

# HOMEWORK 1

CSC311 FALL 2019

Nicole Xin Yue Wang

UtorID: wangnic8

#1004235339

**1. Nearest Neighbors and the Curse of Dimensionality - 30 pts.** In this question, you will verify the claim from lecture that “most” points in a high-dimensional space are far away from each other, and also approximately the same distance.

- (a) [10 pts] First, consider two independent univariate random variables  $X$  and  $Y$  sampled uniformly from the unit interval  $[0, 1]$ . Determine the expectation and variance of the random variable  $Z$ , defined as the squared distance  $Z = (X - Y)^2$ .

$$\begin{aligned}
 \text{SOLN} \quad \mathbb{E}[Z] &= \mathbb{E}[(X - Y)^2] \\
 &= \mathbb{E}[X^2 + Y^2 - 2XY] \\
 &= \mathbb{E}[X^2] + \mathbb{E}[Y^2] - 2\mathbb{E}[XY] \\
 &= \mathbb{E}[X^2] + \mathbb{E}[Y^2] - 2\mathbb{E}[X]\mathbb{E}[Y] \quad \# X, Y \text{ are independent} \\
 \because X \sim \text{Uniform}(0, 1), Y \sim \text{Uniform}(0, 1) \\
 \therefore f(x) = f(y) = \frac{1}{1-0} = 1 \\
 &\quad \downarrow \text{pdf} \downarrow \\
 \Rightarrow \mathbb{E}[X^2] &= \int_0^1 x^2 f(x) dx \quad \mathbb{E}[Y^2] = \int_0^1 y^2 f(y) dy \\
 &= \frac{1}{3}x^3 \Big|_0^1 \quad &= \frac{1}{3}y^3 \Big|_0^1 \\
 &= \frac{1}{3} \\
 \mathbb{E}[XY] &= \int_0^1 x f(x) dx \quad \mathbb{E}[Y] = \int_0^1 y f(y) dy \\
 &= \frac{1}{2}x^2 \Big|_0^1 \quad &= \frac{1}{2}y^2 \Big|_0^1 \\
 &= \frac{1}{2} \\
 \Rightarrow \mathbb{E}[Z] &= \frac{1}{3} + \frac{1}{3} - 2\left(\frac{1}{2}\right)\left(\frac{1}{2}\right) \\
 &= \frac{2}{3} - \frac{1}{2} = \frac{1}{6}
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(Z) &= \mathbb{E}[Z^2] - \mathbb{E}[Z]^2 \\
 &= \mathbb{E}[(X - Y)^4] - \left(\frac{1}{6}\right)^2 \\
 &= \mathbb{E}[X^4] + \mathbb{E}[Y^4] - 4\mathbb{E}[X^3]\mathbb{E}[Y] - 4\mathbb{E}[X]\mathbb{E}[Y^3] + 6\mathbb{E}[X^2]\mathbb{E}[Y^2] - \frac{1}{36} \\
 &= \int_0^1 x^4 dx + \int_0^1 y^4 dy - 4 \int_0^1 x^3 dx \left(\frac{1}{2}\right) - 4 \left(\frac{1}{2}\right) \int_0^1 y^3 dy + 6 \left(\frac{1}{3}\right) \left(\frac{1}{3}\right) - \frac{1}{36} \\
 &= \frac{1}{5} + \frac{1}{5} - 4\left(\frac{1}{4}\right)\left(\frac{1}{2}\right) - 4\left(\frac{1}{2}\right)\left(\frac{1}{4}\right) + \frac{2}{3} - \frac{1}{36} \\
 &= \frac{2}{5} - 1 + \frac{2}{3} - \frac{1}{36} \\
 &= \frac{72 - 180 + 120 - 5}{180} \\
 &= \frac{7}{180}
 \end{aligned}$$

<b>ANS:</b> $\mathbb{E}[Z] = \frac{1}{6}$ , $\text{Var}(Z) = \frac{7}{180}$
-----------------------------------------------------------------------------

- (b) [10 pts] Now suppose we sample two points independently from a unit cube in  $d$  dimensions. Observe that each coordinate is sampled independently from  $[0, 1]$ , i.e. we can view this as sampling random variables  $X_1, \dots, X_d, Y_1, \dots, Y_d$  independently from  $[0, 1]$ . The squared Euclidean distance can be written as  $R = Z_1 + \dots + Z_d$ , where  $Z_i = (X_i - Y_i)^2$ . Using the properties of expectation and variance, determine  $\mathbb{E}[R]$  and  $\text{Var}[R]$ . You may give your answer in terms of the dimension  $d$ , and  $\mathbb{E}[Z]$  and  $\text{Var}[Z]$  (the answers from part (a)).

$$\begin{aligned}\mathbb{E}[R] &= \mathbb{E}[z_1 + \dots + z_d] = \mathbb{E}[z_1] + \dots + \mathbb{E}[z_d] \quad \# \text{ by properties of expectation} \\ &= \underbrace{\frac{1}{6} + \dots + \frac{1}{6}}_{d \text{ times}} \\ &= \frac{d}{6}\end{aligned}$$

$$\begin{aligned}\text{Var}(R) &= \text{Var}(z_1 + \dots + z_d) = \text{Var}(z_1) + \dots + \text{Var}(z_d) \quad \# \text{ by properties of variance} \\ &= \frac{7d}{180} \quad \Rightarrow \text{Var}(X+Y) = \text{Var}(X) + 2\text{Cov}(X,Y) + \text{Var}(Y) \\ &\quad \text{while } X \text{ and } Y \text{ are independent} \Rightarrow \text{Cov}(X,Y) = 0\end{aligned}$$

**ANS:**  $\mathbb{E}[R] = \frac{d}{6}$ ,  $\text{Var}(R) = \frac{7d}{180}$

- (c) [10 pts] Based on your answer to part (b), compare the mean and standard deviation of  $R$  to the maximum possible squared Euclidean distance (i.e. the distance between opposite corners of the cube). Why does this support the claim that in high dimensions, “most points are far away, and approximately the same distance”?

$$\begin{aligned}\max(R) &= \max\left(\sum_{i=1}^d z_i\right) = \sum_{i=1}^d \max(z_i) \\ &= \sum_{i=1}^d \max((x_i - y_i)^2) \\ &= \sum_{i=1}^d (1) \quad \# X \text{ and } Y \text{ are sampled from unit interval } [0,1], \\ &\quad \text{the maximum difference will be } 1-0=1 \\ &= d\end{aligned}$$

In high-dimensions, since both mean and maximum possible values are multiples of  $d$ , the higher dimensions it get, the further the points get.

$\Rightarrow \lim_{d \rightarrow \infty} \mathbb{E}[R] = \lim_{d \rightarrow \infty} \frac{d}{6} \rightarrow \infty \Rightarrow$  therefore points could be very far away in high dimensions

If we define distance between “most” points to be within  $3\sigma$  away from mean, then we have:

$$\begin{aligned}\sigma &= \sqrt{\text{Var}(R)} = \sqrt{\frac{7d}{180}} \\ &\Rightarrow 3\sigma = 3\sqrt{\frac{7d}{180}} \\ &\quad = 0.59\sqrt{d} \Rightarrow \text{very small difference}\end{aligned}$$

Therefore, in high-dimensions, “most” points are far away and approximately the same distance.  $\blacksquare$

2. (a) [6 pts] Write a function `load_data` which loads the data, preprocesses it using a vectorizer ([http://scikit-learn.org/stable/modules/classes.html#module-sklearn.feature\\_extraction.text](http://scikit-learn.org/stable/modules/classes.html#module-sklearn.feature_extraction.text), we suggest you use `CountVectorizer` as it is the simplest in nature), and splits the entire dataset randomly into 70% training, 15% validation, and 15% test examples.

```

1 def load_data():
2     """ Load headlines from clean_real.txt and clean_fake.txt and seperate them
3     into training set, validation set, and testing set. """
4     # read clean_fake.txt and store each headline as an element into an array
5     fake_file = open("clean_fake.txt", "r")
6     fake_arr = fake_file.read().split('\n')
7     fake_file.close()
8     # read clean_real.txt and store each headline as an element into an array
9     real_file = open("clean_real.txt", "r")
10    real_arr = real_file.read().split('\n')
11    real_file.close()
12    #combine the two arrays into one sample array
13    sample_arr = fake_arr + real_arr
14    # target array, 0 means corresponding headline is fake, 1 means real
15    target_arr = [0] * len(fake_arr) + [1] * len(real_arr)
16    # split into three sets
17    vectorizer = CountVectorizer()
18    sample = vectorizer.fit_transform(sample_arr)
19    x_train, x_test, y_train, y_test = train_test_split(
20        sample.toarray(), target_arr, test_size=0.3)
21    x_val, x_test, y_val, y_test = train_test_split(
22        x_test, y_test, test_size=0.5)
23    #return the dictionary
24    return {"train":x_train, "train_target": y_train, "val":x_val,
25            "val_target":y_val, "test":x_test, "test_target": y_test,
26            "feature_names":vectorizer.get_feature_names()}

```

- (b) [6 pts] Write a function `select_model` which trains the decision tree classifier using at least 5 different values of `max_depth`, as well as two different split criteria (information gain and Gini coefficient), evaluates the performance of each one on the validation set, and prints the resulting accuracies of each model. You should use `DecisionTreeClassifier`, but you should write the validation code yourself. Include the output of this function in your solution.

```

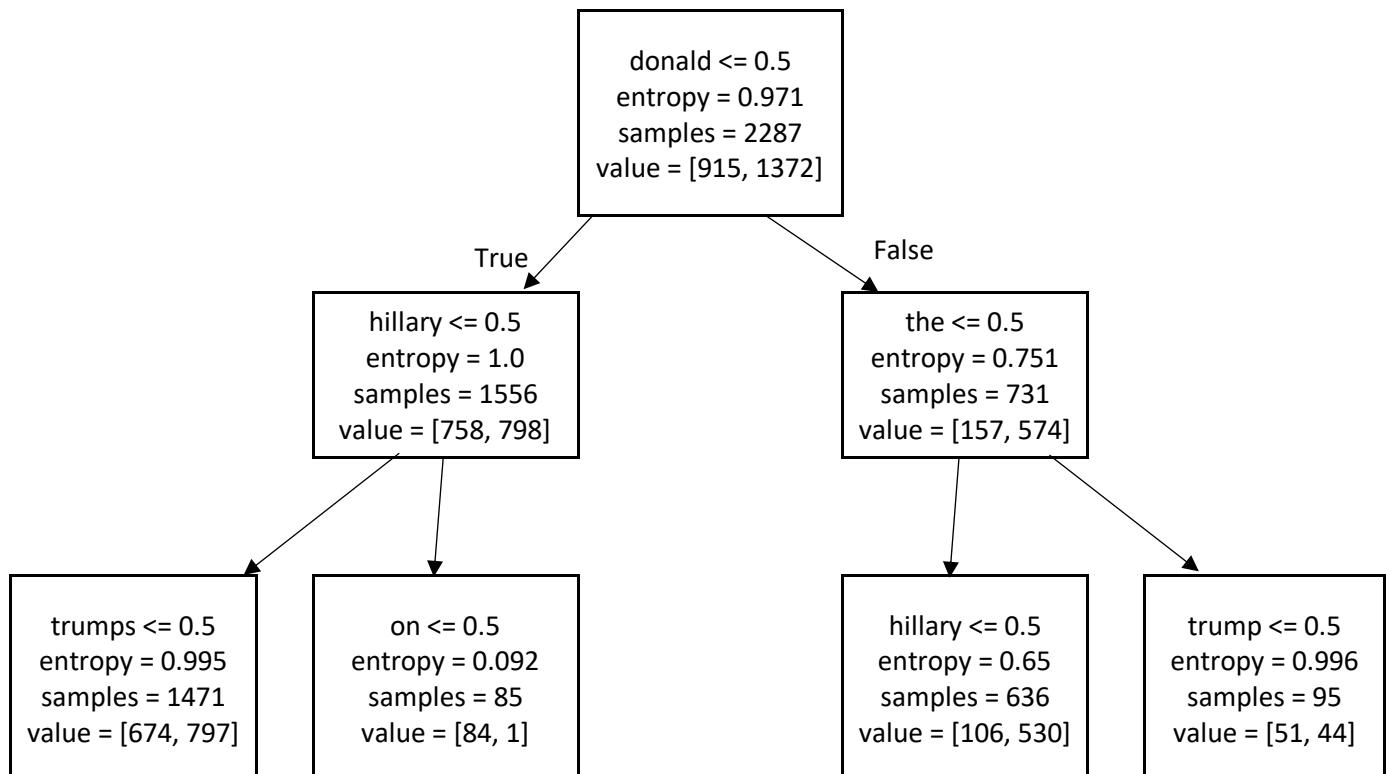
1 def select_model():
2     """ Select the model with the most efficient hyperparameter. """
3     # get data
4     data = load_data()
5     # initialize the variables for storing the model with best accuracy
6     best_acc = 0
7     best_depth = 0
8     best_crit = ""
9     # two criterion to choose
10    for crit in ["gini", "entropy"]:
11        # 5 different types of maximum depth of the tree
12        for i in [4,8,12,16,20]:
13            # train a tree with given criterion and depth
14            clf = DecisionTreeClassifier(criterion=crit, max_depth=i)
15            clf.fit(data["train"], data["train_target"])
16            # accuracies testing on the validation set, store best accuracy
17            acc = clf.score(data["val"], data["val_target"])
18            (best_acc, best_depth, best_crit) = (
19                best_acc, best_depth, best_crit
20                ) if best_acc > acc else (acc, i, crit)
21            print("Score (max_depth=", i, ", criterion=", crit, ") :", acc)
22            print("Best Model(max_depth)=" , best_depth, ", criterion=",
23                  best_crit, "): Score=", best_acc)

```

(b) (cont'd) Output:

```
Score (max_depth= 4 , criterion= gini ) : 0.7081632653061225
Score (max_depth= 8 , criterion= gini ) : 0.7122448979591837
Score (max_depth= 12 , criterion= gini ) : 0.7285714285714285
Score (max_depth= 16 , criterion= gini ) : 0.746938775510204
Score (max_depth= 20 , criterion= gini ) : 0.7489795918367347
Score (max_depth= 4 , criterion= entropy ) : 0.6979591836734694
Score (max_depth= 8 , criterion= entropy ) : 0.7040816326530612
Score (max_depth= 12 , criterion= entropy ) : 0.7326530612244898
Score (max_depth= 16 , criterion= entropy ) : 0.7387755102040816
Score (max_depth= 20 , criterion= entropy ) : 0.753061224489796
Best Model(max_depth= 20 , criterion= entropy ): Score= 0.753061224489796
```

- (c) [6 pts] Now let's stick with the hyperparameters which achieved the highest validation accuracy. Extract and visualize the first two layers of the tree. Your visualization does not have to be an image: it is perfectly fine to display text. It may also be hand-drawn. Include your visualization in your solution pdf.



\* tree is generated by using `sklearn.tree.export_graphviz` to generate a dot file and visualized using <http://www.webgraphviz.com/>

- (d) [12 pts] Write a function `compute_information_gain` which computes the information gain of a split on the training data. That is, compute  $I(Y, x_i)$ , where  $Y$  is the random variable signifying whether the headline is real or fake, and  $x_i$  is the keyword chosen for the split. Report the outputs of this function for the topmost split from the previous part, and for several other keywords.

```

1 def compute_information_gain(keyword):
2     data = load_data()
3     total = len(data["train_target"]) # total number of news
4     # get the index of keyword in the list of features
5     index = data["feature_names"].index(keyword)
6     train = data["train"]
7     real = 0 # number of real news
8     real_key = 0 # number of real news with keyword
9     key = 0 # number of news with keyword
10    for i in range(total):
11        if train[i][index] > 0:
12            key += 1 # count one news with keyword
13            if data["train_target"][i] == 1:
14                real_key += 1 # count one real news with keyword
15            if data["train_target"][i] == 1:
16                real += 1 # count one real news
17        entropy = n_nlogn(real / total) + n_nlogn(
18            (total - real)/ total) # H(Y)
19        p_e_r = real_key / key # real news, with keyword
20        p_e_f = (key-real_key) / key # fake news, with keyword
21        p_ne_r = (real-real_key) / (total-key) # real news, no keyword
22        p_ne_f = (total-real-(key-real_key)) / (total-key) # fake news, no keyword
23        entropy_keyword = (key/total) * (n_nlogn(p_e_r) + n_nlogn(p_e_f)) +
24            (total-key)/total) * (n_nlogn(p_ne_r) + n_nlogn(p_ne_f)) # H(Y|xi)
25        IG = entropy - entropy_keyword # information gain
26        print("Information Gain in splitting at keyword '", keyword, "' is ", IG)
27
28 if __name__ == "__main__":
29     for key in ["donald", "trumps", "hillary", "the", "trump",
30                 "breaking", "prime"]:
31         compute_information_gain(key)

```

Output:

Information Gain in splitting at keyword ' donald ' is 0.051323670040992586 Information Gain in splitting at keyword ' trumps ' is 0.044374383775848236 Information Gain in splitting at keyword ' hillary ' is 0.03859640155331756 Information Gain in splitting at keyword ' the ' is 0.05584727379340548 Information Gain in splitting at keyword ' trump ' is 0.038198158834489626 Information Gain in splitting at keyword ' breaking ' is 0.010519717839329035 Information Gain in splitting at keyword ' prime ' is 3.0293720648333355e-05
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

3. Regression - 40 pts. In this question, you will derive certain properties of linear regression, and experiment with cross validation to tune its regularization parameter.

3.1. Linear regression - 10 pts. Suppose that  $\mathbf{X} \in \mathbb{R}^{n \times m}$  with  $n \geq m$  and  $\mathbf{t} \in \mathbb{R}^n$ , and that  $\mathbf{t}|(\mathbf{X}, \mathbf{w}) \sim \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$ . We know that the maximum likelihood estimate  $\hat{\mathbf{w}}$  of  $\mathbf{w}$  is given by

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}.$$

- (a) Write the log-likelihood implied by the model above, and compute its gradient w.r.t.  $\mathbf{w}$ . By setting it equal to 0, derive the above estimator  $\hat{\mathbf{w}}$ .

Given that  $\mathbf{t}|(\mathbf{X}, \mathbf{w}) \sim \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I})$ , we know that:

$$\begin{aligned} p(t|(\mathbf{X}, \mathbf{w})) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(t_i - \mathbf{x}_i^T \mathbf{w})^2\right\} \\ \text{likelihood} \Rightarrow \log p(t|\mathbf{X}, \mathbf{w}) &= \log \left( \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(t_i - \mathbf{x}_i^T \mathbf{w})^2\right\} \right) \\ &= \sum_{i=1}^n \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(t_i - \mathbf{x}_i^T \mathbf{w})^2\right\} \right) \\ &= \sum_{i=1}^n \left( -\frac{1}{2\sigma^2} (t_i - \mathbf{x}_i^T \mathbf{w})^2 \right) \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \quad \# \text{ by power rule of logarithms} \\ &= -\frac{1}{2\sigma^2} \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \sum_{i=1}^n (t_i - \mathbf{x}_i^T \mathbf{w})^2 \quad \# \text{ by distributive property of summation and} \\ &\quad -\frac{1}{2\sigma^2} \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \text{ is a constant} \\ &= -\frac{1}{2\sigma^2} \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \underbrace{(t - \mathbf{X}\mathbf{w})^T (t - \mathbf{X}\mathbf{w})}_{=(\mathbf{t}^T - \mathbf{w}^T \mathbf{X}^T)} \\ &= -\frac{1}{2\sigma^2} \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) (\mathbf{t}^T \mathbf{t} - \mathbf{t}^T \mathbf{X}\mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{t} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w}) \\ \Rightarrow \nabla [\log p(t|\mathbf{X}, \mathbf{w})] &= \left[ -\frac{1}{2\sigma^2} \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \right] (-\mathbf{t}^T \mathbf{X} - \mathbf{X}^T \mathbf{t} + 2\mathbf{X}^T \mathbf{X}\mathbf{w}) \end{aligned}$$

let  $\nabla = 0$ , then  $-\mathbf{t}^T \mathbf{X} - \mathbf{X}^T \mathbf{t} + 2\mathbf{X}^T \mathbf{X}\mathbf{w} = 0$

$$2\mathbf{X}^T \mathbf{X}\mathbf{w} = \mathbf{t}^T \mathbf{X} + \mathbf{X}^T \mathbf{t}$$

$$\text{since } \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}, \text{ then } \mathbf{t}^T = [t_1, t_2, \dots, t_n]$$

$$\text{since } \mathbf{X} = \begin{bmatrix} \mathbf{x}_{11} & \mathbf{x}_{12} & \dots & \mathbf{x}_{1m} \\ \mathbf{x}_{21} & \mathbf{x}_{22} & \dots & \mathbf{x}_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_{n1} & \mathbf{x}_{n2} & \dots & \mathbf{x}_{nm} \end{bmatrix}, \text{ then } \mathbf{X}^T = \begin{bmatrix} \mathbf{x}_{11} & \mathbf{x}_{21} & \dots & \mathbf{x}_{n1} \\ \mathbf{x}_{12} & \mathbf{x}_{22} & \dots & \mathbf{x}_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_{1m} & \mathbf{x}_{2m} & \dots & \mathbf{x}_{nm} \end{bmatrix}$$

$$\Rightarrow \mathbf{t}^T \mathbf{X} = [t_1 \mathbf{x}_{11} + \dots + t_n \mathbf{x}_{n1}, t_1 \mathbf{x}_{12} + \dots + t_n \mathbf{x}_{n2}, \dots, t_1 \mathbf{x}_{1m} + \dots + t_n \mathbf{x}_{nm}]$$

$$\text{and } \mathbf{X}^T \mathbf{t} = \begin{bmatrix} \mathbf{x}_{11} t_1 + \dots + \mathbf{x}_{n1} t_n \\ \mathbf{x}_{12} t_1 + \dots + \mathbf{x}_{n2} t_n \\ \vdots \\ \mathbf{x}_{1m} t_1 + \dots + \mathbf{x}_{nm} t_n \end{bmatrix} \quad \because \text{both } \mathbf{t}^T \mathbf{X} \text{ and } \mathbf{X}^T \mathbf{t} \text{ are vectors} \\ \therefore \mathbf{t}^T \mathbf{X} = \mathbf{X}^T \mathbf{t}$$

$$\Rightarrow \cancel{2\mathbf{X}^T \mathbf{X}\mathbf{w}} = \cancel{2\mathbf{X}^T \mathbf{t}}$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \quad \blacksquare$$

- (b) Find the distribution of  $\hat{\mathbf{w}}$ , its expectation and covariance matrix. Hint: Read the property of multivariate Gaussian random vectors in preliminaries.pdf on course webpage.

Since  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$ ,  $\mathbf{t}$  is a normal distributed random variable, and  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is a constant matrix, by the property of a multivariate Gaussian random variable,  $\hat{\mathbf{w}}$  is also a multivariate Gaussian random vector.

$$\text{mean: } [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{X}\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t} \\ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

$$\text{covariance: } [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \sigma^2 \mathbf{I} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T]^T = \sigma^2 [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T]^T$$

$$\text{Let } \mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \text{ then } \hat{\mathbf{w}} \sim \mathcal{N}(\mathbf{A}\mathbf{t}, \sigma^2 \mathbf{A}\mathbf{A}^T)$$

3.2. Ridge regression and MAP - 10 pts. Suppose that we have  $t|(\mathbf{X}, \mathbf{w}) \sim \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I})$  and we place a normal prior on  $\mathbf{w}|\mathbf{X}$ , i.e.,  $\mathbf{w} \sim \mathcal{N}(0, \tau^2\mathbf{I})$ . Recall from the first tutorial (also in preliminaries.pdf) that MAP estimate of  $\mathbf{w}$  is given as the maximum of the posterior density

$$\hat{\mathbf{w}}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} \{p(\mathbf{w}|\mathbf{X}, t) \propto p(t|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{X})\}.$$

Here,  $\propto$  notation means *proportional to*, and is used since we dropped the term  $p(t|\mathbf{X})$  in the denominator as it doesn't have  $\mathbf{w}$  in it, thus it doesn't contribute to the maximization problem.

Show that the MAP estimate of  $\mathbf{w}$  given  $(t, \mathbf{X})$  in this context is

$$(3.1) \quad \hat{\mathbf{w}}_{MAP} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top t$$

where  $\lambda = \sigma^2/\tau^2$ .

SOLN

$$\begin{aligned} p(w|\mathbf{X}, t) &\propto p(t|\mathbf{X}, w) p(w|\mathbf{X}) \\ \Rightarrow \log p(w|\mathbf{X}, t) &\propto \log [p(t|\mathbf{X}, w) p(w|\mathbf{X})] \\ &= \log \left( \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2\mathbf{I}}} \exp \left\{ -\frac{1}{2\sigma^2\mathbf{I}} (t_i - \mathbf{x}_i^\top \mathbf{w})^2 \right\} \right) \left( \prod_{i=1}^n \frac{1}{\sqrt{2\pi\tau^2\mathbf{I}}} \exp \left\{ -\frac{1}{2\tau^2\mathbf{I}} (w_i)^2 \right\} \right) \\ &= -\frac{1}{2\sigma^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\sigma^2\mathbf{I}}} \right) (t^\top t - t^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top t + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}) + \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi\tau^2\mathbf{I}}} \exp \left\{ -\frac{1}{2\tau^2\mathbf{I}} (w_i)^2 \right\} \\ &= -\frac{1}{2\sigma^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\sigma^2\mathbf{I}}} \right) (t^\top t - t^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top t + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}) + \sum_{i=1}^n \left[ -\frac{1}{2\tau^2\mathbf{I}} (w_i)^2 \right] \log \left( \frac{1}{\sqrt{2\pi\tau^2\mathbf{I}}} \right) \\ &= -\frac{1}{2\sigma^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\sigma^2\mathbf{I}}} \right) (t^\top t - t^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top t + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}) - \frac{1}{2\tau^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\tau^2\mathbf{I}}} \right) \sum_{i=1}^n w_i^2 \\ &= -\frac{1}{2\sigma^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\sigma^2\mathbf{I}}} \right) (t^\top t - t^\top \mathbf{X} \mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top t + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}) - \frac{1}{2\tau^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\tau^2\mathbf{I}}} \right) \mathbf{w}^\top \mathbf{w} \\ \Rightarrow \nabla &= -\frac{1}{2\sigma^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\sigma^2\mathbf{I}}} \right) (-t^\top \mathbf{X} - \mathbf{X}^\top t + 2\mathbf{X}^\top \mathbf{X} \mathbf{w}) - \frac{1}{2\tau^2\mathbf{I}} \log \left( \frac{1}{\sqrt{2\pi\tau^2\mathbf{I}}} \right) (\mathbf{X}^\top \mathbf{X}) \end{aligned}$$

let  $\nabla = 0$ , then:

$$\begin{aligned} -\frac{1}{2\sigma^2\mathbf{I}} \log \left[ (2\pi\sigma^2\mathbf{I})^{-\frac{1}{2}} \right] (-t^\top \mathbf{X} - \mathbf{X}^\top t + 2\mathbf{X}^\top \mathbf{X} \mathbf{w}) - \frac{1}{2\tau^2\mathbf{I}} \log \left[ (2\pi\tau^2\mathbf{I})^{-\frac{1}{2}} \right] \mathbf{w} &= 0 \\ + \frac{1}{4\sigma^2\mathbf{I}} \log \left( 2\pi\sigma^2\mathbf{I} \right) (-t^\top \mathbf{X} - \mathbf{X}^\top t + 2\mathbf{X}^\top \mathbf{X} \mathbf{w}) &= \frac{-1}{2\tau^2\mathbf{I}} \log \left( 2\pi\tau^2\mathbf{I} \right) \mathbf{w} \\ (-t^\top \mathbf{X} - \mathbf{X}^\top t + 2\mathbf{X}^\top \mathbf{X} \mathbf{w}) &= -\frac{\sigma^2\mathbf{I}}{\tau^2\mathbf{I}} \frac{4 \log (2\pi\sigma^2\mathbf{I})}{2 \log (2\pi\tau^2\mathbf{I})} \mathbf{w} \\ \frac{-2\mathbf{X}^\top t}{\mathbf{w}} + 2\mathbf{X}^\top \mathbf{X} &= -\frac{\sigma^2}{\tau^2} \mathbf{I} \frac{2 \log (2\pi\sigma^2\mathbf{I})}{\log (2\pi\sigma^2\mathbf{I})} \\ \frac{-2\mathbf{X}^\top t}{\mathbf{w}} &= 2\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I} \frac{4 \log (2\pi\sigma^2\mathbf{I})}{\log (2\pi\sigma^2\mathbf{I})} \\ \mathbf{w} &= \left( \mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I} \frac{\log (2\pi\sigma^2\mathbf{I})}{\log (2\pi\sigma^2\mathbf{I})} \right)^{-1} \mathbf{X}^\top t \\ \hat{\mathbf{w}}_{MAP} &= \left( \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top t \quad \# \text{let } \lambda = \frac{\sigma^2}{\tau^2}, \text{ and drop } \frac{\log (2\pi\sigma^2\mathbf{I})}{\log (2\pi\sigma^2\mathbf{I})} \text{ since} \\ &\quad \text{it has very small affect comparing to } \frac{\sigma^2}{\tau^2}. \text{ and we can adjust } \lambda. \quad \blacksquare \end{aligned}$$

### 3.3. Cross validation - 20 pts.

(b) Code for the 6 functions:

```

1 import numpy as np
2 from numpy.linalg import pinv
3
4 def shuffle_data(data):
5     """ Takes in a data that contains a (t,X) pair and returns its randomly
6     permuted version along the samples. """
7     # generate an array that contains a random order of indices for data
8     random_order = np.random.permutation(len(data["X"]))
9     # return de reordered X and t
10    return {"X": data["X"][random_order], "t": data["t"][random_order]}
11

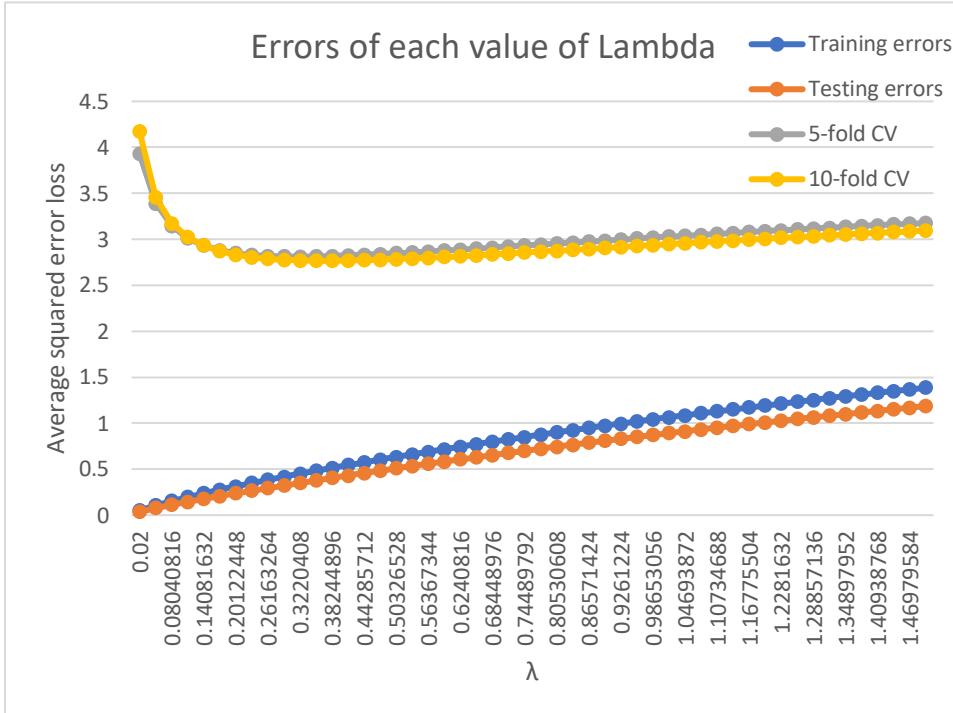
```

```

12 def split_data(data, num_fold, fold):
13     """
14         Takes in data that is a (t,X) pair, folds each row into num_fold equal folds
15         and the foldth will be the validation fold. Returns the validation fold as
16         a 1d array and the rest folds as a 2d array.
17     """
18     # fold the samples and target
19     folded_x = np.array_split(data["X"], num_fold)
20     folded_t = np.array_split(data["t"], num_fold)
21     # pop the fold to be used as validation set
22     data_fold = {"X": folded_x.pop(fold-1), "t": folded_t.pop(fold-1)}
23     # add the rest of the folds into return value
24     data_rest = {"X":np.concatenate(folded_x), "t":np.concatenate(folded_t)}
25     return data_fold, data_rest
26
27
28 def train_model(data, lambd):
29     """
30         Takes a set of data and return the coefficient of ridge regression with
31         penalty level lambd.
32     """
33     x = np.array(data["X"])
34     t = data["t"]
35     x_trans = x.transpose()
36     return np.matmul(pinv(np.matmul(x_trans, x)) +
37                     np.dot(lambd,np.identity(len(x_trans))), np.matmul(x_trans, t))
38
39 def predict(data, model):
40     """ Takes in data and model coefficient and predict the outcome. """
41     return np.matmul(data["X"], model)
42
43 def loss(data, model):
44     """ Takes in the data and model, and return the average squared error loss.
45     """
46     t_xw = data["t"] - predict(data, model)
47     return t_xw.dot(t_xw) / len(data["t"])
48
49
50 def cross_validation(data, num_folds, lambd_seq):
51     """ Takes in the training data, number of folds, and a sequence of  $\lambda$  and
52         return a vector of 50 cross validation errors. """
53     cv_error = []
54     data = shuffle_data(data)
55     for i in range(len(lambd_seq)):
56         lambd = lambd_seq[i]
57         cv_loss_lmd = 0
58         for fold in range(1, num_folds+1):
59             val_cv, train_cv = split_data(data, num_folds, fold)
60             model = train_model(train_cv, lambd)
61             cv_loss_lmd += loss(val_cv, model)
62         cv_error.append(cv_loss_lmd / num_folds)
63     return cv_error
64
65 if __name__ == "__main__":
66     data_train = {"X": np.genfromtxt("data_train_X.csv", delimiter=','),
67                   "t": np.genfromtxt("data_train_y.csv", delimiter=',')}
68     data_test = {"X": np.genfromtxt("data_test_X.csv", delimiter=','),
69                  "t": np.genfromtxt("data_test_y.csv", delimiter=',')}
70     lambd_seq = np.linspace(0.02, 1.5, num=50)
71

```

(d) Ploting



Training errors and testing errors follows a similar pattern while the two cross validations follows a similar patten while more fold makes slightly smaller errors than less folds. The training and testing errors are way smaller than cross validations. However, this is probably because error evaluations on traning and testing datas are on the exact same set of datas used for traning, while cross validations uses some set of sample datas for training, and then error losts are calculated by using a validation set which the model never seen before. If a new set of samples are to be tested on the model trained by the training and testing datas, the errors would increase.