

▼ Machine Learning MIIS (lab 1)

(Due to October 26)

The goal of this task is to familiarize with the Colab environment and with [scikit-learn](#), and to gain some experience on learning simple predictive models.

We will focus on the training error only and consider a training dataset, i.e., we will not look at generalization error on testing data.

Consider the LIBSVM repository of datasets: <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

```
# # Example of downloading a dataset
# !wget -t inf https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary/ala
```

```
# # Example of loading a dataset
# from sklearn.datasets import load_svmlight_file
# X_train, y_train = load_svmlight_file("ala")
# print(X_train)
```

▼ 1. Regression Task

Instructions

Choose a regression dataset and apply linear regression on a random subset of the training set of increasing size. You should select training sets that include more and more data points.

1.1. *Plot the approximation error (square loss) on the training set as a function of the number of samples N (i.e. data points in the training set). [You can use [matplotlib](#) for plotting. You can average curves over several permutations of the samples in the training dataset to obtain smoother curves, or even use errorbars.]*

1.2. *Plot the cpu-time as a function of N .*

1.3. *Explain in detail the behaviour of both curves and relate this behavior with what you would expect from theory.*

1.4. Explore how the learned weights change as a function of N . Can you find an interpretation for the learned weights? [You can, for example, plot the value of each weight (in a different colour) as a function of N].

▼ Solution

We selected the Geographical Analysis Spatial Data (`space_ga`) regression dataset, which can be found [here](#). It consists of just over 3 thousand examples reporting votes for the 1980 US presidential election, by county.

The dependent variable `VOTES` describes the number of votes cast.

The dataset includes the following 6 features:

- `POP`: no. people of voting age (18+)
- `EDUCATION`: no. people with 12th grade or higher education
- `HOUSES`: no. of owner-occupied houses
- `INCOME`: aggregate income
- `XCOORD`: X spatial coordinate of county
- `YCOORD`: Y spatial coordinate of county

```
# Import sklearn stuff
from sklearn.datasets import load_svmlight_file
from sklearn.utils import shuffle
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
```

```
# Import process_time() for cpu time
from time import process_time
```

```
# Import mean for easy averaging of model coefficients
from statistics import mean
```

```
# Import matplotlib for plotting
```

```

import matplotlib.pyplot as plt
%matplotlib inline

# Set figure and font size
plt.rcParams["figure.figsize"] = (10,8)
plt.rcParams.update({'font.size': 18})

# Download dataset
!wget -t inf https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/regression/space_ga

--2021-10-26 09:59:11--  https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/regression/space_ga
Resolving www.csie.ntu.edu.tw (www.csie.ntu.edu.tw)... 140.112.30.26
Connecting to www.csie.ntu.edu.tw (www.csie.ntu.edu.tw)|140.112.30.26|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 565473 (552K)
Saving to: 'space_ga.1'

space_ga.1          100%[=====>] 552.22K   620KB/s   in 0.9s

2021-10-26 09:59:13 (620 KB/s) - 'space_ga.1' saved [565473/565473]

```

```

# Load dataset
X, y = load_svmlight_file("space_ga")

# Print dataset info
print(f'X shape : {X.shape}')
print(f'X format : {type(y)}')

X shape : (3107, 6)
X format : <class 'numpy.ndarray'>

```

```

def linreg(X, y, size_increment):
    """
    Selects subsets of increasing size, applies linear regression to each subset
    and returns square loss, cpu times and model coefficients for each subset

    NOTE: cpu time refers to time taken for model fitting

```

```

:param numpy.ndarray X: matrix of features
:param numpy.ndarray y: matrix of labels
:param int size_increment: size increment for each subset

:return dict results: loss cpu time and model coefficients for each subset
indexed by subset size
"""
results = {}

X, y = shuffle(X, y)

curr_subset_size = size_increment

n_subsets = X.shape[0] // size_increment

for n in range(n_subsets):
    # Get the first `curr_subset_size` rows from shuffled dataset
    X_subset, y_subset = X[:curr_subset_size], y[:curr_subset_size]

    start = process_time()
    model = LinearRegression().fit(X_subset, y_subset)
    end = process_time()

    results[curr_subset_size] = {}
    results[curr_subset_size]["cpu_time"] = end - start

    y_pred = model.predict(X)
    results[curr_subset_size]["mse"] = mean_squared_error(y, y_pred)

    results[curr_subset_size]["coef"] = model.coef_

    curr_subset_size += size_increment

    # Last subset should take remaining samples if division wasn't whole
    if n == n_subsets-2:
        curr_subset_size = X.shape[0]

return results

```

```

def avg_results_linreg(X, y, size_increment, n):
    """
    Runs linear regression on random subsets of increasing sizes n times, and
    computes average square loss and cpu time for each subset size

    NOTE: cpu time refers to time taken for model fitting

    :param numpy.ndarray X: matrix of features
    :param numpy.ndarray y: matrix of labels
    :param int size_increment: size increment for each subset
    :param int n: no. of times lin reg should be applied to random subsets

    :return dict avg_results: average loss and cpu time for each subset indexed by
    subset size
    """
    avg_results = {}

    for i in range(n):
        results = linreg(X, y, size_increment)

        # Store results in lists by subset size
        for size in results:
            if size in avg_results:
                avg_results[size]["mse"].append(results[size]["mse"])
                avg_results[size]["cpu_time"].append(results[size]["cpu_time"])
            else:
                avg_results[size] = {}
                avg_results[size]["mse"] = [results[size]["mse"]]
                avg_results[size]["cpu_time"] = [results[size]["cpu_time"]]

    # Average losses and cpu times for each subset size
    for size in avg_results.keys():
        losses = avg_results[size]["mse"]
        avg_results[size]["mse"] = sum(losses) / len(losses)

        cpu_times = avg_results[size]["cpu_time"]
        avg_results[size]["cpu_time"] = sum(cpu_times) / len(cpu_times)

    return avg_results

```

```
results = avg_results_linreg(X, y, size_increment=500, n=50)
```

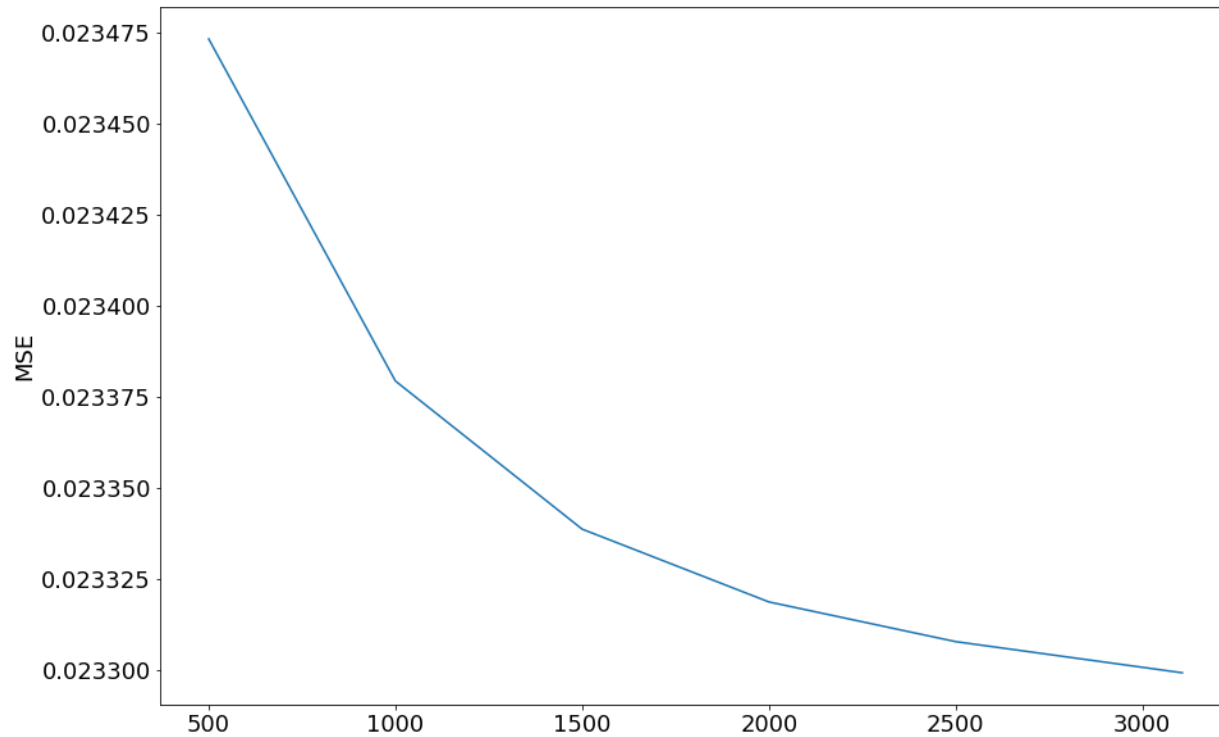
```
results
```

```
{500: {'cpu_time': 0.0016711040799999388, 'mse': 0.023473169762590288},  
1000: {'cpu_time': 0.001753435719999974, 'mse': 0.02337923207647167},  
1500: {'cpu_time': 0.0018810207800001066, 'mse': 0.02333855936282825},  
2000: {'cpu_time': 0.0019144412400000554, 'mse': 0.02331857160154746},  
2500: {'cpu_time': 0.002028778220000049, 'mse': 0.02330766979687296},  
3107: {'cpu_time': 0.0021679923799999834, 'mse': 0.023299114264797728}}
```

1.1. Plot the square loss as a function of N

```
# Parse results into lists for plotting  
sizes = []  
losses = []  
for size, result in results.items():  
    sizes.append(size)  
    losses.append(result["mse"])  
  
# Plot square loss against N  
plt.xlabel("N")  
plt.ylabel("MSE")  
plt.plot(sizes, losses)
```

[<matplotlib.lines.Line2D at 0x7f06c6a3cbd0>]

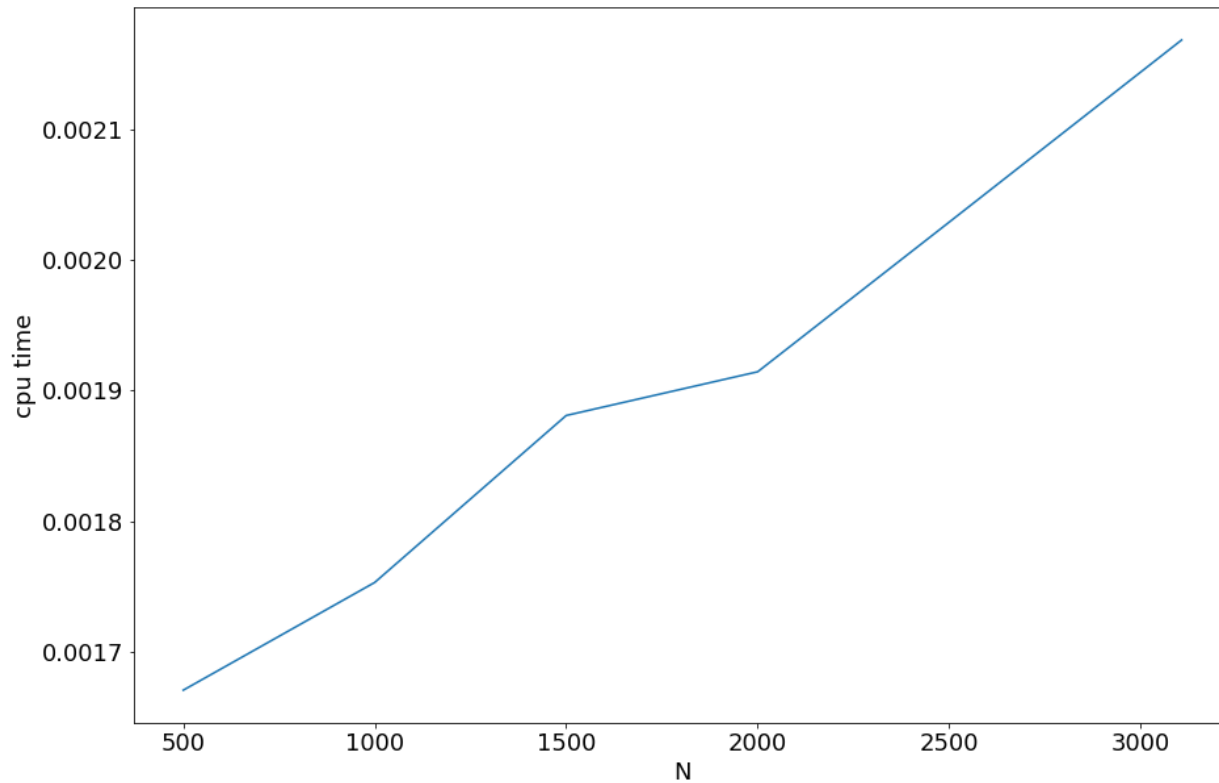


1.2. Plot the cpu-time as a function of N .

```
# Parse results into lists for plotting
sizes = []
times = []
for size,result in results.items():
    sizes.append(size)
    times.append(result["cpu_time"])

# Plot cpu time against N
plt.xlabel("N")
plt.ylabel("cpu time")
plt.plot(sizes, times)
```

```
[<matplotlib.lines.Line2D at 0x7f06c69a8750>]
```



1.3. Explain in detail the behaviour of both curves and relate this behavior with what you would expect from theory.

We ran linear regression on subsets of increasing sizes where each subset contains 500 more examples than the previous. We repeated this operation 50 times and averaged the results.

The first curve shows the average loss as a function of the number of training samples. As we can observe, the larger training set the lower the loss value. This behaviour is expected, since fitting a model on a larger training set allows us to obtain to better coefficients,

which in turn cause the model to generate more accurate predictions. Since the loss function computes the distance between the gold labels and model outputs, the more accurate the predictions are the lower the loss value will be. In this case, the decrease is steeper between subsets of 500 and 1,500 examples, after which the loss continues to decrease, albeit at a lower rate. However it should be noted that the numeric change in loss is actually very slight, only ranging between around .0234 to around .0233.

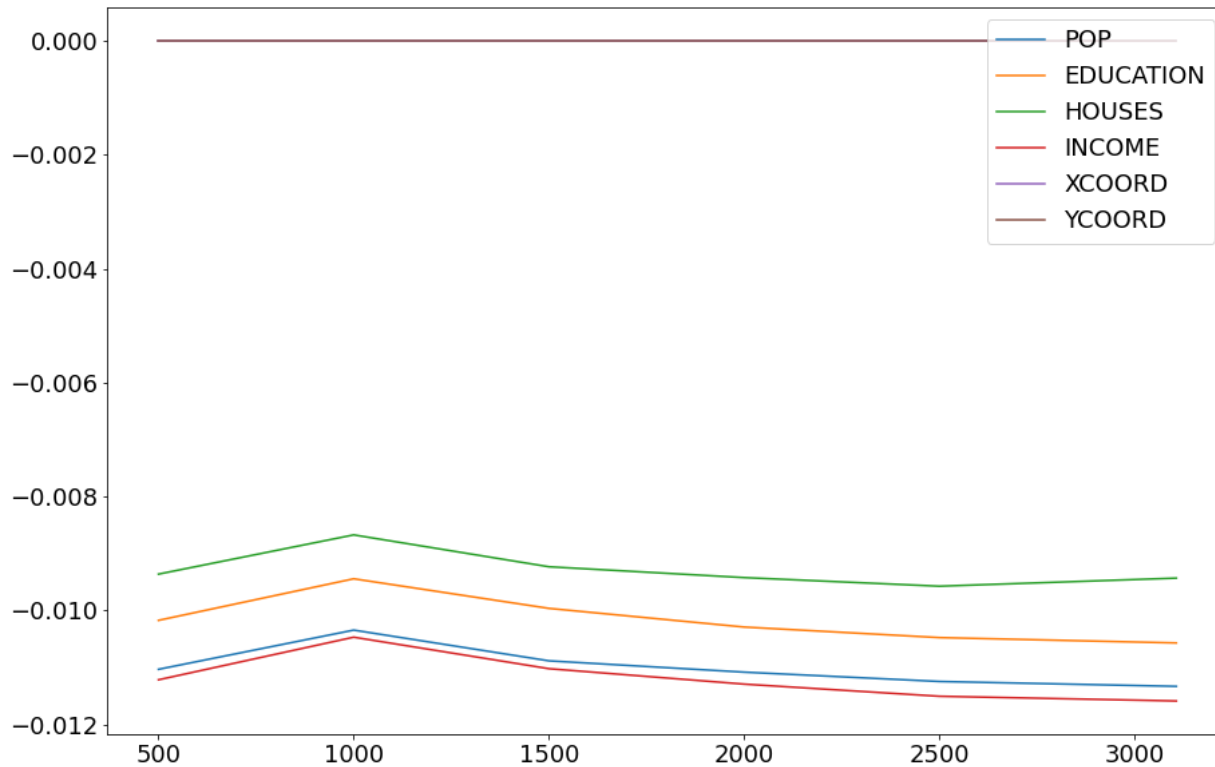
The second curve shows the average cpu time for model fitting as a function of the number of training samples. It can be observed that larger training sets correspond to longer cpu time. This is expected, because a larger training set contains more examples, and the model requires more calculations to fit the weights. More expensive computations trivially lead to an increased cpu time.

1.4. Explore how the learned weights change as a function of N . Can you find an interpretation for the learned weights?

[You can, for example, plot the value of each weight (in a different colour) as a function of N].

```
# Run linreg once
results = linreg(X, y, size_increment=500)
# Parse results into lists for plotting
sizes = []
coefs = []
for size,result in results.items():
    sizes.append(size)
    coefs.append(result["coef"])
```

```
# Plot coefs against N
plot = plt.plot(sizes, coefs)
labels = ["POP", "EDUCATION", "HOUSES", "INCOME", "XCOORD", "YCOORD"]
plt.legend(plot, labels, loc="upper right")
plt.show()
```



We ran linear regression on subsets of increasing size where each subset contains 500 more examples than the previous. Each line in the figure above shows how the corresponding weight changes when the training size grows, according to the top right legend.

We can observe that the weight associated with `YCOORD` does not significantly change with N , remaining close to .00018. This can be interpreted as a weak positive correlation with `VOTES`, corresponding to a .00018 (approx) increase in cast votes for every 1-unit increase in `YCOORD`. If we assume that `YCOORD` places the county on a point between two cardinal points on the US map, for example South vs North, then the higher the `YCOORD` the more North the county is located. Then, this correlation could be interpreted as describing a tendency for Northern counties to cast slightly more votes.

The other weights are noticeably all negative, albeit very small. The changes for these are much more noticeable, and they all follow a similar pattern. Starting from a certain value, there is an increase between 500 and 1000 examples. Then all the weights slowly descend

until they reach a value that in most cases is lower than the one from the first model. In general, all of these can be interpreted as describing a slightly negative correlation with VOTES, i.e. a decrease in VOTES for every 1-unit increase wrt corresponding features.

▼ 2. Classification Task

Instructions

Choose a classification dataset and apply logistic regression. Repeat the previous steps using the mean accuracy instead of the squared loss.

2.1. *Plot the mean accuracy on the training set as a function of the number of samples N , i.e., data points in the training set. [You can use the scikit-learn implementation of [LogisticRegression](#), but check carefully the default parameter values. Do not use regularization].*

2.2. *Plot the cpu-time as a function of N .*

2.3. *Explain in detail the behaviour of both curves and relate this behavior with what you would expect from theory.*

2.4. *Explore how the learned weights change as a function of N . Can you find an interpretation for the learned weights?*

▼ Solution

We chose the German Credit Data dataset in its scaled numeric format (`german.numeric_scale`), which can be found [here](#). It consists of 1,000 examples classifying people as good or bad credit risks using 24 features.

The dependent variable `y` describes whether the person is good or bad risk.

The original dataset contained 20 features, but some of its features were categorical. To obtain a numeric version of the dataset, some of these categorical features were turned into numeric indicators, for the purpose of making the dataset suitable for algorithms unable to use categorical variables.

We opted for the scaled version to favor convergence of our logistic regression model. The features were thus scaled to lie between -1 and 1.

Some examples of features:

- `AGE`
- `CREDIT_AMOUNT`

- `DURATION_CURR_EMPL`: duration (years) of current employment
- `DURATION_CURR_RES`: duration (years) of current residence
- `HOUSING_RENT`: if person rents the property they live in
- `HOUSING_OWN`: if person owns the property they live in
- `NO_CARDS`: no. of existing credits at this bank

```
# Import sklearn stuff
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

```
# Download dataset
!wget -t inf https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary/german.numer_scale
```

```
--2021-10-26 10:16:55-- https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary/german.numer_scale
Resolving www.csie.ntu.edu.tw (www.csie.ntu.edu.tw)... 140.112.30.26
Connecting to www.csie.ntu.edu.tw (www.csie.ntu.edu.tw)|140.112.30.26|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 165393 (162K)
Saving to: 'german.numer_scale'

german.numer_scale  100%[=====>] 161.52K   303KB/s   in 0.5s

2021-10-26 10:16:57 (303 KB/s) - 'german.numer_scale' saved [165393/165393]
```

```
# Load dataset
X, y = load_svmlight_file("german.numer_scale")

# Print dataset info
print(f'X shape : {X.shape}')
print(f'X format : {type(y)}')
```

```
X shape : (1000, 24)
X format : <class 'numpy.ndarray'>
```

```
def logreg(X, y, size_increment):
    """
```

Selects subsets of increasing size, applies logistic regression to each subset and returns accuracy, cpu times and model coefficients for each subset

NOTE: cpu time refers to time taken for model fitting

```
:param scipy.sparse.csr.csr_matrix X: matrix of features
:param scipy.sparse.csr.csr_matrix y: matrix of labels
:param int size_increment: size increment for each subset

:return dict results: accuracy cpu time and model coefficients for each subset
indexed by subset size
"""
results = {}

X, y = shuffle(X, y)

curr_subset_size = size_increment

n_subsets = X.shape[0] // size_increment

for n in range(n_subsets):
    # Get the first `curr_subset_size` rows from shuffled dataset
    X_subset, y_subset = X[:curr_subset_size], y[:curr_subset_size]

    start = process_time()
    model = LogisticRegression().fit(X_subset, y_subset)
    end = process_time()

    results[curr_subset_size] = {}
    results[curr_subset_size]["cpu_time"] = end - start

    y_pred = model.predict(X)
    results[curr_subset_size]["accuracy"] = accuracy_score(y, y_pred)

    results[curr_subset_size]["coef"] = model.coef_

    curr_subset_size += size_increment

# Last subset should take remaining samples if division wasn't whole
```

```

    if n == n_subsets-2:
        curr_subset_size = X.shape[0]

return results

```

```

def avg_results_logreg(X, y, size_increment, n):
    """
    Runs logistic regression on random subsets of increasing sizes n times, and
    computes average square loss and cpu time for each subset size

    NOTE: cpu time refers to time taken for model fitting

    :param scipy.sparse.csr.csr_matrix X: matrix of features
    :param scipy.sparse.csr.csr_matrix y: matrix of labels
    :param int size_increment: size increment for each subset
    :param int n: no. of times log reg should be applied to random subsets

    :return dict avg_results: average accuracy cpu time and coefficients for each
    subset indexed by subset size
    """
    avg_results = {}

    for i in range(n):
        results = logreg(X, y, size_increment)

        # Store results in lists by subset size
        for size in results:
            if size in avg_results:
                avg_results[size]["accuracy"].append(results[size]["accuracy"])
                avg_results[size]["cpu_time"].append(results[size]["cpu_time"])
            else:
                avg_results[size] = {}
                avg_results[size]["accuracy"] = [results[size]["accuracy"]]
                avg_results[size]["cpu_time"] = [results[size]["cpu_time"]]

    # Average accuracies, cpu times and coefficients for each subset size
    for size in avg_results.keys():
        accuracies = avg_results[size]["accuracy"]

```

```

    avg_results[size]["accuracy"] = sum(accuracies) / len(accuracies)

    cpu_times = avg_results[size]["cpu_time"]
    avg_results[size]["cpu_time"] = sum(cpu_times) / len(cpu_times)

    return avg_results

results = avg_results_logreg(X, y, size_increment=100, n=10)

results

{100: {'accuracy': 0.7348000000000001, 'cpu_time': 0.011179752399999642},
 200: {'accuracy': 0.7464000000000001, 'cpu_time': 0.012454438400000356},
 300: {'accuracy': 0.7623, 'cpu_time': 0.015037584000000237},
 400: {'accuracy': 0.7713, 'cpu_time': 0.015927863699998566},
 500: {'accuracy': 0.7729999999999999, 'cpu_time': 0.015920925300000073},
 600: {'accuracy': 0.7773, 'cpu_time': 0.017762035199999103},
 700: {'accuracy': 0.7785, 'cpu_time': 0.018243005299999736},
 800: {'accuracy': 0.7810999999999999, 'cpu_time': 0.019502057399999728},
 900: {'accuracy': 0.7854000000000001, 'cpu_time': 0.01992974199999864},
1000: {'accuracy': 0.7859999999999999, 'cpu_time': 0.02242733419999965}}
```

2.1. Plot the mean accuracy as a function of N .

```

# Parse results into lists for plotting
sizes = []
accuracies = []
for size, result in results.items():
    sizes.append(size)
    accuracies.append(result["accuracy"])

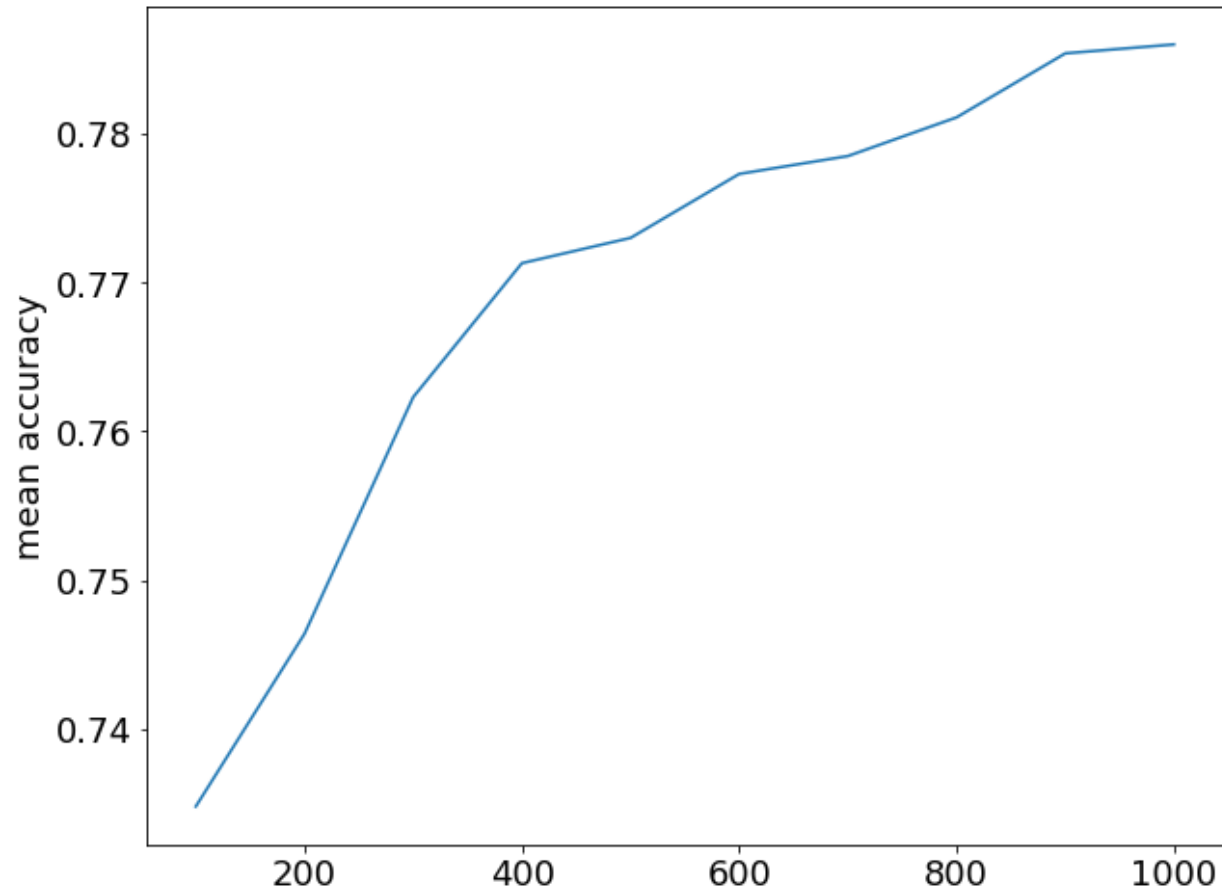
```

```

# Plot accuracy against N
plt.xlabel("N")
plt.ylabel("mean accuracy")
plt.plot(sizes, accuracies)

```

```
[<matplotlib.lines.Line2D at 0x7f06c1569790>]
```



2.2. Plot the cpu-time as a function of N .

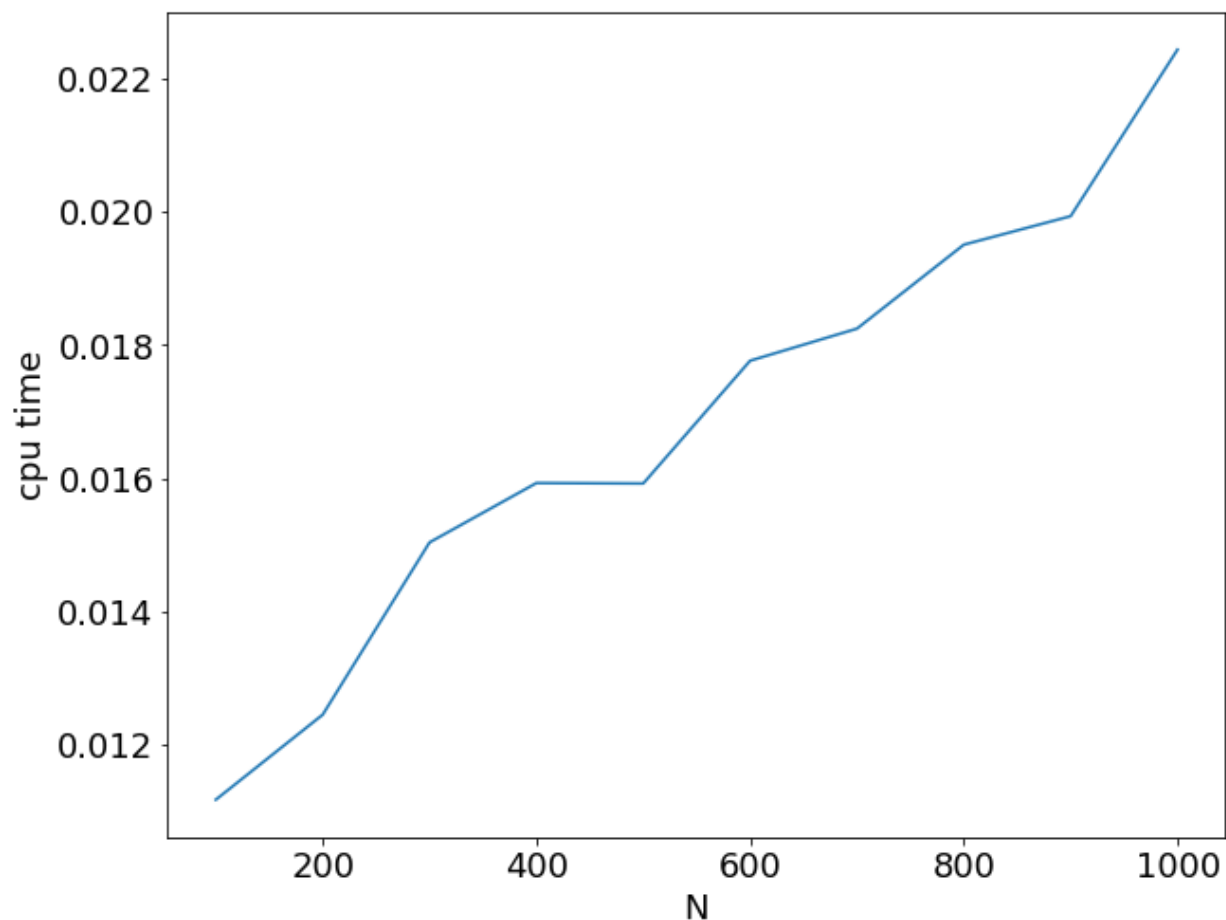
```
# Parse results into lists for plotting
sizes = []
times = []
for size,result in results.items():
    sizes.append(size)
    times.append(result["cpu_time"])
```

```
# Plot mean cpu time against N
plt.xlabel("N")
plt.ylabel("cpu time")
```



```
plt.plot(sizes, times)
```

```
[<matplotlib.lines.Line2D at 0x7f06c14bb690>]
```



2.3. Explain in detail the behaviour of both curves and relate this behavior with what you would expect from theory.

We ran logistic regression on subsets of increasing sizes where each subset contains 100 more examples than the previous. We repeated this operation 50 times and averaged the results.

The first curve describes the average accuracy as a function of the number of training samples. As we can observe, the accuracy is higher when the model is trained on larger subsets. With a subset of 100 examples we reach around 73% accuracy, while with the whole

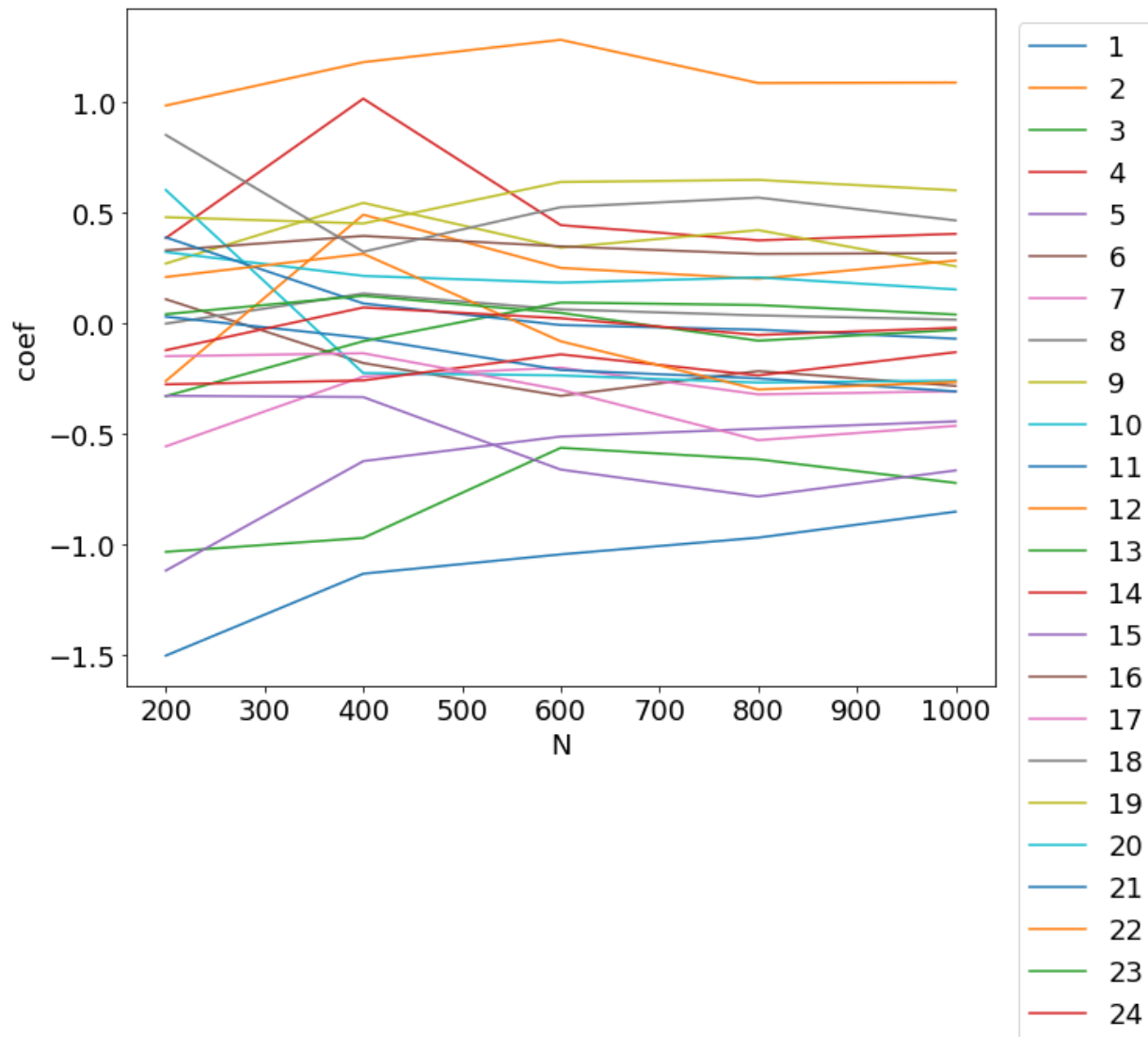
dataset we obtain around 78.6%. This improvement is expected, and due to the fact that the model is able to better fit weights when more data is provided. With better weights, predictions can be more accurate.

The second curve displays cpu time for model fitting as a function of the number of training samples. The plot shows a clear increase of cpu time required by the model for larger subsets. The growth is steep between sets of 100 and 400, after which we see a flat line up to subsets of around 500 examples, but this can be disregarded as most processes are subject to a certain degree of randomness in their scheduling. The growth then proceeds at a similar rate as before. As expected, more numerous calculations are required for subsets with more data points. Because of this, the cpu time is expected to be higher for larger datasets, and the behaviour described by our plot is consistent with this.

2.4. Explore how the learned weights change as a function of N . Can you find an interpretation for the learned weights?

```
# Run logreg once
results = logreg(X, y, size_increment=200)
# Parse results into lists for plotting
sizes = []
coefs = []
for size, result in results.items():
    sizes.append(size)
    coefs.append(result["coef"][0])
```

```
# Plot coefs against N
plot = plt.plot(sizes, coefs)
labels = range(1,25)
plt.legend(plot, labels, loc="upper right", bbox_to_anchor=(1.2, 1))
plt.xlabel("N")
plt.ylabel("coef")
plt.show()
```



Considering that there are numerous features, we increased the size step for each subset from 100 to 200, and ran logistic regression on 5 subsets instead of 10, so as to reduce the amount of points needed for plotting.

We can firstly observe that weights lie within a $[-1.5, 1.1]$ range. Moreover, we can see that for the first model 13 out of the 24 feature weights are positive while the rest are negative. For the last model, 11 feature weights are positive, and the other half are negative.

Out of the positive feature weights in the first model, weight 2 was the highest and remained high with increasing N s. 11 and 16 became negative for the last model.

Out of the negative weights 7 and 15 decreased, while the rest all increased.

Unfortunately, details of how the numeric dataset was obtained are missing, so it is quite hard to interpret these weights because it is unclear what the feature they are associated with describe. However, in general the value of a weight represents the numeric change in the probability predicted by our model corresponding with a 1-unit increase wrt the feature associated, when all other features remain the same. For instance, let's assume that weight 19 is associated with `AGE`. Then, we have that this weight is equal to around .3 for the last model. Under this assumption, a for every year added to the age of the person, there is a 30% increase in the probability of predicting "good risk", keeping all other features fixed.