# Assignment 1: Linear Regression & KNN

Ash

2018/12/06

## Moore's Law

### Extract the date and base speed from the data

```
1  #Loading the data
2  data = pd.read_csv('/Users/ash/Downloads/benchmarks.csv')
3
4  #Extracting the date
5  data["cpu"], data["date"], data["somestuff"] = data['testID'].str.split('-'
       , 2).str
6  data["date"] = pd.to_datetime(data["date"], yearfirst=True)
7  data.dropna(subset=['date'], inplace=True)
```

Since there are some entries without date to extract from and there are only a few of them (578 missing values for some cpus in the 90s, compared to the total of 136995 entries - only 0.4%) I decided to drop these entries from analysis.

### Semi-log plots

There are 73 benchmarks in total, some of which were only used for a specific period. Just to get a feel of different benchmarks, we plot a few of them in figure 1 and figure 2 in $log_2$ plots.

```
1  uniques = data['benchName'].unique()
2  m = 4
3  n = 1
4  nplot = m*n
5  count = 1
6  for _ in uniques[:nplot]:
7      many_cpus = data[data['benchName'] == _]
8      x = np.array(many_cpus['date'], dtype='datetime64[us]')
9      y = np.log2(many_cpus['base'])
10     plt.subplot(m,n,count)
11     plt.scatter(x, y, s=3, label=_)
12     plt.legend(loc=2)
13     plt.xlabel('Time')
14     plt.ylabel('$log_{2}$ of base speed')
15     count += 1
16 plt.show()
17 print(data.groupby('benchName').base.nunique().sort_values())
```
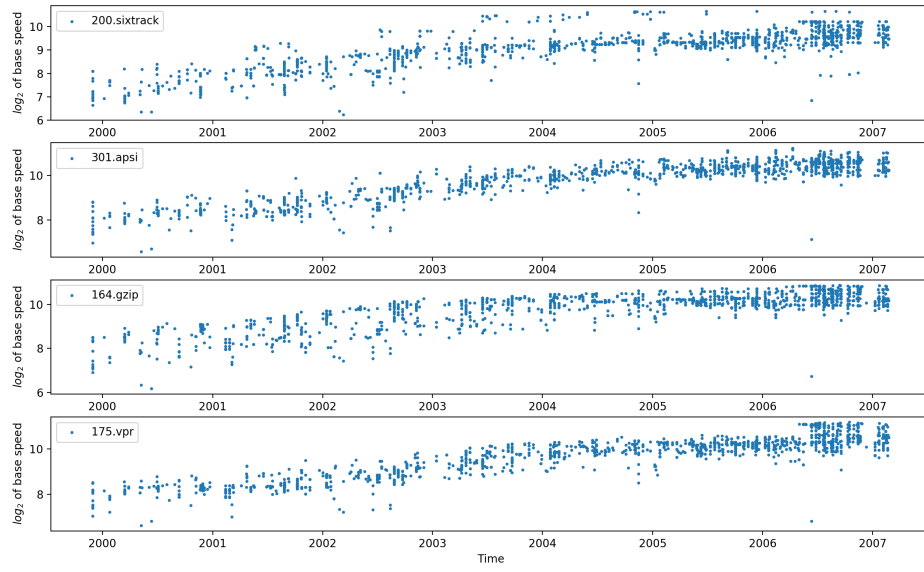
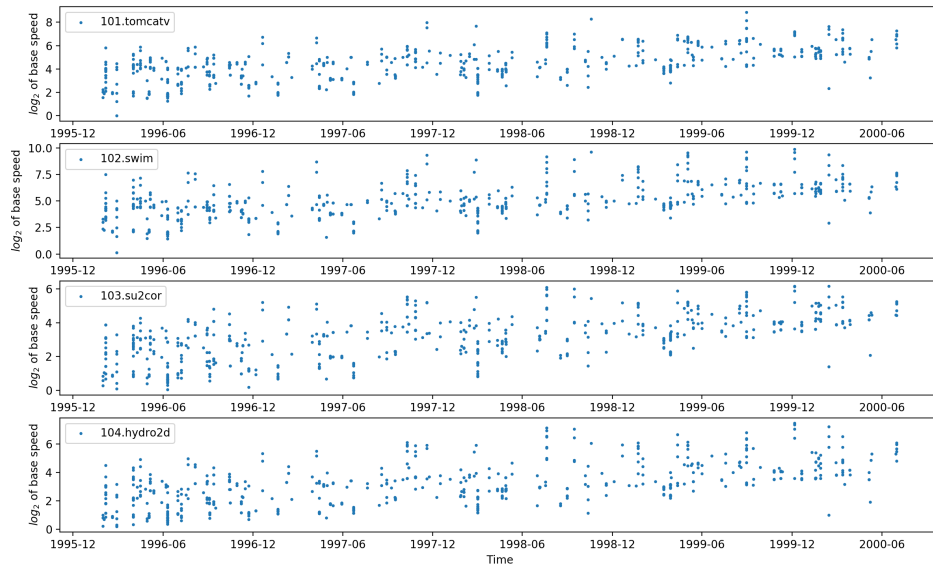Figure 1: Semi-log of base speed for some benchmarks - 1



Figure 2: Semi-log of base speed for some benchmarks - 2

Also plotting the semi-log graph for "178.galgel" in figure 3. This benchmark was chosen because it has the most base speed data.
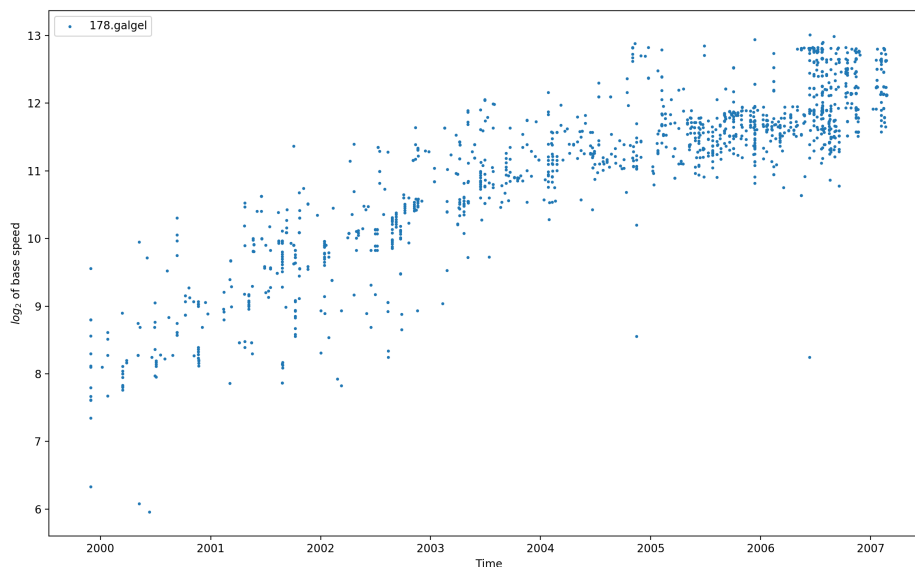


Figure 3: Semi-log of base speed for 178.galgel

**Fitting linear models**

To fit a linear model for the log of base speed of each year (since we are examining Moore's law) we need to average the data for each year because ordinary least square will minimize the mean squared error between a linear regression form relating the year to the log of base speed:

$$base = \theta_1 year + \theta_0 \qquad (or\ \ y = \theta_1 x + \theta_0)$$

which means minimizing:

$$error = (y_{truth} - y_{pred})^2 = (y_{truth} - \theta_1 x - \theta_0)^2$$

Linear regression is convenient because it has a closed form solution for its coefficients, and it also has a nice statistical interpretation: if we assume that the likelihood function (distribution of $\theta$ over the given data) and the prior (distribution of $y$ and $x$ themselves) are normally distributed, and all data points are identically and independently distributed (i.i.d.), we can consider the linear regression problem as updating the prior to find the posteriors for the whole dataset, which returns the maximum likelihood result

3

of the parameters $\theta$s).

Fitting a linear regression by using scikitlearn for the averages of 178.galgel:

```
1  many_cpus = data [ data [ 'benchName' ] == '178.galgel' ]
2  many_cpus [ 'year' ] = many_cpus [ 'date' ] . dt . year
3  x = [ _ for _ in range ( many_cpus [ 'year' ] . min ( ) , many_cpus [ 'year' ] . max ( ) ) ]
4  y = [ ]
5  for _ in range ( many_cpus [ 'year' ] . min ( ) , many_cpus [ 'year' ] . max ( ) ) :
6      dumb = many_cpus [ many_cpus [ 'year' ] == _ ]
7      y . append ( dumb [ 'base' ] . median ( ) )
8  x = np . array ( x ) . reshape ( -1 , 1 )
9  y = np . log2 ( y )
10
11 model = linear_model . LinearRegression ( )
12 model . fit ( x , y )
13 print ( model . coef_ )
14 print ( model . intercept_ )
15 pred = model . predict ( x )
16 plt . scatter ( x , y )
17 plt . plot ( x , pred , color='red' )
18 plt . xlabel ( 'Time' )
19 plt . ylabel ( '$log_{2}$ of base speed' )
20 plt . show ( )
```
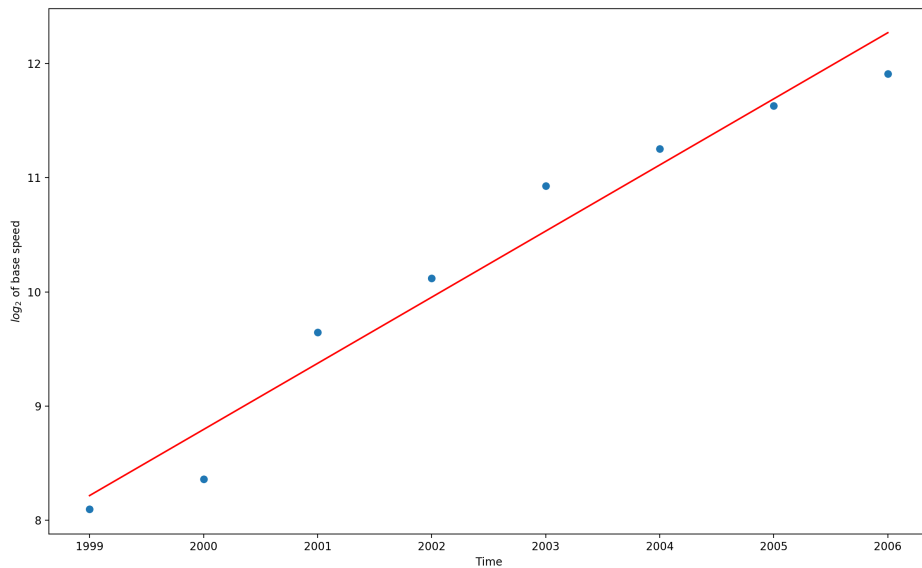


Figure 4: Semi-log of base speed for 178.galgel with best fit line ($\theta_1 = 0.5787779$)

We can also do linear regression models for all the available benchmarks like in figure 5. Eye-balling the graph we can see that generally the slope of the best fit lines are

4

similar, but by some more recent benchmarks the slope is steeper suggesting a more rapid growth of computational power.
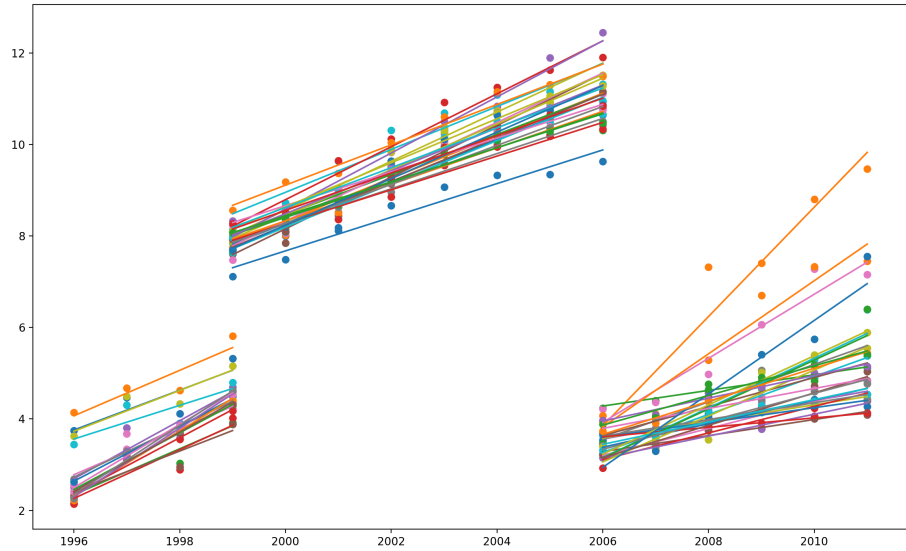


Figure 5: Semi-log of base speed for all benchmarks with best fit lines

## Moore's law

Assume at year 0 the computational power (expressed through base speed in this case) is $a$. If Moore's law is right, then the base speed will double every year, which means:
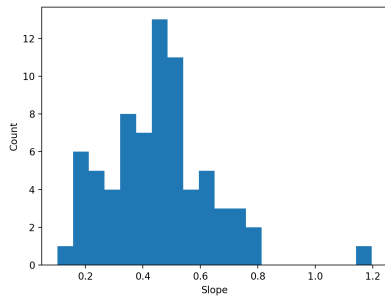
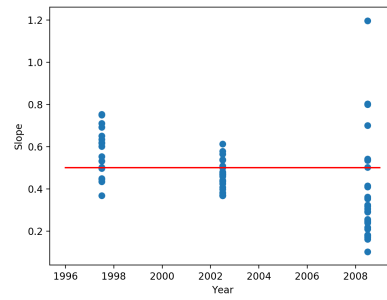$$base_{year_n} = 2^{\frac{n}{2}} base_{year_0} = 2^{\frac{n}{2}} a$$

or

$$log_2(base_{year_n}) = log_2\left(2^{\frac{n}{2}}\right) + log_2(a) = 0.5n + log_2 a$$

so the coefficient (slope) of the best fit line for the semi-log plot will be roughly 0.5 if it obeys Moore's law. For the 178.galgel benchmark we can see that the slope is 0.58, which means it grew faster that Moore predicted. However this is just one benchmark and if we want a closer picture we need to investigate further. If we fit a linear regression models for all 73 benchmarks and extract their slopes:

```
1 def test_moore(cpu, benchmark):
2     many_cpus = cpu[cpu['benchName'] == benchmark]
3     many_cpus['year'] = many_cpus['date'].dt.year
```

(a) Histogram of all benchmark slopes for best fit lines

(b) Graph of all benchmark slopes for best fit lines over year

Figure 6: Slopes of all benchmark

```
4       x = [_ for _ in range(many_cpus['year'].min(), many_cpus['year'].max())
        ]
5       y = []
6       for _ in range(many_cpus['year'].min(), many_cpus['year'].max()):
7           dumb = many_cpus[many_cpus['year'] == _]
8           y.append(dumb['base'].median())
9       x = np.array(x).reshape(-1, 1)
10      y = np.log2(y)
11      model = linear_model.LinearRegression()
12      model.fit(x, y)
13      return model.coef_
14
15
16 uniques = data['benchName'].unique()
17 coef_ = []
18 for _ in uniques:
19     coef_.append(test_moore(data, _))
20
21 print(np.mean(coef_))
22 print(np.median(coef_))
23 plt.hist(np.array(coef_), bins=20)
24 plt.show()
```

After running the above script we will find that the mean of the benchmark growth rates is about 0.4565 and the median is about 0.4493, which is quite close to Moore's law (in figure 6a the histogram of the slopes also concurs). However if we plot the slopes by its time of being used (figure 6b) we will see that in 2008 the slops are more spread out comparing to the other two time marks, suggesting that benchmarks used in 2008 is more inconsistent in terms of its growth of base speed over time, and fewer benchmarks are above the 0.5 red line.

6

## MNIST Digits

### Load MNIST and plot examples

Load the abridged version of the MNIST dataset by scikitlearn then plot some examples (presented in figure 7):

```python
import numpy as np
import math
import matplotlib.pyplot as plt
from sklearn.datasets import load_digits
from sklearn.neighbors import KNeighborsClassifier
import matplotlib

#Plotting function
def plot_digits(digits, m, n):
    ndigits = m*n
    count = 1
    for _ in range(ndigits):
        plt.subplot(m,n,count)
        im = digits.images[_]
        plt.imshow(im, cmap=matplotlib.cm.binary,
                    interpolation="nearest")
        plt.axis("off")
        plt.title(str(digits.target[_]))
        count += 1
    plt.show()

#Loading and plotting some examples from dataset
mnist = load_digits()
plot_digits(mnist, 5, 3)
```

### KNN for two classes and error rate

KNN is a non-parametric machine learning algorithm that is used for classification of multiple classes by measuring distances between a new instances and its "neighbors" in the high-dimensional feature space and classify the new instances by a majority vote (or a weighted vote) of the classes of its neighbors. Probabilistically, KNN is assigning each training instance a Gaussian distribution of variance $\sigma^2$ in the feature space then model each class by a mixture of Gaussians. For a new instance, the probability of it being in each class is considered the probability value of those mixture of Gaussians in the specific place that new instance is located in the feature space. As $\sigma^2$ approach 0, only the nearest neighbor of the new instance has an effect on its classification (since only with the nearest neighbor class the new instance has any probability mass because other Gaussians variance shrink arbitrarily close to 0) but if $\sigma^2 \neq 0$ then other neighbors has some effects on the new instance classification. Using KNeighborsClassifier class in scikitlearn we build the KNN model for the small MNIST dataset with a binary classification of 2 or 5:
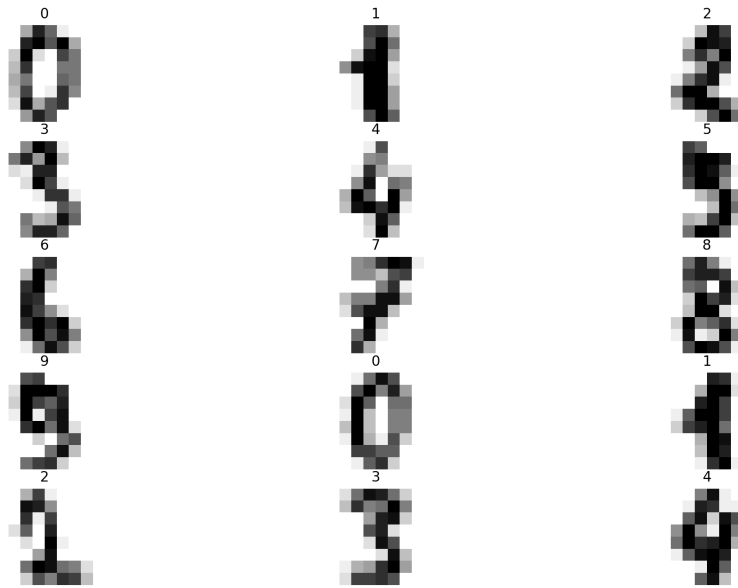
```python
#Loading Model
```

Figure 7: Small MNIST examples

```
2  k = 10
3  weights = 'distance'
4  knn = KNeighborsClassifier(k, weights=weights)
5
6  # #Choosing two numbers to build model
7  first = 2
8  second = 5
9
10 X_train_1st = mnist.data[mnist.target == first]
11 y_train_1st = mnist.target[mnist.target == first]
12 X_train_2nd = mnist.data[mnist.target == second]
13 y_train_2nd = mnist.target[mnist.target == second]
14 X = np.concatenate((X_train_1st/255, X_train_2nd/255), axis=0)
15 y = np.concatenate((y_train_1st, y_train_2nd), axis=0)
16
17 #Shuffling to prevent any artificial sequential order
18 shuffling = np.arange(X.shape[0])
19 np.random.shuffle(shuffling)
20 X = X[shuffling]
21 y = y[shuffling]
22
23 #Train/test split
24 X_train = X[:math.floor(0.8*X.shape[0]),:]
25 X_test = X[math.floor(0.8*X.shape[0]):,:]
26 y_train = y[:math.floor(0.8*X.shape[0])]
```

```
27  y_test = y[math.floor(0.8*X.shape[0]):]
28
29
30  knn.fit(X_train,y_train)
31  pred = knn.predict(X_test)
32  acc = [pred == y_test]
33  print('Accuracy:', sum(acc[0])/y_test.shape[0])
```

In this model we get an accuracy of 100% because it a relatively small dataset and we have a high-dimensional feature space ($8 \times 8 = 64$ features). Let us test the KNN model with multi-class classification:

```
1   #Shuffling to prevent any artificial sequential order
2   X = mnist.data
3   y = mnist.target
4   shuffling = np.arange(X.shape[0])
5   np.random.shuffle(shuffling)
6   X = X[shuffling]
7   y = y[shuffling]
8
9   #Train/test split
10  X_train = X[:math.floor(0.8*X.shape[0]),:]
11  X_test = X[math.floor(0.8*X.shape[0]):,:]
12  y_train = y[:math.floor(0.8*X.shape[0])]
13  y_test = y[math.floor(0.8*X.shape[0]):]
14
15
16  knn.fit(X_train,y_train)
17  pred = knn.predict(X_test)
18  acc = [pred == y_test]
19  print('Accuracy:', sum(acc[0])/y_test.shape[0])
```

The accuracy for 10 classes is around 98.33% which is pretty impressive for such a simple model. Still, this is a smaller version of MNIST, and in practice KNN is rarely used with dataset that has either a large number of features or many instances in the training set because KNN needs to store the whole training set to make inference, which is very memory-costly.

### Full MNIST

Loading the full MNIST dataset (60,000 training examples of 28x28 images for 10 classes) with keras, we build the same KNN model for binary classification of 1 or 3:

```
1   from keras.datasets import mnist
2   (x_train, y_train), (x_test, y_test) = mnist.load_data()
3
4   #Loading Model
5   k = 10
6   weights = 'distance'
7   knn = KNeighborsClassifier(k, weights=weights)
8
9   # #Choosing two numbers to build model
```

```
10  first = 1
11  second = 3
12
13  X_train_1st = x_train [ y_train == first ]
14  y_train_1st = y_train [ y_train == first ]
15  X_train_2nd = x_train [ y_train == second ]
16  y_train_2nd = y_train [ y_train == second ]
17  X = np.concatenate (( X_train_1st /255 , X_train_2nd /255) , axis=0)
18  y = np.concatenate (( y_train_1st , y_train_2nd ) , axis=0)
19
20
21  X_test_1st = x_test [ y_test == first ]
22  y_test_1st = y_test [ y_test == first ]
23  X_test_2nd = x_test [ y_test == second ]
24  y_test_2nd = y_test [ y_test == second ]
25  X_test = np.concatenate (( X_test_1st /255 , X_test_2nd /255) , axis=0)
26  Y_test = np.concatenate (( y_test_1st , y_test_2nd ) , axis=0)
27
28  #Shuffling to prevent any artificial sequential order
29
30  shuffling = np.arange(X.shape [0])
31  np.random.shuffle(shuffling)
32  X = X[shuffling]
33  y = y[shuffling]
34
35
36  knn.fit (X.reshape (X.shape [0] , 28*28) ,y)
37  pred = knn.predict ( X_test.reshape ( X_test.shape [0] , 28*28))
38  acc = [pred == Y_test]
39  print ('Accuracy: ', sum(acc [0])/Y_test.shape [0])
```

It took around 4 minutes for my machine to build the KNN model and make predictions and achieved an impressive 99.94% accuracy on the test set, but again this feature space has 28x28 dimensions so it's possible that in this high-dimensional space these two digits is just separated really distinctively. To test this hypothesis we will use principle component analysis to reduce the dimension of the original to 2 and 3 and plot them to see if they are really separated in these spaces:

```
1  from sklearn.decomposition import PCA
2  pca = PCA( n_components=2)
3
4  new_X = pca.fit_transform ( X_test.reshape ( X_test.shape [0] , 28*28))
5  print ( pca.explained_variance_ratio_ )
6  X_2 = new_X [ Y_test == 1]
7  X_5 = new_X [ Y_test == 3]
8
9  plt.scatter ( X_2 [: ,0] , X_2 [: ,1] , color='red ')
10  plt.scatter ( X_5 [: ,0] , X_5 [: ,1] , color='blue ')
11  plt.show ()
```

In 2 dimensions, even though the explained variance is only 21% and 12% for the 1st and 2nd dimension respectively, we can see in figure 8 that they are quite separated. In 3 dimensions, the 3rd dimension only adds a 5% explained variance but we can see

that they look even more separated. It is very likely that in high dimensional space (where 62% of the remaining variance is located) they are clustered in completely different groups, and therefore the KNN model was able to predict them with such high accuracy.

```python
pca = PCA(n_components=3)

new_X = pca.fit_transform(X_test.reshape(X_test.shape[0], 28*28))
print(pca.explained_variance_ratio_)
X_2 = new_X[Y_test == 1]
X_5 = new_X[Y_test == 3]

ax = plt.axes(projection='3d')
ax.scatter3D(X_2[:,0], X_2[:,1], X_2[:,2], 'red')
ax.scatter3D(X_5[:,0], X_5[:,1], X_5[:,2], 'red')
plt.show()
```
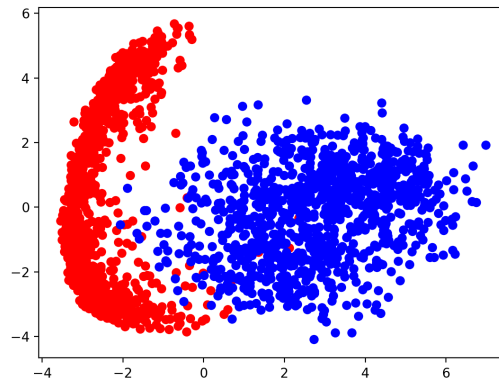


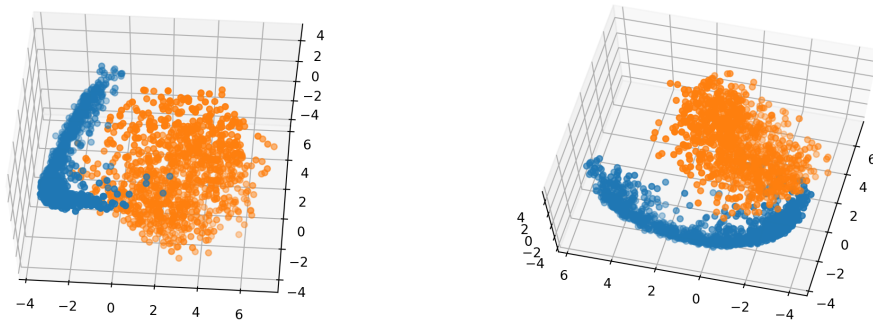Figure 8: 1s and 3s of the original MNIST in 2 dimensions, reduced by PCA

Figure 9: 1s and 3s of the original MNIST in 3 dimensions, reduced by PCA