

The code for calculations and generating the tree is in the Jupyter notebook in the codes directory.

Q1-1:

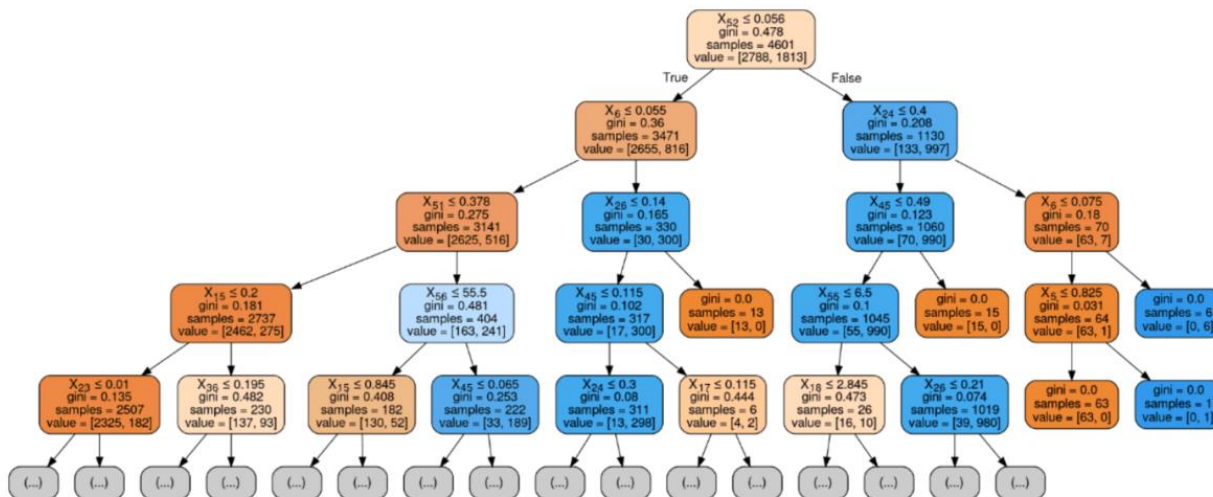
There are 1813 spam emails and 2788 not spam emails.

Q1-2:

The fitted decision tree model is:

```
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
                      max_depth=None, max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=2,
                      min_weight_fraction_leaf=0.0, presort='deprecated',
                      random_state=None, splitter='best')
```

I used two libraries to generate the graphics for the tree. Both were too large to fit here in the report. I used portion of the tree for our report here. In the Jupyter, I have full tree with both libraries.



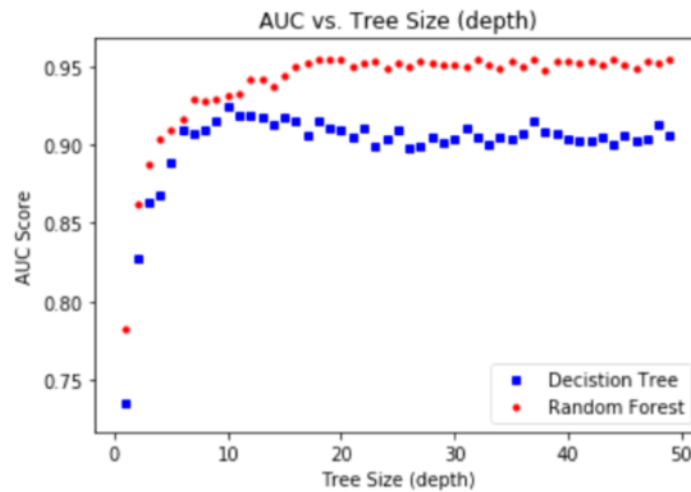
Q1-3:

For this section, I calculated the AUC for both Random Forest Tree and Decision Tree for various tree size (depth). The fitted decision tree is as in Q1-2 and the random forest is:

```
RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                      criterion='gini', max_depth=None, max_features='auto',
                      max_leaf_nodes=None, max_samples=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=2,
                      min_weight_fraction_leaf=0.0, n_estimators=100,
                      n_jobs=None, oob_score=False, random_state=None,
                      verbose=0, warm_start=False)
```

The result is shown the figure in the next page.

The ultimate accuracy of random forest (96%) is higher than decision tree (92%). The results show that while random forest accuracy increased by increasing the tree size (depth) up to 20, from that size, the algorithm started to flatten out the accuracy did not change (or in case of decision tree the accuracy worsen then flatten out). The decision tree's accuracy increased up to depth increase to 10 and from there, the accuracy decreased.



Q2-1:

The fitted linear regression model and the **fitting error** is as below:

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
The intercept [7.38412739]
The slope [[0.02128314]]
The root mean square error 3.243561781943323
The r2 0.6861858695095129
```

Q2-2:

The strategy is to first evaluate the Ridge regression model with only $\lambda = 0$. Eliminating λ will make the model, a typical polynomial regression model. I will test various polynomial degrees (2 to 10). We will pick the best fit data (least RMSE). The expectation is to witness the RMSE goes down by increasing the polynomial degree and then at some point the fit get worsen. The strategy to select λ for ridge regression is explained in the **next question**.

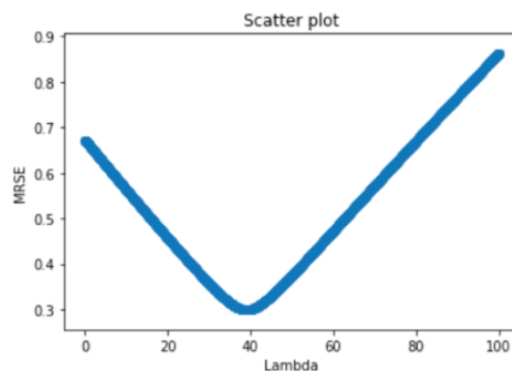
The best polynomial degree for my model is degree 5.

```
Ridge(alpha=0, copy_X=True, fit_intercept=True, max_iter=None, normalize=False,
      random_state=None, solver='auto', tol=0.001)

best polynomial degree 5
best RMSE 0.3225830042568235
```

Q2-3:

For this section, we used the 5 KFold technique. And to define the best λ we did 10000 iteration by starting λ with 0.01 and increase λ with iteration by 0.01. I got the result below:



HW7

This result shows the best lambda is around 40 for the polynomial model of degree 5. The model (polynomial degree 5 and ridge regression lambda 39.03) predicted the coefficient for 400 as 17.43 while the linear regression predicted 15.90 as below.

```
model = LinearRegression()
model.fit(X, y)
y_pred = model.predict([[400]])
print(y_pred)
```

```
[[15.89738509]]
```

```
polynomial_features = PolynomialFeatures(degree=bestDegree)
X_poly = polynomial_features.fit_transform(X)

model = Ridge(alpha=39.03)
model.fit(X_poly, y)
y_pred = model.predict(polynomial_features.fit_transform([[400]]))
print(y_pred)
```

```
[[17.42891061]]
```

By checking the scatter plot for coefficient and temperature, the results for polynomial is more accurate compare to linear regression.

Q3-1 The regularized linear regression is.

$$\hat{\beta}(\lambda) = \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_2^2 \right\}$$

$$\cong \arg \min_{\beta} \left\{ \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^n \beta_j^2 \right\}$$

the criterion to be minimized can be written:

$$f(\beta, \lambda) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

To minimize, we take the derivative of $f(\beta, \lambda)$ by β :

$$\frac{\partial f(\beta, \lambda)}{\partial \beta} = \frac{\partial}{\partial \beta} (y^2 - 2X\beta y + X^T X \beta^2 + \lambda \beta^2)$$

$$= 2X^T \beta - 2Xy + 2\lambda \beta$$

$$= 2(X^T X)\beta - 2X^T y + 2\lambda I \beta = 0$$

$$\Rightarrow 2(X^T X)\beta + 2\lambda I \beta = 2X^T y$$

$$\Rightarrow (X^T X)\beta + \lambda I \beta = X^T y$$

$$\Rightarrow \hat{\beta}_{\text{ridge}} = \beta = (X^T X + \lambda I)^{-1} X^T y \quad \boxed{\text{closed form}}$$

Let assume $\hat{\beta}_{\text{ridge}}$ has Gaussian distribution then to define the distribution we need the mean and variance $N(\mu, \sigma^2)$

$$E[\hat{\beta}_{\text{ridge}}] = E[(X^T X + \lambda I)^{-1} X^T y] = (X^T X + \lambda I)^{-1} X^T E[y]$$

$$= (X^T X + \lambda I)^{-1} X^T X \beta$$

$$\leftarrow E[y] = E[X^T \beta + \epsilon] = X^T \beta + 0 \quad (1)$$

Gaussian noise with 0 as mean

$$X^T X^T = X^{T^2} = X^T X \quad (2)$$

$$\begin{aligned}
 \text{Var}[\hat{\beta}_{\text{Ridge}}] &= \text{Var}[(X^T X + \lambda I)^{-1} X^T y] \\
 &= (X^T X + \lambda I)^{-1} X^T \text{Var}[y] X (X^T X + \lambda I)^{-1} \\
 &= (X^T X + \lambda I)^{-1} X^T \sigma^2 (X^T X + \lambda I)^{-1} X^T \\
 &= \sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}
 \end{aligned}$$

$$\text{Var}[ax] = a^2 \text{Var}[x]$$

The distribution of $\hat{\beta}_{\text{Ridge}}$ will be $N(\mathbb{E}[\hat{\beta}_{\text{Ridge}}], \text{Var}[\hat{\beta}_{\text{Ridge}}])$

$$\begin{aligned}
 \text{Q3-2: Bias}(\hat{y}) &= \mathbb{E}[X^T \hat{\beta}_{\text{Ridge}}] - X^T \beta \\
 &= X^T \mathbb{E}[\hat{\beta}_{\text{Ridge}}] - X^T \beta \\
 &= X^T (\mathbb{E}[\hat{\beta}_{\text{Ridge}}] - \beta)
 \end{aligned}$$

$$\text{Bias} = \mathbb{E}[\hat{y}] - y$$

$\beta = \beta^*$ is considered constant

$$\begin{aligned}
 &= X^T ((X^T X + \lambda I)^{-1} X^T X \beta - \beta) \\
 &= X^T ((X^T X + \lambda I)^{-1} (X^T X + \lambda I - \lambda I) \beta - \beta) \\
 &= X^T (I - \lambda (X^T X + \lambda I)^{-1} \beta - \beta) \\
 &= X^T (\beta - \lambda (X^T X + \lambda I)^{-1} \beta - \beta) \\
 &= -X^T \lambda (X^T X + \lambda I)^{-1} \beta
 \end{aligned}$$

Q3-3

$$E[(X^T \hat{\beta}_{\text{Ridge}} - E[X^T \hat{\beta}_{\text{Ridge}}])^2] = \text{Var}(X^T \hat{\beta}_{\text{Ridge}})$$

$$= X^T \text{Var}(\hat{\beta}_{\text{Ridge}}) X =$$

$$= X^T (\sigma^2 (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1}) X$$

$$\text{Var}(X) = E[(X - E[X])^2]$$

$$\text{Var}(aX) = a^2 \text{Var}(X)$$

according Q3-1

$$\text{Q3-4: } MSE(\hat{\beta}_{\text{Ridge}}) = \text{Var}[X^T \hat{\beta}_{\text{Ridge}}] + \text{bias}[X^T \hat{\beta}_{\text{Ridge}}]^2$$

Results from
b & cwhen $\lambda \rightarrow \infty$

$$\text{Bias}[X^T \hat{\beta}_{\text{Ridge}}] = -X \lambda (X^T X + \lambda I)^{-1} \beta$$

$$\underset{\lambda \rightarrow \infty}{=} \frac{-X \lambda \beta}{X^T X + \lambda I} \approx \frac{-X \beta}{I} \quad (1)$$

this shows that MSE will depend on the coefficient β .

$$\text{Var}[X^T \hat{\beta}_{\text{Ridge}}] = 0 \quad (2)$$

 $\lambda \rightarrow \infty$ (1), (2) \Rightarrow when λ increase, Bias increases and Variance decreaseswhen $\lambda \rightarrow 0$

$$\underset{\lambda \rightarrow 0}{\text{Bias}[X^T \hat{\beta}_{\text{Ridge}}]} = 0 \quad (1), \quad \underset{\lambda \rightarrow 0}{\text{Var}[X^T \hat{\beta}_{\text{Ridge}}]} = \infty \quad (2)$$

(1), (2) \Rightarrow when λ decreases, Bias decreases and Variance increasesall λ are in the denominator

$$\lim_{x \rightarrow \infty} \frac{1}{x} = 0 \quad \lim_{x \rightarrow 0} \frac{1}{x} = \infty$$