceBook's Detectron project and YOLOv3 model from Joseph Redmon, when different error sources were added. The models from FaceBook's Detectron project were FasterRCNN,

MaskRCNN and RetinaNet.

```
[1]: import re
      from abc import ABC, abstractmethod

      import matplotlib.pyplot as plt
      from PIL import Image

      from dpemu import runner
      from dpemu.dataset_utils import load_coco_val_2017
      from dpemu.filters.image import JPEG_Compression
      from dpemu.ml_utils import run_ml_module_using_cli, load_yolov3
      from dpemu.nodes import Array, Series
      from dpemu.plotting_utils import print_results_by_model, visualize_scores
      from dpemu.utils import get_project_root
```

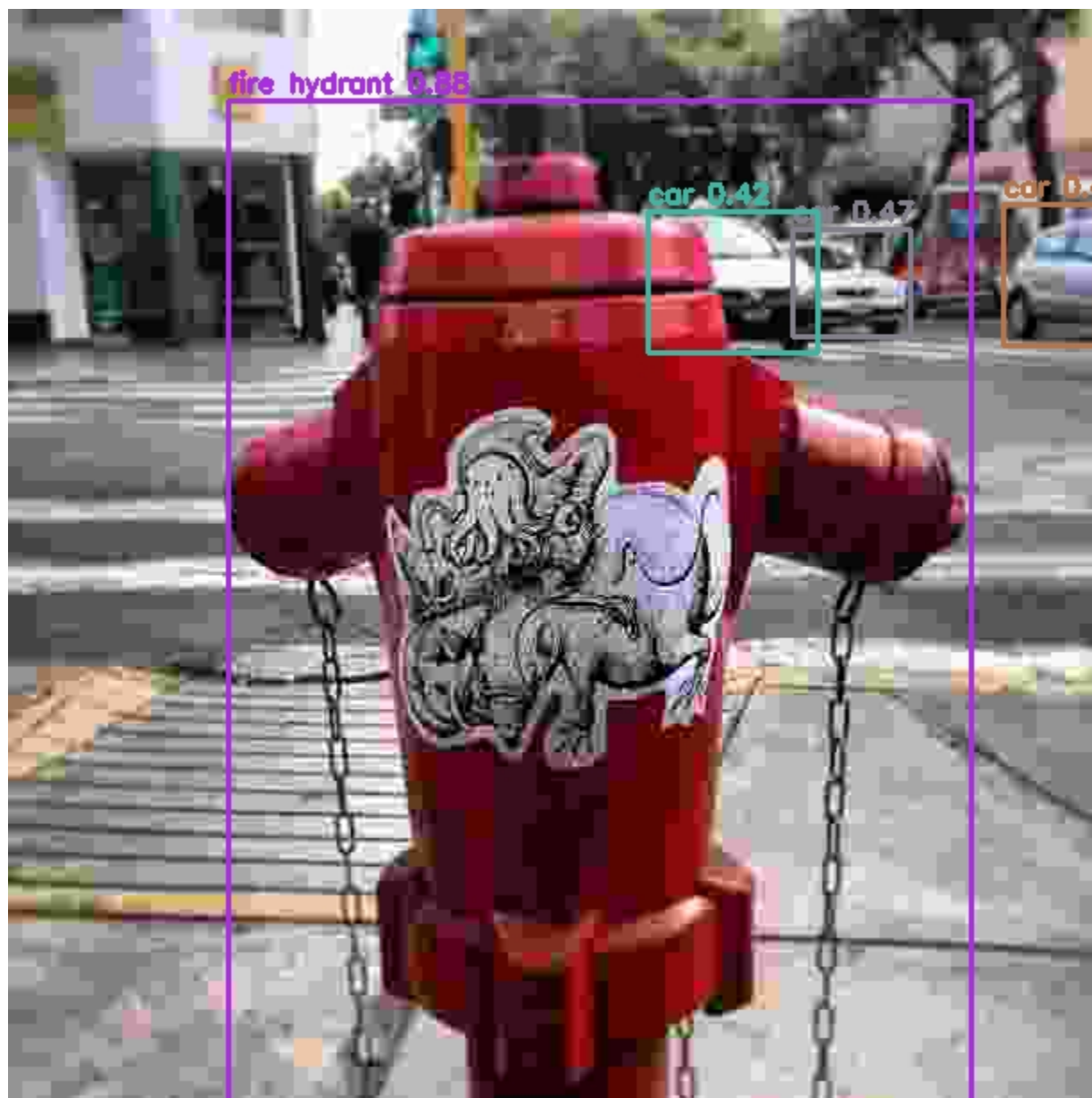
We used 118 287 jpg images (COCO train2017) as the train set and 5000 images (COCO val2017) as the test set to calculate the mAP-50 scores.

```
[ ]: def get_data():
      imgs, _, _, img_filenames = load_coco_val_2017()
      return imgs, img_filenames
```

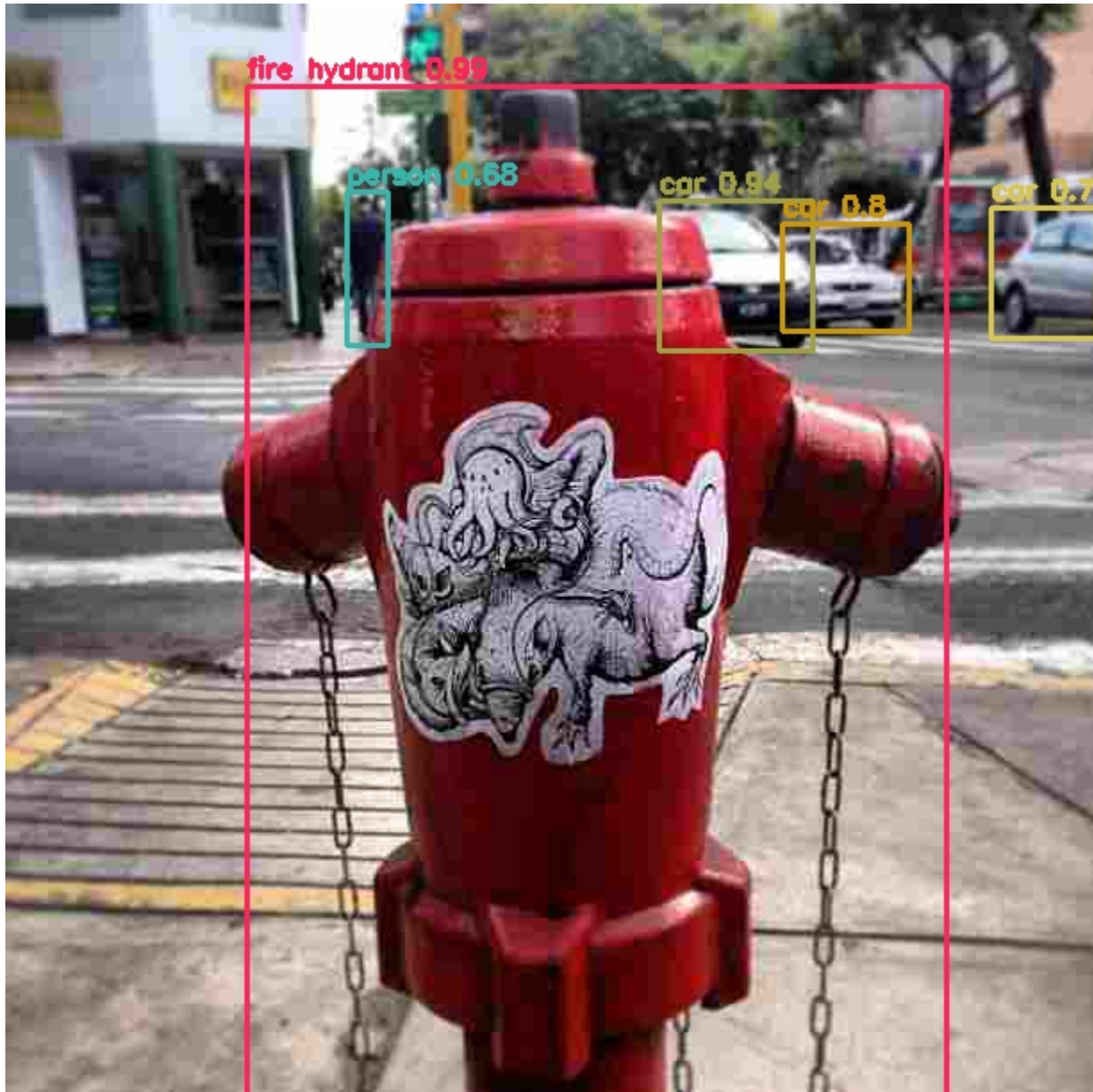
```
[2]: def get_err_root_node():
      err_node = Array()
      err_root_node = Series(err_node)
      err_node.addfilter(JPEG_Compression("quality"))
      return err_root_node
```

4.7.1 Examples from run_yolo_example.py

quality: 5



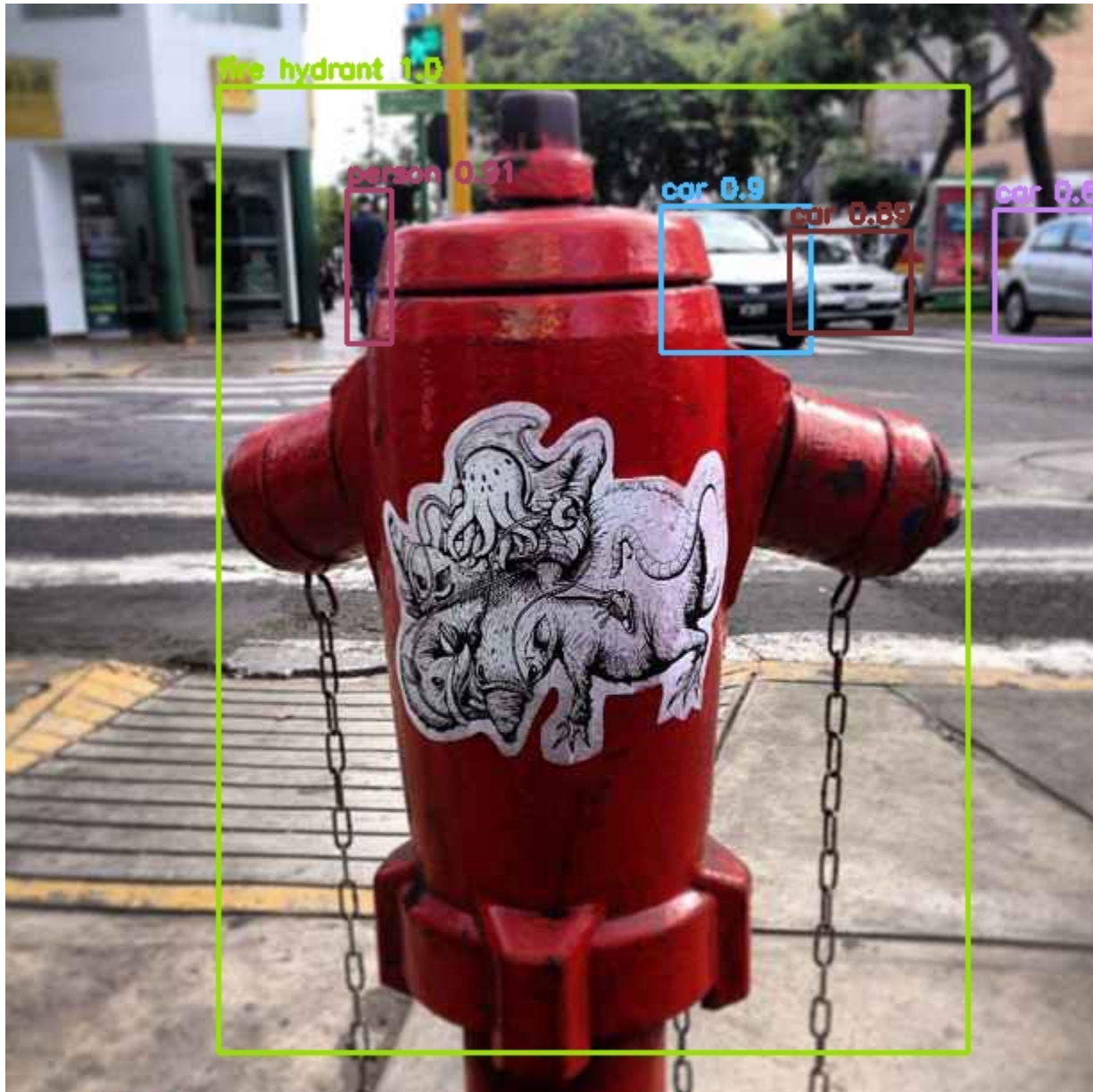
quality: 15



quality: 25



quality: 40



```
[3]: def get_err_params_list():
      return [{"quality": q} for q in range(5, 45, 5)]
```

```
[4]: class Preprocessor:

      def run(self, _, imgs, params):
          img_filenames = params["img_filenames"]

          for i, img_arr in enumerate(imgs):
              img = Image.fromarray(img_arr)
              path_to_img = f"{get_project_root()}/tmp/val2017/" + img_filenames[i]
              img.save(path_to_img, "jpeg", quality=100)
```

(continues on next page)

(continued from previous page)

```
return None, imgs, {}
```

Detectron's model zoo had pretrained weights for FasterRCNN, MaskRCNN and RetinaNet. YOLOv3's weights were trained by us, using the Kale cluster of University of Helsinki. The training took approximately five days when two NVIDIA Tesla V100 GPUs were used.

```
[5]: class YOLOv3Model:

    def run(self, _, imgs, params):
        path_to_yolov3_weights, path_to_yolov3_cfg = load_yolov3()

        cline = f"{get_project_root()}/libs/darknet/darknet detector map {get_project_
→root()}/data/coco.data \
                {path_to_yolov3_cfg} {path_to_yolov3_weights}"
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"\(mAP@0.50\) = (\d+\.\d+)", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

class AbstractDetectronModel(ABC):

    def run(self, _, imgs, params):
        path_to_cfg = self.get_path_to_cfg()
        url_to_weights = self.get_url_to_weights()

        cline = f"""{get_project_root()}/libs/Detectron/tools/test_net.py \
                --cfg {path_to_cfg} \
                TEST.WEIGHTS {url_to_weights} \
                NUM_GPUS 1 \
                TEST.DATASETS '("coco_2017_val",)' \
                MODEL.MASK_ON False \
                OUTPUT_DIR {get_project_root()}/tmp \
                DOWNLOAD_CACHE {get_project_root()}/tmp"""
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"IoU=0.50      \|| area=    all \|| maxDets=100 ] = (\d+\.\d+)
→", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

    @abstractmethod
    def get_path_to_cfg(self):
        pass

    @abstractmethod
    def get_url_to_weights(self):
        pass

class FasterRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
→faster_rcnn_X-101-64x4d-FPN_1x.yaml"
```

(continues on next page)

(continued from previous page)

```

def get_url_to_weights(self):
    return (
        "https://dl.fbaipublicfiles.com/detectron/35858015/12_2017_baselines/"
        "e2e_faster_rcnn_X-101-64x4d-FPN_1x.yaml.01_40_54.1xc565DE/output/train/"
        "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
    )

class MaskRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
↪mask_rcnn_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36494496/12_2017_baselines/"
            "e2e_mask_rcnn_X-101-64x4d-FPN_1x.yaml.07_50_11.fkwVtEvg/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
        )

class RetinaNetModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/
↪retinanet_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36768875/12_2017_baselines/"
            "retinanet_X-101-64x4d-FPN_1x.yaml.08_34_37.FSXgMpzP/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/retinanet/model_final.pkl"
        )

```

```

[6]: def get_model_params_dict_list():
    return [
        {"model": FasterRCNNModel, "params_list": [{}]},
        {"model": MaskRCNNModel, "params_list": [{}]},
        {"model": RetinaNetModel, "params_list": [{}]},
        {"model": YOLOv3Model, "params_list": [{}]},
    ]

```

```

[7]: def visualize(df):
    visualize_scores(
        df,
        score_names=["mAP-50"],
        is_higher_score_better=[True],
        err_param_name="quality",
        title="Object detection with JPEG compression"
    )
    plt.show()

```

```
[8]: def main():
    imgs, img_filenames = get_data()

    df = runner.run(
        train_data=None,
        test_data=imgs,
        preproc=Preprocessor,
        preproc_params={"img_filenames": img_filenames},
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(),
        n_processes=1
    )

    print_results_by_model(df)
    visualize(df)
```

```
[9]: main()

loading annotations into memory...
Done (t=0.54s)
creating index...
index created!

0%|          | 0/8 [00:00<?, ?it/s]

12%|         | 1/8 [48:22<5:38:35, 2902.14s/it]

25%|        | 2/8 [1:37:02<4:50:44, 2907.49s/it]

38%|       | 3/8 [2:26:01<4:03:06, 2917.21s/it]

50%|      | 4/8 [3:15:03<3:14:58, 2924.60s/it]

62%|     | 5/8 [4:04:14<2:26:37, 2932.38s/it]
```

(continues on next page)

(continued from previous page)

75%| | 6/8 [4:53:27<1:37:56, 2938.49s/it]

88%| | 7/8 [5:42:35<49:01, 2941.60s/it]

100%|| 8/8 [6:31:51<00:00, 2945.84s/it]

FasterRCNN #1

	mAP-50	quality	time_err	time_pre	time_mod
0	0.092	5	64.709	122.873	808.970
1	0.277	10	66.442	127.487	809.532
2	0.407	15	68.678	134.865	811.471
3	0.479	20	70.196	137.164	813.800
4	0.523	25	70.890	139.103	814.348
5	0.550	30	72.191	142.072	812.131
6	0.566	35	73.478	144.595	812.938
7	0.576	40	74.204	145.890	813.014

MaskRCNN #1

	mAP-50	quality	time_err	time_pre	time_mod
0	0.105	5	64.709	122.873	796.518
1	0.301	10	66.442	127.487	809.272
2	0.431	15	68.678	134.865	810.367
3	0.498	20	70.196	137.164	808.612
4	0.539	25	70.890	139.103	811.646
5	0.563	30	72.191	142.072	811.534
6	0.576	35	73.478	144.595	810.514
7	0.587	40	74.204	145.890	812.432

RetinaNet #1

	mAP-50	quality	time_err	time_pre	time_mod
0	0.122	5	64.709	122.873	1007.092
1	0.305	10	66.442	127.487	1009.791
2	0.412	15	68.678	134.865	1014.549
3	0.466	20	70.196	137.164	1011.720
4	0.500	25	70.890	139.103	1015.210
5	0.521	30	72.191	142.072	1014.208
6	0.535	35	73.478	144.595	1009.012
7	0.543	40	74.204	145.890	1010.712

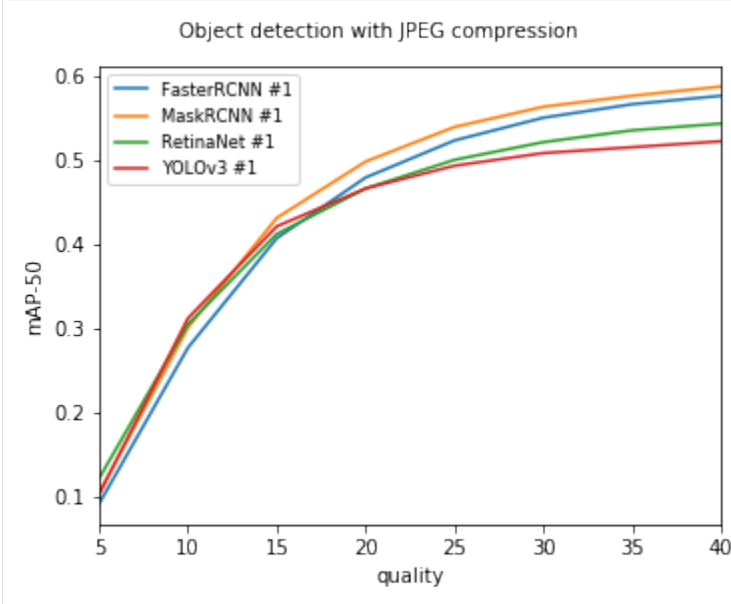
YOLOv3 #1

	mAP-50	quality	time_err	time_pre	time_mod
0	0.103	5	64.709	122.873	98.066
1	0.312	10	66.442	127.487	92.068
2	0.421	15	68.678	134.865	92.472
3	0.466	20	70.196	137.164	92.607

(continues on next page)

(continued from previous page)

4	0.493	25	70.890	139.103	92.315
5	0.508	30	72.191	142.072	92.577
6	0.515	35	73.478	144.595	92.671
7	0.522	40	74.204	145.890	92.401



The notebook for this case study can be found [here](#).

4.8 Object detection: Reduced resolution

Note: *Object detection requirements.*

Warning: Runtimes can be several hours even on clusters.

We compared the performance of models from FaceBook’s Detectron project and YOLOv3 model from Joseph Redmon, when different error sources were added. The models from FaceBook’s Detectron project were FasterRCNN, MaskRCNN and RetinaNet.

```
[1]: import re
from abc import ABC, abstractmethod

import matplotlib.pyplot as plt
from PIL import Image

from dpemu import runner
from dpemu.dataset_utils import load_coco_val_2017
from dpemu.filters.image import Resolution
from dpemu.ml_utils import run_ml_module_using_cli, load_yolov3
from dpemu.nodes import Array, Series
from dpemu.plotting_utils import print_results_by_model, visualize_scores
from dpemu.utils import get_project_root
```

We used 118 287 jpg images (COCO train2017) as the train set and 5000 images (COCO val2017) as the test set to calculate the mAP-50 scores.

```
[2]: def get_data():  
      imgs, _, _, img_filenames = load_coco_val_2017()  
      return imgs, img_filenames
```

```
[3]: def get_err_root_node():  
      err_node = Array()  
      err_root_node = Series(err_node)  
      err_node.addfilter(Resolution("k"))  
      return err_root_node
```

4.8.1 Examples from run_yolo_example.py

k: 1



k: 2



k: 3



k: 4



```
[4]: def get_err_params_list():  
      return [{"k": k} for k in range(1, 5)]
```

```
[5]: class Preprocessor:  
  
      def run(self, _, imgs, params):  
          img_filenames = params["img_filenames"]  
  
          for i, img_arr in enumerate(imgs):  
              img = Image.fromarray(img_arr)  
              path_to_img = f"{get_project_root()}/tmp/val2017/" + img_filenames[i]  
              img.save(path_to_img, "jpeg", quality=100)
```

(continues on next page)

(continued from previous page)

```
return None, imgs, {}
```

Detectron's model zoo had pretrained weights for FasterRCNN, MaskRCNN and RetinaNet. YOLOv3's weights were trained by us, using the Kale cluster of University of Helsinki. The training took approximately five days when two NVIDIA Tesla V100 GPUs were used.

```
[6]: class YOLOv3Model:

    def run(self, _, imgs, params):
        path_to_yolov3_weights, path_to_yolov3_cfg = load_yolov3()

        cline = f"{get_project_root()}/libs/darknet/darknet detector map {get_project_
→root()}/data/coco.data \
                {path_to_yolov3_cfg} {path_to_yolov3_weights}"
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"\(mAP@0.50\) = (\d+\.\d+)", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

class AbstractDetectronModel(ABC):

    def run(self, _, imgs, params):
        path_to_cfg = self.get_path_to_cfg()
        url_to_weights = self.get_url_to_weights()

        cline = f"""{get_project_root()}/libs/Detectron/tools/test_net.py \
                --cfg {path_to_cfg} \
                TEST.WEIGHTS {url_to_weights} \
                NUM_GPUS 1 \
                TEST.DATASETS '("coco_2017_val",)' \
                MODEL.MASK_ON False \
                OUTPUT_DIR {get_project_root()}/tmp \
                DOWNLOAD_CACHE {get_project_root()}/tmp"""
        out = run_ml_module_using_cli(cline, show_stdout=False)

        match = re.search(r"IoU=0.50      \|| area=    all \|| maxDets=100 ] = (\d+\.\d+)
→", out)
        return {"mAP-50": round(float(match.group(1)), 3)}

    @abstractmethod
    def get_path_to_cfg(self):
        pass

    @abstractmethod
    def get_url_to_weights(self):
        pass

class FasterRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
→faster_rcnn_X-101-64x4d-FPN_1x.yaml"
```

(continues on next page)

(continued from previous page)

```

def get_url_to_weights(self):
    return (
        "https://dl.fbaipublicfiles.com/detectron/35858015/12_2017_baselines/"
        "e2e_faster_rcnn_X-101-64x4d-FPN_1x.yaml.01_40_54.1xc565DE/output/train/"
        "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
    )

class MaskRCNNModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/e2e_
↪mask_rcnn_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36494496/12_2017_baselines/"
            "e2e_mask_rcnn_X-101-64x4d-FPN_1x.yaml.07_50_11.fkwVtEvg/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/generalized_rcnn/model_final.
↪pk1"
        )

class RetinaNetModel(AbstractDetectronModel):

    def get_path_to_cfg(self):
        return f"{get_project_root()}/libs/Detectron/configs/12_2017_baselines/
↪retinanet_X-101-64x4d-FPN_1x.yaml"

    def get_url_to_weights(self):
        return (
            "https://dl.fbaipublicfiles.com/detectron/36768875/12_2017_baselines/"
            "retinanet_X-101-64x4d-FPN_1x.yaml.08_34_37.FSXgMpzP/output/train/"
            "coco_2014_train%3Acoco_2014_valminusminival/retinanet/model_final.pkl"
        )

```

```

[7]: def get_model_params_dict_list():
    return [
        {"model": FasterRCNNModel, "params_list": [{}]},
        {"model": MaskRCNNModel, "params_list": [{}]},
        {"model": RetinaNetModel, "params_list": [{}]},
        {"model": YOLOv3Model, "params_list": [{}]},
    ]

```

```

[8]: def visualize(df):
    visualize_scores(
        df,
        score_names=["mAP-50"],
        is_higher_score_better=[True],
        err_param_name="k",
        title="Object detection with reduced resolution"
    )
    plt.show()

```



```
[9]: def main():
    imgs, img_filenames = get_data()

    df = runner.run(
        train_data=None,
        test_data=imgs,
        preproc=Preprocessor,
        preproc_params={"img_filenames": img_filenames},
        err_root_node=get_err_root_node(),
        err_params_list=get_err_params_list(),
        model_params_dict_list=get_model_params_dict_list(),
        n_processes=1
    )

    print_results_by_model(df)
    visualize(df)
```

```
[10]: main()

loading annotations into memory...
Done (t=0.52s)
creating index...
index created!

0%|          | 0/4 [00:00<?, ?it/s]

25%|          | 1/4 [48:04<2:24:14, 2884.76s/it]

50%|          | 2/4 [1:35:50<1:35:58, 2879.18s/it]

75%|          | 3/4 [2:23:16<47:48, 2868.95s/it]

100%|| 4/4 [3:09:35<00:00, 2842.14s/it]

FasterRCNN #1
```

	mAP-50	k	time_err	time_pre	time_mod
0	0.636	1	43.079	188.690	778.357
1	0.448	2	43.463	189.020	773.277
2	0.192	3	43.073	197.176	763.480
3	0.070	4	43.302	151.681	758.630

MaskRCNN #1

	mAP-50	k	time_err	time_pre	time_mod
0	0.643	1	43.079	188.690	773.331
1	0.448	2	43.463	189.020	771.375
2	0.188	3	43.073	197.176	759.624
3	0.084	4	43.302	151.681	758.115

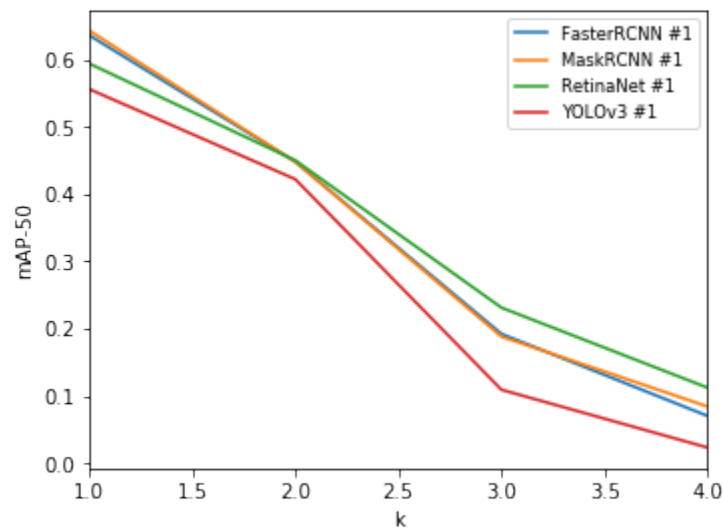
RetinaNet #1

	mAP-50	k	time_err	time_pre	time_mod
0	0.594	1	43.079	188.690	1007.471
1	0.450	2	43.463	189.020	1003.882
2	0.231	3	43.073	197.176	996.645
3	0.112	4	43.302	151.681	982.990

YOLOv3 #1

	mAP-50	k	time_err	time_pre	time_mod
0	0.556	1	43.079	188.690	89.957
1	0.422	2	43.463	189.020	81.409
2	0.109	3	43.073	197.176	79.648
3	0.023	4	43.302	151.681	79.119

Object detection with reduced resolution



The notebook for this case study can be found [here](#).

4.9 Spoken commands example

This example uses an audio classifier model from a Tensorflow tutorial: https://www.tensorflow.org/tutorials/sequences/audio_recognition

N.B. This script downloads a large (2.3GB) speech commands dataset!

```
[1]: import sys
      sys.path.append('.')
      from pathlib import Path
      import tarfile
```

(continues on next page)

(continued from previous page)

```

import shutil
import pandas as pd
from scipy.io.wavfile import read, write
from sklearn.metrics import confusion_matrix
from dpemu.nodes.series import Series
from dpemu.nodes.tuple import Tuple
from dpemu.filters.sound import ClipWAV
from dpemu.filters.common import ApplyToTuple
from dpemu.plotting_utils import visualize_confusion_matrix

```

First we download the dataset unless it is already present. If you have downloaded and extracted the dataset into a different directory, change the `data_dir` variable accordingly.

```

[2]: data_url = "https://storage.googleapis.com/download.tensorflow.org/data/speech_
      ↪ commands_v0.02.tar.gz"
      fname = "speech_commands_v0.02.tar.gz"
      data_dir = Path.home() / "datasets/speech_data"

      if not data_dir.exists():
          data_dir.mkdir(parents=True)
          !wget {data_url} -P {data_dir}
          tarfile.open(data_dir / fname, "r:gz").extractall(data_dir)

```

```

[3]: trained_categories = ["yes", "no", "up", "down", "left", "right", "on", "off", "stop",
      ↪ "go"]
      labels = ["_silence_", "_unknown_", "yes", "no", "up", "down", "left", "right", "on",
      ↪ "off", "stop", "go"]

      test_set_rel_paths = !cat {data_dir / "testing_list.txt"}
      test_set_files = [data_dir / p for p in test_set_rel_paths]
      test_categories = !cut -d '/' -f1 {data_dir / "testing_list.txt"} | sort -u

      len(test_set_files), len(test_categories), len(trained_categories)

```

```

[3]: (11005, 35, 10)

```

In order to download the speech commands dataset to the correct place, we need to set the variables `dpemu_path` and `example_path`.

```

[4]: dpemu_path = Path.cwd().parents[1]
      example_path = dpemu_path / "examples/speech_commands"

```

Choose a category in which to generate errors. Later on we will generate errors in all of the test set categories.

```

[5]: category = "stop"
      data_subset_dir = data_dir / category

      fs = list(data_subset_dir.iterdir())
      wavs = [read(f) for f in data_subset_dir.iterdir()]

```

Create an error generating tree and generate errors in the category chosen above.

```

[6]: wav_node = Tuple()
      wav_node.addfilter(ApplyToTuple(ClipWAV("dyn_range"), 1))
      root_node = Series(wav_node)

```

(continues on next page)

(continued from previous page)

```
err_params = {"dyn_range": .2}
clipped = root_node.generate_error(wavs, err_params)
```

Now we arbitrarily choose a speech command example from the data subset. To try another audio clip, change the index.

```
[7]: example_index = 123
```

```
[8]: clipped_filename = data_dir / 'clipped.wav'
write(clipped_filename, 16000, clipped[example_index][1])
```

```
[9]: !aplay {fs[example_index]}

Playing WAVE '/home/jpssilve/datasets/speech_data/stop/3ec05c3d_nohash_0.wav' :
↳ Signed 16 bit Little Endian, Rate 16000 Hz, Mono
```

```
[10]: !aplay {clipped_filename}

Playing WAVE '/home/jpssilve/datasets/speech_data/clipped.wav' : Signed 16 bit Little
↳ Endian, Rate 16000 Hz, Mono
```

Define a function to filter out irrelevant output (e.g. Python deprecation warnings):

```
[11]: def filter_scores(output):
        return [line for line in output if "score" in line or ".wav" in line]
```

Run the model on the clean clip selected above:

```
[12]: scores_clean = !python {example_path}/label_wav.py \
--graph={example_path}/trained_model/my_frozen_graph.pb \
--labels={example_path}/trained_model/conv_labels.txt \
--wav={fs[example_index]}

filter_scores(scores_clean)

[12]: ['stop (score = 0.54378)',
      'off (score = 0.19993)',
      '_unknown_ (score = 0.07233)']
```

Run the model on the corresponding errorified clip:

```
[13]: scores_clipped = !python {example_path}/label_wav.py \
--graph={example_path}/trained_model/my_frozen_graph.pb \
--labels={example_path}/trained_model/conv_labels.txt \
--wav={clipped_filename}

filter_scores(scores_clipped)

[13]: ['stop (score = 0.22963)',
      'down (score = 0.16858)',
      '_unknown_ (score = 0.11415)']
```

You can also run the model on an entire directory of .wav files in one go:

```
[14]: scores_clean_dir = !python {example_path}/label_wav_dir.py \
--graph={example_path}/trained_model/my_frozen_graph.pb \
--labels={example_path}/trained_model/conv_labels.txt \
```

(continues on next page)

(continued from previous page)

```

--wav_dir={data_subset_dir}

filter_scores(scores_clean_dir)

[14]: ['0f46028a_nohash_4.wav',
      'stop (score = 0.84888)',
      'up (score = 0.10150)',
      '_unknown_ (score = 0.02897)',
      '095847e4_nohash_0.wav',
      'stop (score = 0.83839)',
      'up (score = 0.10791)',
      'down (score = 0.01377)',
      'f8ba7c0e_nohash_1.wav',
      'stop (score = 0.99616)',
      'down (score = 0.00215)',
      '_unknown_ (score = 0.00114)',
      '4cee0c60_nohash_1.wav',
      'stop (score = 0.94652)',
      'up (score = 0.04828)',
      '_unknown_ (score = 0.00210)',
      '52e228e9_nohash_1.wav',
      'stop (score = 0.98153)',
      'down (score = 0.00989)',
      'up (score = 0.00290)',
      '42f81601_nohash_0.wav',
      'stop (score = 0.95047)',
      'up (score = 0.02973)',
      '_unknown_ (score = 0.01149)',
      'bc065a17_nohash_1.wav',
      'stop (score = 0.51887)',
      'down (score = 0.33725)',
      '_unknown_ (score = 0.13488)',
      '692a88e6_nohash_1.wav',
      'stop (score = 0.83974)',
      'up (score = 0.15001)',
      '_unknown_ (score = 0.00472)',
      '96a48d28_nohash_0.wav',
      'stop (score = 0.99714)',
      '_unknown_ (score = 0.00153)',
      'up (score = 0.00119)',
      '763188c4_nohash_0.wav',
      'stop (score = 0.92912)',
      '_unknown_ (score = 0.04616)',
      'go (score = 0.02025)',
      '53fd1780_nohash_0.wav',
      'stop (score = 0.26836)',
      'down (score = 0.19330)',
      '_unknown_ (score = 0.13893)',
      'e9323bd9_nohash_0.wav',
      'stop (score = 0.69810)',
      'up (score = 0.15704)',
      '_unknown_ (score = 0.06251)',
      '686d030b_nohash_4.wav',
      'stop (score = 0.99679)',
      'up (score = 0.00229)',
      'down (score = 0.00053)',
      'fc3ba625_nohash_0.wav',

```

(continues on next page)

(continued from previous page)

```

'stop (score = 0.94688)',
'_unknown_ (score = 0.02761)',
'up (score = 0.02019)',
'c4e00ee9_nohash_1.wav',
'stop (score = 0.61116)',
'up (score = 0.25302)',
'off (score = 0.04122)',
'66774579_nohash_0.wav',
'stop (score = 0.71052)',
'off (score = 0.07688)',
'up (score = 0.07429)',
'ee07dcb9_nohash_0.wav',
'stop (score = 0.90859)',
'up (score = 0.03769)',
'_unknown_ (score = 0.01578)',
'4634529e_nohash_1.wav',
'stop (score = 0.29116)',
'up (score = 0.19151)',
'off (score = 0.13839)',
'8f3f252c_nohash_0.wav',
'stop (score = 0.86171)',
'up (score = 0.10564)',
'off (score = 0.01070)',
'3e31dffe_nohash_4.wav',
'stop (score = 0.96517)',
'_unknown_ (score = 0.03305)',
'up (score = 0.00064)',
'3d794813_nohash_4.wav',
'stop (score = 0.91730)',
'up (score = 0.05711)',
'_unknown_ (score = 0.02363)',
'f15a354c_nohash_0.wav',
'stop (score = 0.99641)',
'up (score = 0.00274)',
'off (score = 0.00026)',
'c71e3acc_nohash_0.wav',
'stop (score = 0.69207)',
'up (score = 0.10310)',
'go (score = 0.07329)',
'004ae714_nohash_0.wav',
'stop (score = 0.90009)',
'off (score = 0.02894)',
'_unknown_ (score = 0.02361)',
'a16013b7_nohash_4.wav',
'stop (score = 0.57909)',
'up (score = 0.28567)',
'down (score = 0.03399)',
'c22ebf46_nohash_0.wav',
'stop (score = 0.92027)',
'up (score = 0.06529)',
'_unknown_ (score = 0.01240)',
'2a89ad5c_nohash_0.wav',
'stop (score = 0.56984)',
'up (score = 0.14682)',
'_unknown_ (score = 0.06884)',
'b2ae3928_nohash_0.wav',
'stop (score = 0.98627)',

```

(continues on next page)

(continued from previous page)

```
'_unknown_ (score = 0.01248)',
'down (score = 0.00076)',
'37dca74f_nohash_3.wav',
'stop (score = 0.78066)',
'up (score = 0.07329)',
'_unknown_ (score = 0.07164)',
'3bb68054_nohash_1.wav',
'stop (score = 0.99162)',
'go (score = 0.00331)',
'up (score = 0.00219)',
'a6f2fd71_nohash_1.wav',
'stop (score = 0.59300)',
'up (score = 0.39114)',
'_unknown_ (score = 0.00468)',
'893705bb_nohash_1.wav',
'stop (score = 0.44148)',
'up (score = 0.21584)',
'go (score = 0.08655)',
'46114b4e_nohash_1.wav',
'stop (score = 0.86626)',
'up (score = 0.10812)',
'down (score = 0.00908)',
'32561e9e_nohash_0.wav',
'stop (score = 0.91559)',
'up (score = 0.03920)',
'_unknown_ (score = 0.01798)',
'513aeddf_nohash_4.wav',
'stop (score = 0.95504)',
'_unknown_ (score = 0.04000)',
'go (score = 0.00360)',
'0137b3f4_nohash_3.wav',
'stop (score = 0.74314)',
'up (score = 0.22082)',
'off (score = 0.01753)',
'85851131_nohash_1.wav',
'stop (score = 0.98796)',
'up (score = 0.01117)',
'_unknown_ (score = 0.00050)',
'28612180_nohash_0.wav',
'up (score = 0.42242)',
'stop (score = 0.14160)',
'down (score = 0.09028)',
'e07dd7d4_nohash_0.wav',
'stop (score = 0.34194)',
'up (score = 0.33417)',
'_unknown_ (score = 0.13248)',
'01bb6a2a_nohash_1.wav',
'stop (score = 0.92179)',
'up (score = 0.03660)',
'_unknown_ (score = 0.01734)',
'645ed69d_nohash_3.wav',
'stop (score = 0.99787)',
'up (score = 0.00141)',
'_unknown_ (score = 0.00037)',
'34d5aa5a_nohash_1.wav',
'stop (score = 0.83214)',
'up (score = 0.03625)',
```

(continues on next page)

(continued from previous page)

```

'down (score = 0.02931)',
'333784b7_nohash_3.wav',
'stop (score = 0.97966)',
'up (score = 0.01852)',
'_unknown_ (score = 0.00124)',
'9a69672b_nohash_4.wav',
'stop (score = 0.89391)',
'up (score = 0.07879)',
'go (score = 0.01087)',
'31f01a8d_nohash_4.wav',
'stop (score = 0.97557)',
'up (score = 0.01867)',
'off (score = 0.00205)',
'0d6d7360_nohash_1.wav',
'stop (score = 0.70121)',
'up (score = 0.14588)',
'_unknown_ (score = 0.04814)',
'4a1e736b_nohash_1.wav',
'stop (score = 0.98115)',
'_unknown_ (score = 0.01325)',
'up (score = 0.00284)',
'3b4f8f24_nohash_0.wav',
'stop (score = 0.98360)',
'down (score = 0.00743)',
'_unknown_ (score = 0.00367)',
'982babaf_nohash_1.wav',
'stop (score = 0.98975)',
'_unknown_ (score = 0.00548)',
'up (score = 0.00266)',
'7fd25f7c_nohash_1.wav',
'stop (score = 0.99790)',
'_unknown_ (score = 0.00121)',
'up (score = 0.00055)',
'a7200079_nohash_3.wav',
'stop (score = 0.98081)',
'up (score = 0.00732)',
'_unknown_ (score = 0.00548)',
'af6fbbf5_nohash_0.wav',
'stop (score = 0.99353)',
'_unknown_ (score = 0.00320)',
'up (score = 0.00177)',
'e882abb2_nohash_1.wav',
'stop (score = 0.89639)',
'_unknown_ (score = 0.04231)',
'up (score = 0.04003)',
'7ff4fc72_nohash_0.wav',
'stop (score = 0.65718)',
'down (score = 0.19441)',
'_unknown_ (score = 0.05682)',
'80c45ed6_nohash_0.wav',
'stop (score = 0.94357)',
'down (score = 0.02360)',
'up (score = 0.02062)',
'fc2411fe_nohash_1.wav',
'stop (score = 0.89233)',
'up (score = 0.03672)',
'down (score = 0.02436)',

```

(continues on next page)

(continued from previous page)

```

'54ad8f22_nohash_3.wav',
'stop (score = 0.51013)',
'down (score = 0.17139)',
'go (score = 0.08875)',
'692a88e6_nohash_3.wav',
'stop (score = 0.93051)',
'up (score = 0.06421)',
'_unknown_ (score = 0.00301)',
'893705bb_nohash_6.wav',
'stop (score = 0.44148)',
'up (score = 0.21584)',
'go (score = 0.08655)',
'171edea9_nohash_2.wav',
'stop (score = 0.99756)',
'up (score = 0.00147)',
'_unknown_ (score = 0.00069)',
'f0522ff4_nohash_4.wav',
'stop (score = 0.99686)',
'_unknown_ (score = 0.00174)',
'up (score = 0.00112)',
'824e8ce5_nohash_1.wav',
'stop (score = 0.86248)',
'up (score = 0.06937)',
'_unknown_ (score = 0.01894)',
'a9ca1818_nohash_4.wav',
'stop (score = 0.84088)',
'up (score = 0.14219)',
'_unknown_ (score = 0.00756)',
'48a9f771_nohash_2.wav',
'stop (score = 0.72441)',
'up (score = 0.08968)',
'off (score = 0.05914)',
'6c429c7b_nohash_1.wav',
'up (score = 0.35406)',
'stop (score = 0.30395)',
'off (score = 0.16752)',
'f035e2ea_nohash_3.wav',
'stop (score = 0.99878)',
'_unknown_ (score = 0.00053)',
'up (score = 0.00030)',
'b06c19b0_nohash_0.wav',
'stop (score = 0.99874)',
'_unknown_ (score = 0.00081)',
'up (score = 0.00021)',
'9a356ab9_nohash_0.wav',
'down (score = 0.30055)',
'_unknown_ (score = 0.19656)',
'stop (score = 0.15799)',
'0cd323ec_nohash_1.wav',
'stop (score = 0.41698)',
'off (score = 0.14234)',
'down (score = 0.12729)',
'f19c1390_nohash_0.wav',
'stop (score = 0.99386)',
'up (score = 0.00474)',
'down (score = 0.00070)',
'435695e3_nohash_0.wav',

```

(continues on next page)

(continued from previous page)

```

'stop (score = 0.99527)',
'_unknown_ (score = 0.00271)',
'up (score = 0.00107)',
'179a61b7_nohash_0.wav',
'stop (score = 0.45093)',
'_unknown_ (score = 0.24628)',
'go (score = 0.06381)',
'190821dc_nohash_0.wav',
'stop (score = 0.85112)',
'up (score = 0.06795)',
'go (score = 0.03845)',
'82951cf0_nohash_1.wav',
'stop (score = 0.37752)',
'up (score = 0.32371)',
'_unknown_ (score = 0.05701)',
'bd76a7fd_nohash_4.wav',
'stop (score = 0.98600)',
'up (score = 0.01023)',
'_unknown_ (score = 0.00227)',
'b4ea0d9a_nohash_2.wav',
'stop (score = 0.95925)',
'up (score = 0.03561)',
'_unknown_ (score = 0.00178)',
'e4be0cf6_nohash_0.wav',
'stop (score = 0.43057)',
'up (score = 0.36197)',
'off (score = 0.13706)',
'626e323f_nohash_0.wav',
'stop (score = 0.99196)',
'up (score = 0.00500)',
'down (score = 0.00119)',
'3589bc72_nohash_3.wav',
'stop (score = 0.98538)',
'go (score = 0.00724)',
'_unknown_ (score = 0.00485)',
'9a7c1f83_nohash_0.wav',
'stop (score = 0.97974)',
'_unknown_ (score = 0.01620)',
'go (score = 0.00237)',
'90804775_nohash_2.wav',
'stop (score = 0.45099)',
'_unknown_ (score = 0.12778)',
'off (score = 0.09222)',
'1e412fac_nohash_0.wav',
'stop (score = 0.97331)',
'_unknown_ (score = 0.01048)',
'down (score = 0.00528)',
'72e382bd_nohash_2.wav',
'stop (score = 0.99472)',
'down (score = 0.00212)',
'_unknown_ (score = 0.00174)',
'37d38e44_nohash_0.wav',
'stop (score = 0.87186)',
'down (score = 0.06694)',
'go (score = 0.02199)',
'322d17d3_nohash_3.wav',
'stop (score = 0.99799)',

```

(continues on next page)

(continued from previous page)

```

'up (score = 0.00112)',
'_unknown_ (score = 0.00045)',
'a045368c_nohash_4.wav',
'stop (score = 0.99478)',
'up (score = 0.00325)',
'off (score = 0.00098)',
'b69002d4_nohash_0.wav',
'stop (score = 0.99936)',
'up (score = 0.00028)',
'_unknown_ (score = 0.00019)',
'a7200079_nohash_2.wav',
'stop (score = 0.85205)',
'down (score = 0.12575)',
'_unknown_ (score = 0.00612)',
'5b09db89_nohash_3.wav',
'stop (score = 0.78082)',
'up (score = 0.13418)',
'_unknown_ (score = 0.03140)',
'fa446c16_nohash_3.wav',
'stop (score = 0.99090)',
'up (score = 0.00493)',
'_unknown_ (score = 0.00252)',
'b4ea0d9a_nohash_1.wav',
'stop (score = 0.95521)',
'up (score = 0.03692)',
'_unknown_ (score = 0.00382)',
'493392c6_nohash_0.wav',
'stop (score = 0.97289)',
'up (score = 0.01633)',
'_unknown_ (score = 0.00663)',
'ca4eeab0_nohash_0.wav',
'stop (score = 0.81987)',
'down (score = 0.05813)',
'up (score = 0.05399)',
'f34e6f44_nohash_0.wav',
'stop (score = 0.39819)',
'_unknown_ (score = 0.22448)',
'up (score = 0.10589)',
'92e17cc4_nohash_1.wav',
'stop (score = 0.98871)',
'up (score = 0.00868)',
'_unknown_ (score = 0.00104)',
'36050ef3_nohash_1.wav',
'stop (score = 0.87487)',
'up (score = 0.05612)',
'go (score = 0.02334)',
'6a014b29_nohash_1.wav',
'stop (score = 0.79478)',
'off (score = 0.06393)',
'_unknown_ (score = 0.04445)',
'5b26c81b_nohash_0.wav',
'stop (score = 0.92717)',
'off (score = 0.03468)',
'up (score = 0.02418)',
'88d009d2_nohash_0.wav',
'stop (score = 0.23585)',
'up (score = 0.15615)',

```

(continues on next page)

(continued from previous page)

```

'_unknown_ (score = 0.11477)',
'ab00c4b2_nohash_1.wav',
'stop (score = 0.98693)',
'up (score = 0.00912)',
'down (score = 0.00194)',
'a7acbbeb_nohash_2.wav',
'stop (score = 0.91138)',
'up (score = 0.08291)',
'_unknown_ (score = 0.00190)',
'87d5e978_nohash_1.wav',
'stop (score = 0.99362)',
'up (score = 0.00550)',
'_unknown_ (score = 0.00048)',
'15dd287d_nohash_2.wav',
'stop (score = 0.60951)',
'down (score = 0.11640)',
'go (score = 0.10724)',
'69a1a79f_nohash_2.wav',
'stop (score = 0.99213)',
'up (score = 0.00532)',
'_unknown_ (score = 0.00173)',
'd3831f6a_nohash_0.wav',
'stop (score = 0.99898)',
'up (score = 0.00088)',
'_unknown_ (score = 0.00011)',
'd9b8fab2_nohash_1.wav',
'stop (score = 0.87669)',
'off (score = 0.03555)',
'_unknown_ (score = 0.03319)',
'2dc4f05d_nohash_2.wav',
'stop (score = 0.72243)',
'_unknown_ (score = 0.16989)',
'down (score = 0.03805)',
'5af0ca83_nohash_0.wav',
'stop (score = 0.93210)',
'up (score = 0.03181)',
'down (score = 0.01050)',
'1dc86f91_nohash_2.wav',
'stop (score = 0.99897)',
'_unknown_ (score = 0.00045)',
'up (score = 0.00029)',
'89947bd7_nohash_0.wav',
'stop (score = 0.50612)',
'up (score = 0.47685)',
'off (score = 0.00776)',
'b528edb3_nohash_1.wav',
'stop (score = 0.71747)',
'down (score = 0.13705)',
'_unknown_ (score = 0.06468)',
'c1d39ce8_nohash_0.wav',
'stop (score = 0.66247)',
'go (score = 0.17296)',
'_unknown_ (score = 0.07923)',
'b93528e3_nohash_0.wav',
'stop (score = 0.32753)',
'up (score = 0.17626)',
'off (score = 0.11351)',

```

(continues on next page)

(continued from previous page)

```

'c4a7a867_nohash_0.wav',
'stop (score = 0.98290)',
'up (score = 0.00738)',
'_unknown_ (score = 0.00650)',
'6d1dcca6_nohash_0.wav',
'stop (score = 0.68633)',
'up (score = 0.08010)',
'off (score = 0.07813)',
'26e573a9_nohash_0.wav',
'stop (score = 0.81587)',
'down (score = 0.05165)',
'up (score = 0.04305)',
'5188de0d_nohash_0.wav',
'stop (score = 0.22536)',
'_unknown_ (score = 0.18580)',
'go (score = 0.12804)',
'0585b66d_nohash_3.wav',
'stop (score = 0.94474)',
'up (score = 0.02064)',
'_unknown_ (score = 0.01619)',
'1b835b87_nohash_1.wav',
'stop (score = 0.94834)',
'up (score = 0.01997)',
'down (score = 0.01052)',
'aff582a1_nohash_3.wav',
'stop (score = 0.99678)',
'_unknown_ (score = 0.00310)',
'down (score = 0.00006)',
'f9273a21_nohash_1.wav',
'stop (score = 0.95208)',
'up (score = 0.04081)',
'_unknown_ (score = 0.00239)',
'4a4e28f1_nohash_0.wav',
'stop (score = 0.26048)',
'down (score = 0.16748)',
'go (score = 0.13965)',
'1ed0b13d_nohash_3.wav',
'stop (score = 0.99921)',
'_unknown_ (score = 0.00047)',
'up (score = 0.00021)',
'3ec05c3d_nohash_0.wav',
'stop (score = 0.54378)',
'off (score = 0.19993)',
'_unknown_ (score = 0.07233)',
'7846fd85_nohash_3.wav',
'stop (score = 0.92796)',
'down (score = 0.02070)',
'_unknown_ (score = 0.01309)',
'8dc18a75_nohash_0.wav',
'stop (score = 0.78939)',
'up (score = 0.19312)',
'_unknown_ (score = 0.00795)',
'7846fd85_nohash_4.wav',
'stop (score = 0.89709)',
'down (score = 0.03361)',
'_unknown_ (score = 0.01997)',
'f5626af6_nohash_3.wav',

```

(continues on next page)

(continued from previous page)

```

'stop (score = 0.94139)',
'_unknown_ (score = 0.04825)',
'up (score = 0.00626)',
'ffd2ba2f_nohash_3.wav',
'stop (score = 0.99623)',
'_unknown_ (score = 0.00219)',
'up (score = 0.00138)',
'513aeddf_nohash_2.wav',
'stop (score = 0.89436)',
'_unknown_ (score = 0.05919)',
'go (score = 0.02806)',
'551e42e8_nohash_0.wav',
'stop (score = 0.45233)',
'_unknown_ (score = 0.16313)',
'off (score = 0.11665)',
'26e573a9_nohash_1.wav',
'stop (score = 0.93726)',
'down (score = 0.02824)',
'up (score = 0.01093)',
'9be15e93_nohash_3.wav',
'stop (score = 0.99688)',
'up (score = 0.00160)',
'_unknown_ (score = 0.00101)',
'264f471d_nohash_1.wav',
'stop (score = 0.80808)',
'_unknown_ (score = 0.10039)',
'up (score = 0.03336)',
'b959cd0c_nohash_0.wav',
'stop (score = 0.24488)',
'_unknown_ (score = 0.11815)',
'left (score = 0.10235)',
'bd76a7fd_nohash_2.wav',
'stop (score = 0.99836)',
'up (score = 0.00083)',
'_unknown_ (score = 0.00038)',
'74241b28_nohash_1.wav',
'stop (score = 0.97935)',
'_unknown_ (score = 0.01288)',
'up (score = 0.00481)',
'6ef407da_nohash_1.wav',
'stop (score = 0.98262)',
'up (score = 0.01171)',
'_unknown_ (score = 0.00243)',
'51eefcc6_nohash_0.wav',
'stop (score = 0.64022)',
'go (score = 0.16679)',
'no (score = 0.06918)',
'a827e3a1_nohash_3.wav',
'stop (score = 0.76328)',
'off (score = 0.15193)',
'up (score = 0.02800)',
'0d82fd99_nohash_3.wav',
'stop (score = 0.92439)',
'up (score = 0.05283)',
'_unknown_ (score = 0.00782)',
'5efb758c_nohash_0.wav',
'stop (score = 0.79350)',

```

(continues on next page)

(continued from previous page)

```

'go (score = 0.05066)',
'up (score = 0.04322)',
'6094340e_nohash_1.wav',
'stop (score = 0.98467)',
'down (score = 0.00778)',
'_unknown_ (score = 0.00234)',
'067f61e2_nohash_3.wav',
'stop (score = 0.99034)',
'up (score = 0.00893)',
'_unknown_ (score = 0.00047)',
'54d9ccb5_nohash_0.wav',
'stop (score = 0.98414)',
'up (score = 0.01113)',
'down (score = 0.00169)',
'01b4757a_nohash_0.wav',
'stop (score = 0.67216)',
'up (score = 0.13359)',
'_unknown_ (score = 0.05075)',
'953felad_nohash_2.wav',
'stop (score = 0.78502)',
'up (score = 0.10579)',
'down (score = 0.05232)',
'af790082_nohash_0.wav',
'stop (score = 0.96869)',
'up (score = 0.02890)',
'_unknown_ (score = 0.00105)',
'9a7c1f83_nohash_4.wav',
'stop (score = 0.95722)',
'go (score = 0.02449)',
'down (score = 0.00590)',
'94de6a6a_nohash_1.wav',
'up (score = 0.88006)',
'_unknown_ (score = 0.03862)',
'stop (score = 0.02416)',
'332d33b1_nohash_0.wav',
'stop (score = 0.44896)',
'down (score = 0.21524)',
'no (score = 0.15847)',
'674ca5ea_nohash_0.wav',
'stop (score = 0.93532)',
'up (score = 0.02308)',
'off (score = 0.01062)',
'b97c9f77_nohash_0.wav',
'stop (score = 0.65481)',
'down (score = 0.10290)',
'_unknown_ (score = 0.09691)',
'bfd26d6b_nohash_1.wav',
'stop (score = 0.83370)',
'up (score = 0.15832)',
'_unknown_ (score = 0.00376)',
'f632210f_nohash_1.wav',
'stop (score = 0.68990)',
'up (score = 0.12806)',
'down (score = 0.05177)',
'4290ca61_nohash_0.wav',
'stop (score = 0.86563)',
'off (score = 0.03857)',

```

(continues on next page)

(continued from previous page)

```

'up (score = 0.03350)',
'893705bb_nohash_8.wav',
'stop (score = 0.71492)',
'up (score = 0.16641)',
'go (score = 0.04935)',
'3bfd30e6_nohash_0.wav',
'stop (score = 0.63599)',
'go (score = 0.14249)',
'up (score = 0.10391)',
'8ff44869_nohash_1.wav',
'stop (score = 0.79370)',
'_unknown_ (score = 0.06695)',
'down (score = 0.06271)',
'439c84f4_nohash_3.wav',
'stop (score = 0.93857)',
'go (score = 0.02124)',
'down (score = 0.01481)',
'd264f7b6_nohash_2.wav',
'stop (score = 0.35819)',
'up (score = 0.19114)',
'down (score = 0.08839)',
'b414c653_nohash_4.wav',
'up (score = 0.58615)',
'stop (score = 0.25748)',
'off (score = 0.13660)',
'0d85a428_nohash_0.wav',
'stop (score = 0.16377)',
'up (score = 0.14177)',
'_unknown_ (score = 0.14092)',
'525eaa62_nohash_2.wav',
'stop (score = 0.99972)',
'up (score = 0.00016)',
'_unknown_ (score = 0.00011)',
'6aafb34f_nohash_0.wav',
'stop (score = 0.77719)',
'up (score = 0.11667)',
'_unknown_ (score = 0.03379)',
'226537ab_nohash_0.wav',
'stop (score = 0.97408)',
'up (score = 0.02085)',
'_unknown_ (score = 0.00265)',
'559bc36a_nohash_1.wav',
'stop (score = 0.97451)',
'go (score = 0.01067)',
'_unknown_ (score = 0.00638)',
'df280250_nohash_1.wav',
'stop (score = 0.98016)',
'down (score = 0.00753)',
'go (score = 0.00371)',
'b5cf6ea8_nohash_4.wav',
'stop (score = 0.99810)',
'up (score = 0.00149)',
'_unknown_ (score = 0.00032)',
'587f3271_nohash_0.wav',
'stop (score = 0.85410)',
'_unknown_ (score = 0.06047)',
'up (score = 0.02200)',

```

(continues on next page)

(continued from previous page)

```

'1ecfb537_nohash_4.wav',
'stop (score = 0.99523)',
'_unknown_ (score = 0.00262)',
'left (score = 0.00094)',
'321aba74_nohash_1.wav',
'stop (score = 0.99719)',
'_unknown_ (score = 0.00182)',
'up (score = 0.00048)',
'5170b77f_nohash_2.wav',
'stop (score = 0.97167)',
'up (score = 0.01566)',
'_unknown_ (score = 0.00619)',
'd78858d9_nohash_0.wav',
'stop (score = 0.52191)',
'up (score = 0.22323)',
'left (score = 0.05919)',
'113b3fbc_nohash_0.wav',
'stop (score = 0.64668)',
'down (score = 0.20922)',
'no (score = 0.04130)',
'3d9200b9_nohash_0.wav',
'stop (score = 0.61945)',
'up (score = 0.12795)',
'off (score = 0.09970)',
'06f6c194_nohash_4.wav',
'stop (score = 0.58292)',
'up (score = 0.39636)',
'off (score = 0.00788)',
'093f65a1_nohash_0.wav',
'stop (score = 0.96773)',
'_unknown_ (score = 0.00939)',
'up (score = 0.00694)',
'e41a903b_nohash_0.wav',
'stop (score = 0.99393)',
'_unknown_ (score = 0.00234)',
'up (score = 0.00195)',
'1b88bf70_nohash_0.wav',
'stop (score = 0.89692)',
'up (score = 0.04895)',
'off (score = 0.02443)',
'1ffd513b_nohash_0.wav',
'stop (score = 0.99402)',
'up (score = 0.00317)',
'down (score = 0.00159)',
'ecef25ba_nohash_0.wav',
'stop (score = 0.18026)',
'left (score = 0.13037)',
'down (score = 0.12699)',
'837a0f64_nohash_2.wav',
'stop (score = 0.99877)',
'_unknown_ (score = 0.00084)',
'up (score = 0.00034)',
'64220627_nohash_1.wav',
'stop (score = 0.48925)',
'_unknown_ (score = 0.12916)',
'right (score = 0.11773)',
'bd8412df_nohash_1.wav',

```

(continues on next page)

(continued from previous page)

```

'stop (score = 0.97986)',
'_unknown_ (score = 0.01157)',
'up (score = 0.00352)',
'ed3c2d05_nohash_0.wav',
'stop (score = 0.98379)',
'up (score = 0.01037)',
'_unknown_ (score = 0.00274)',
'578d3efb_nohash_3.wav',
'stop (score = 0.99270)',
'up (score = 0.00675)',
'_unknown_ (score = 0.00050)',
'3df9a3d4_nohash_0.wav',
'stop (score = 0.27674)',
'up (score = 0.25178)',
'no (score = 0.10935)',
'2f0ce4d9_nohash_2.wav',
'stop (score = 0.86411)',
'off (score = 0.06190)',
'_unknown_ (score = 0.03669)',
'b3bdded5_nohash_2.wav',
'stop (score = 0.88871)',
'up (score = 0.09937)',
'_unknown_ (score = 0.00315)',
'a8e25ebb_nohash_0.wav',
'stop (score = 0.48994)',
'up (score = 0.12429)',
'no (score = 0.10686)',
'ad63d93c_nohash_1.wav',
'stop (score = 0.51262)',
'up (score = 0.38220)',
'_unknown_ (score = 0.05622)',
'4c3cddb8_nohash_4.wav',
'stop (score = 0.99406)',
'up (score = 0.00474)',
'_unknown_ (score = 0.00085)',
'0132a06d_nohash_2.wav',
'stop (score = 0.98721)',
'up (score = 0.01022)',
'_unknown_ (score = 0.00182)',
'7846fd85_nohash_1.wav',
'stop (score = 0.92152)',
'down (score = 0.04298)',
'no (score = 0.00939)',
'a04817c2_nohash_1.wav',
'stop (score = 0.98223)',
'up (score = 0.01197)',
'down (score = 0.00265)',
'605ed0ff_nohash_0.wav',
'stop (score = 0.45783)',
'up (score = 0.13047)',
'_unknown_ (score = 0.08471)',
'90e72357_nohash_2.wav',
'stop (score = 0.98573)',
'_unknown_ (score = 0.00698)',
'go (score = 0.00461)',
'2da58b32_nohash_4.wav',
'stop (score = 0.79965)',

```

(continues on next page)

(continued from previous page)

```

'off (score = 0.06937)',
'up (score = 0.04585)',
'28ce0c58_nohash_1.wav',
'go (score = 0.48640)',
'up (score = 0.16947)',
'_unknown_ (score = 0.15874)',
'24a3e589_nohash_2.wav',
'stop (score = 0.98158)',
'up (score = 0.01708)',
'_unknown_ (score = 0.00080)',
'2579e514_nohash_1.wav',
'stop (score = 0.48661)',
'up (score = 0.15354)',
'down (score = 0.10605)',
'437455be_nohash_0.wav',
'stop (score = 0.33172)',
'up (score = 0.12149)',
'_unknown_ (score = 0.10527)',
'1b4c9b89_nohash_0.wav',
'stop (score = 0.98440)',
'up (score = 0.00508)',
'_unknown_ (score = 0.00465)',
'f798ac78_nohash_4.wav',
'stop (score = 0.95439)',
'up (score = 0.03564)',
'_unknown_ (score = 0.00394)',
'29229c21_nohash_1.wav',
'stop (score = 0.99749)',
'up (score = 0.00134)',
'_unknown_ (score = 0.00095)',
'd78858d9_nohash_2.wav',
'stop (score = 0.70006)',
'up (score = 0.20123)',
'left (score = 0.03805)',
'5ebc1cda_nohash_6.wav',
'stop (score = 0.63192)',
'_unknown_ (score = 0.15333)',
'go (score = 0.12291)',
'ab46af55_nohash_0.wav',
'stop (score = 0.99821)',
'up (score = 0.00156)',
'_unknown_ (score = 0.00010)',
'2296blaf_nohash_1.wav',
'stop (score = 0.92353)',
'_unknown_ (score = 0.02751)',
'down (score = 0.02737)',
'9a69672b_nohash_2.wav',
'stop (score = 0.99818)',
'up (score = 0.00158)',
'_unknown_ (score = 0.00013)',
'964e8cfd_nohash_1.wav',
'stop (score = 0.95619)',
'up (score = 0.02845)',
'_unknown_ (score = 0.00523)',
'cb72dfb6_nohash_0.wav',
'stop (score = 0.93942)',
'up (score = 0.03724)',

```

(continues on next page)

(continued from previous page)

```

'go (score = 0.00980)',
'9be15e93_nohash_4.wav',
'stop (score = 0.99567)',
'up (score = 0.00360)',
'_unknown_ (score = 0.00044)',
'5ebc1cda_nohash_1.wav',
'stop (score = 0.75682)',
'_unknown_ (score = 0.09401)',
'down (score = 0.04084)',
'1e9b215e_nohash_1.wav',
'stop (score = 0.46735)',
'down (score = 0.11926)',
'_unknown_ (score = 0.11014)',
'28ce0c58_nohash_4.wav',
'_unknown_ (score = 0.35277)',
'stop (score = 0.25628)',
'up (score = 0.22787)',
'0135f3f2_nohash_1.wav',
'down (score = 0.64915)',
'stop (score = 0.13520)',
'go (score = 0.06183)',
'742d6431_nohash_0.wav',
'stop (score = 0.51185)',
'up (score = 0.45758)',
'off (score = 0.01466)',
'563aa4e6_nohash_3.wav',
'stop (score = 0.98674)',
'_unknown_ (score = 0.00840)',
'up (score = 0.00273)',
'bbaa7946_nohash_0.wav',
'stop (score = 0.99574)',
'down (score = 0.00201)',
'_unknown_ (score = 0.00110)',
'a7acbbeb_nohash_1.wav',
'stop (score = 0.78802)',
'up (score = 0.11744)',
'down (score = 0.04498)',
'ec21c46b_nohash_0.wav',
'stop (score = 0.97416)',
'_unknown_ (score = 0.01517)',
'up (score = 0.00404)',
'18ffa72d_nohash_1.wav',
'off (score = 0.16264)',
'_silence_ (score = 0.11814)',
'yes (score = 0.09154)',
'32ad5b65_nohash_0.wav',
'stop (score = 0.99422)',
'up (score = 0.00445)',
'_unknown_ (score = 0.00075)',
'fd32732a_nohash_0.wav',
'stop (score = 0.98304)',
'_unknown_ (score = 0.00850)',
'down (score = 0.00307)',
'a2b16113_nohash_0.wav',
'stop (score = 0.91704)',
'up (score = 0.07075)',
'off (score = 0.00562)',

```

(continues on next page)

(continued from previous page)

```

'92b0a735_nohash_0.wav',
'stop (score = 0.96610)',
'up (score = 0.01413)',
'down (score = 0.00977)',
'6124b431_nohash_1.wav',
'stop (score = 0.91341)',
'up (score = 0.04878)',
'down (score = 0.01150)',
'7fb8d703_nohash_2.wav',
'stop (score = 0.98021)',
'up (score = 0.01911)',
'_unknown_ (score = 0.00029)',
'5e3dde6b_nohash_1.wav',
'stop (score = 0.87034)',
'up (score = 0.06751)',
'_unknown_ (score = 0.04405)',
'ce7a8e92_nohash_1.wav',
'stop (score = 0.13540)',
'up (score = 0.11850)',
'_unknown_ (score = 0.10476)',
'c0fb6812_nohash_0.wav',
'stop (score = 0.67461)',
'up (score = 0.11249)',
'down (score = 0.05729)',
'c79159aa_nohash_4.wav',
'stop (score = 0.98290)',
'up (score = 0.00918)',
'_unknown_ (score = 0.00575)',
'b3327675_nohash_1.wav',
'stop (score = 0.49829)',
'down (score = 0.13202)',
'_unknown_ (score = 0.11668)',
'87070229_nohash_4.wav',
'stop (score = 0.98843)',
'_unknown_ (score = 0.00946)',
'go (score = 0.00120)',
'a5dlbecc_nohash_2.wav',
'stop (score = 0.42095)',
'_unknown_ (score = 0.13196)',
'up (score = 0.09977)',
'0ff728b5_nohash_3.wav',
'stop (score = 0.97009)',
'_unknown_ (score = 0.02019)',
'go (score = 0.00626)',
'964e8cfd_nohash_4.wav',
'stop (score = 0.98068)',
'up (score = 0.00856)',
'_unknown_ (score = 0.00304)',
'017c4098_nohash_4.wav',
'stop (score = 0.99057)',
'up (score = 0.00858)',
'down (score = 0.00050)',
'b72e58c9_nohash_0.wav',
'stop (score = 0.98696)',
'_unknown_ (score = 0.00550)',
'up (score = 0.00223)',
'ef77b778_nohash_2.wav',

```

(continues on next page)

(continued from previous page)

```
'stop (score = 0.93676)',
'up (score = 0.06078)',
'_unknown_ (score = 0.00106)',
'0f3f64d5_nohash_1.wav',
'stop (score = 0.84753)',
'up (score = 0.11528)',
'down (score = 0.01605)',
'131e738d_nohash_4.wav',
'stop (score = 0.86223)',
'_unknown_ (score = 0.04718)',
'go (score = 0.03105)',
'bab36420_nohash_1.wav',
'stop (score = 0.99879)',
'up (score = 0.00114)',
'_unknown_ (score = 0.00003)',
'27c24504_nohash_0.wav',
'stop (score = 0.91839)',
'up (score = 0.03990)',
'_unknown_ (score = 0.01542)',
'21e8c417_nohash_0.wav',
'no (score = 0.16717)',
'go (score = 0.14407)',
'down (score = 0.13581)',
'54d9ccb5_nohash_1.wav',
'stop (score = 0.98426)',
'up (score = 0.01385)',
'_unknown_ (score = 0.00076)',
'9448c397_nohash_4.wav',
'stop (score = 0.43997)',
'_unknown_ (score = 0.21405)',
'go (score = 0.05974)',
'c7dc7278_nohash_0.wav',
'stop (score = 0.99714)',
'up (score = 0.00277)',
'_unknown_ (score = 0.00007)',
...]
```

That was not pretty! We'd better define some helper functions to extract the model's guesses from that messy output:

```
[15]: def get_guesses(scores):
    scores = filter_scores(scores)
    if len(scores) % 4 != 0:
        raise ValueError(f"Expected scores list to have a length divisible by 4 after_
→filtering but got length {len(scores)}")
    num_files = len(scores) / 4
    fnames = scores[0::4]
    guesses = [guess.split(' ')[0] for guess in scores[1::4]]
    return zip(fnames, guesses)

def score_directory(directory):
    scores = !python {example_path}/label_wav_dir.py \
        --graph={example_path}/trained_model/my_frozen_graph.pb \
        --labels={example_path}/trained_model/conv_labels.txt \
        --wav_dir={directory}
    return filter_scores(scores)
```

Define a function to generate errors in all wav files in a given directory. If an inclusion list is provided, only files on

the list will be processed.

```
[16]: def errorify_directory(data_root_dir, dir_name, tree_root, err_params, inclusion_
    ↳list=None):
    clean_data_dir = data_root_dir / dir_name
    if not clean_data_dir.exists():
        raise ValueError(f"Directory {clean_data_dir} does not exist.")
    err_data_dir = data_root_dir / (dir_name + "_err")
    if not err_data_dir.exists():
        err_data_dir.mkdir()
    if not inclusion_list:
        inclusion_list = [f for f in clean_data_dir.iterdir() if ".wav" in str(f)]
    for file in inclusion_list:
        fname = file.name
        wav = read(file)
        clipped = tree_root.generate_error([wav], err_params)[0]
        err_file_path = err_data_dir / fname
        write(err_file_path, clipped[0], clipped[1])
    return err_data_dir
```

Define a function to generate errors in all wav files on a list. The function is needed when files from multiple categories are present on the list. To facilitate comparisons between clean and errorified data, the clean files the list can be automatically copied to suitably named directories. To do this, provide the parameter `copy_clean=True`.

```
[17]: def errorify_list(data_files, categories, tree_root, err_params, copy_clean=False):
    data_root_dir = data_files[0].parents[1]
    for cat in categories:
        files_in_cat = [f for f in data_files if (cat + "/" in str(f))]
        print("category:", cat)
        print(f"{len(files_in_cat)}")
        errorify_directory(data_root_dir, cat, tree_root, err_params, inclusion_
    ↳list=files_in_cat)
        if copy_clean:
            copy_dir = data_root_dir / (cat + "_clean")
            copy_dir.mkdir(exist_ok=True)
            for file in files_in_cat:
                shutil.copy(file, copy_dir)
```

Define a function to compare the model's guesses on clean and errorified data. The results are returned in a Pandas dataframe.

```
[18]: def compare(data_root, category, clean_ext="_clean", err_ext="_err"):
    scores_clean = score_directory(data_root / (category + clean_ext))
    guesses_clean = get_guesses(scores_clean)
    scores_err = score_directory(data_root / (category + err_ext))
    guesses_err = get_guesses(scores_err)
    df_clean = pd.DataFrame(guesses_clean, columns=["file", "clean_guess"])
    df_err = pd.DataFrame(guesses_err, columns=["file", "err_guess"])
    res = pd.merge(df_clean, df_err, on="file", how="inner")
    res['true_label'] = category
    return res
```

Generate errors in all test set audio clips.

```
[19]: errorify_list(test_set_files, trained_categories, root_node, err_params, copy_
    ↳clean=True)
```

```
category: yes
419
category: no
405
category: up
425
category: down
406
category: left
412
category: right
396
category: on
396
category: off
402
category: stop
411
category: go
402
```

Run model on clean and errorified data.

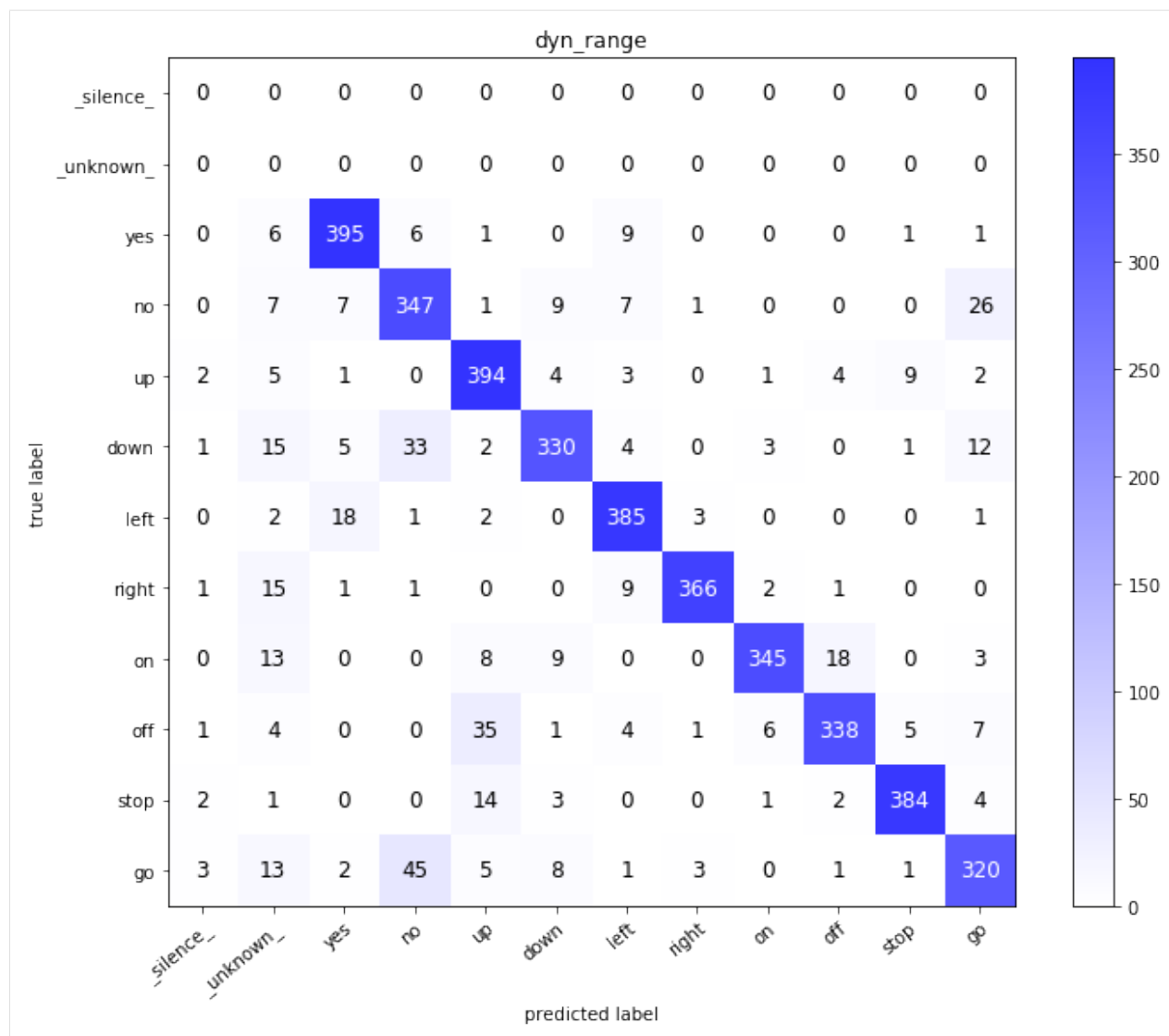
```
[20]: results = [compare(data_dir, cat) for cat in trained_categories]
df = pd.concat(results)
```

Create confusion matrices for clean and errorified data, respectively.

```
[21]: cm_clean = confusion_matrix(df['true_label'], df['clean_guess'], labels=labels)
cm_err = confusion_matrix(df['true_label'], df['err_guess'], labels=labels)
```

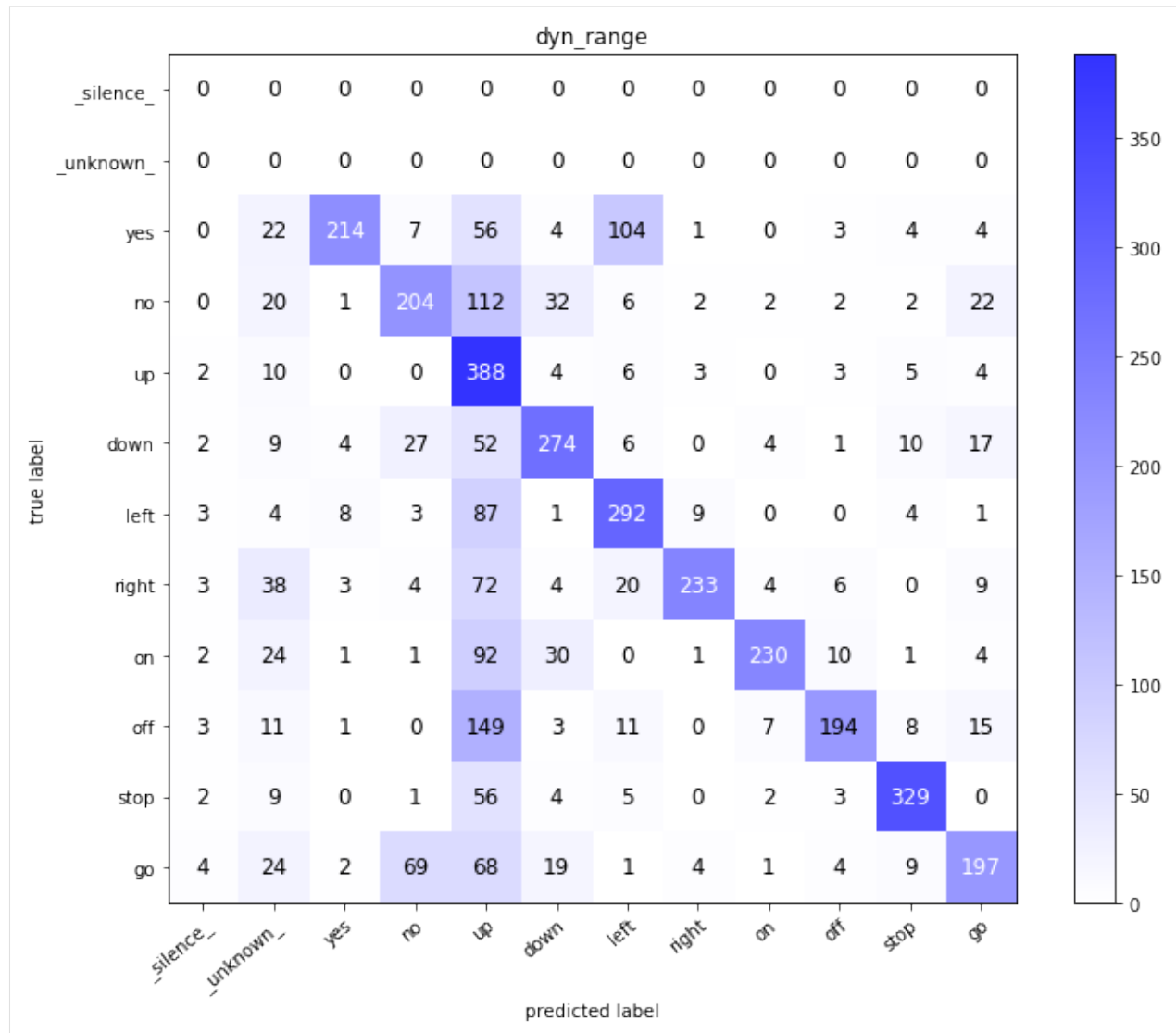
Visualize the confusion matrix for the clean data.

```
[22]: visualize_confusion_matrix(df, cm_clean, 0, labels, "dyn_range", "true_label", "clean_
↪guess")
```

Visualize the confusion matrix for the errorified data.

```
[23]: visualize_confusion_matrix(df, cm_err, 0, labels, "dyn_range", "true_label", "err_
      ↪ guess")
```



The notebook for this case study can be found [here](#).

MODULE DOCUMENTATION

- modindex