# Machine Learning Assignment 1

Kyriacos Xanthos
CID: 01389741

February 3, 2022

## Question 1

In this question we are given a dataset of the form $\mathcal{D} = \{x_i, z_i, y_i\}_{i=1}^{100}$ and we consider a model of the form:

$$Y_i = aZ_i + (b_1 + b_2 Z_i)\cos(X_i) + (c_1 + c_2 Z_i) X_i + (d_1 + d_2 Z_i) X_i^2 + (e_1 + e_2 Z_i) X_i^3 + \epsilon_i \quad (1)$$

where $\epsilon_i$ are independent random variables with mean 0.

### 1.

In this part we consider ridge regression. Note that ridge regression is an extension to linear regression by introducing an $l_2$ penalty on the coefficients to a standard least squares problem. Mathematically we are trying to solve the optimization problem:

$$\underset{\boldsymbol{\beta}}{\operatorname{argmin}} \|\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\beta}\|^2 + \lambda \left(\boldsymbol{\beta}^T \boldsymbol{\beta}\right) \quad (2)$$

where in our example $\boldsymbol{\beta} = \begin{bmatrix} a & b_1 & b_2 & c_1 & c_2 & d_1 & d_2 & e_1 & e_2 \end{bmatrix}^T$, $\lambda$ is the regularization parameter, and $\boldsymbol{y}$ and $\boldsymbol{X}$ are defined as:

$$\boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad, \quad \boldsymbol{X}^T = \begin{bmatrix} Z_1 & Z_2 & \dots & Z_n \\ \cos(X_1) & \cos(X_2) & \dots & \cos(X_n) \\ Z_1\cos(X_1) & Z_2\cos(X_2) & \dots & Z_n\cos(X_n) \\ X_1 & X_2 & \dots & X_n \\ Z_1 X_1 & Z_2 X_2 & \dots & Z_n X_n \\ X_1^2 & X_2^2 & \dots & X_n^2 \\ Z_1 X_1^2 & Z_2 X_2^2 & \dots & Z_n X_n^2 \\ X_1^3 & X_2^3 & \dots & X_n^3 \\ Z_1 X_1^3 & Z_2 X_2^3 & \dots & Z_n X_n^3 \end{bmatrix}$$

where $n$ is 100 for our dataset.

We will use the function `Ridge` from the package `sklearn` in python which solves equation 2 and estimates the coefficients in the following form:

$$\hat{\beta} = \left(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}\right)^{-1}\mathbf{X}^T\mathbf{y} \quad (3)$$

We set the `fit_intercept` parameter to `False` so that it does not include an intercept and use the `svd` solver which performs singular value Decomposition at a low computational cost. We will assess our model performance by using the root mean squared error, defined as:

$$\text{rmse} = \left(\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2\right)^{1/2} \quad (4)$$

where $\hat{y} = X\beta$. With `sklearn` this is calculated using the function `mean_squared_error` with the argument `squared=False`.

$\lambda$ is actually a hyperparameter and we need to tune it so that we get the optimum value. To tune it we use cross validation. The idea is that we split the data $\mathcal{D}$ randomly, into K-folds. We then fit the model K times, we validate on $k^{th}$ fold and train on the remaining $K - 1$ folds. By averaging across multiple folds of the data we decrease the variability in the model fit and validation error caused by the randomness of the splits. Figure 1 shows how the model RMSE varies for different parameters $\lambda$ and for $K = 5$ folds. It is clear that many of the folds have not reached a minimum and suggest that the minimum is less than 0. However, we cannot have $\lambda < 0$ because then the matrix $\left(X^T X + \lambda I\right)^{-1}$ would not be invertible.
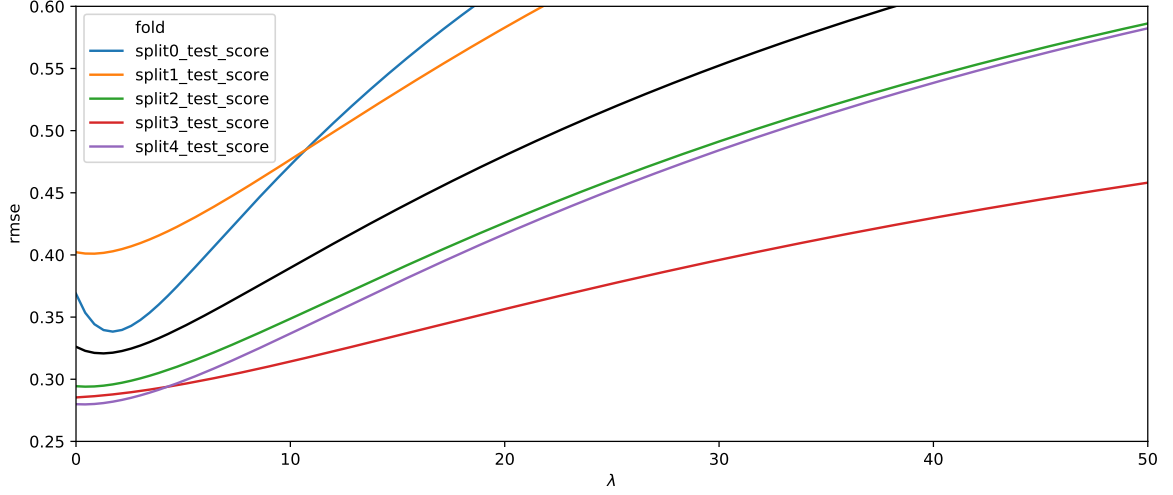


Figure 1: RMSE as defined in 4 against different values of hyperparameter $\lambda$ for different folds.

We utilise cross validation using the `GridSearchCV` function from `sklearn` to search for the optimum value of $\lambda$ that minimises the RMSE. This value is found to be $\lambda = 1.2156$ and looking at figure 2 we see that it indeed minimises the RMSE. We interpret this result with a grain of salt since the splits in figure 1 did not show a clear minimum before zero but this might be because our dataset only has 100 values and so the splits just consist of 20 points.
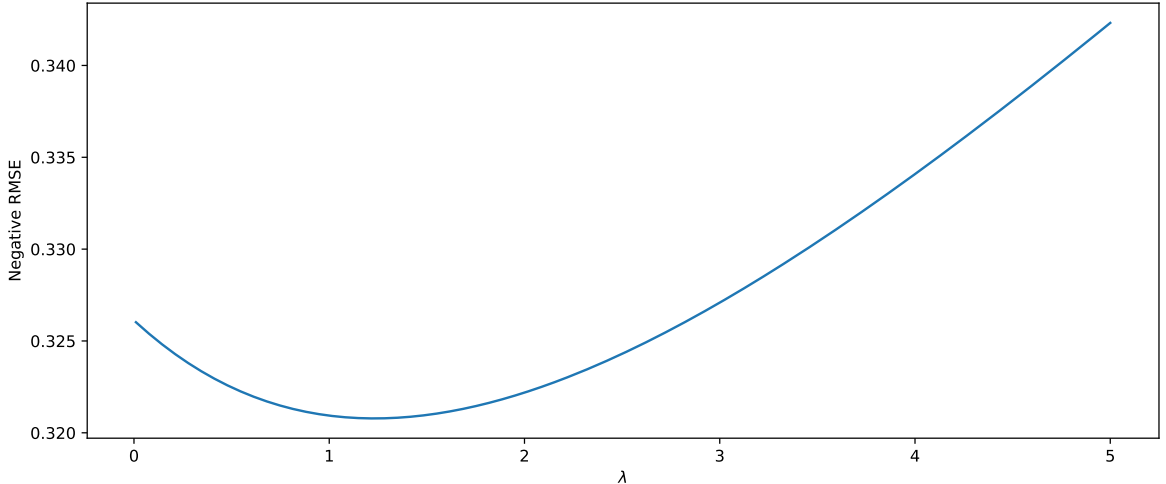


Figure 2: RMSE as defined in 4 against different values of hyperparameter $\lambda$.

Using this optimum value of $\lambda$ we can re-run the Ridge regression and obtain the model summary in figure 3. The model seems to predict the true $y$ with a high accuracy, with RMSE $= 0.286$ and the residual plot shows that there is not any particular bias. The coefficients are found to be:

$$\hat{\beta} = \begin{bmatrix} 0.0159 & 0.9993 & 0.8435 & 0.0242 & 0.01536 & -0.0006 & 0.1834 & 0.01453 & -0.0006 \end{bmatrix}^T$$
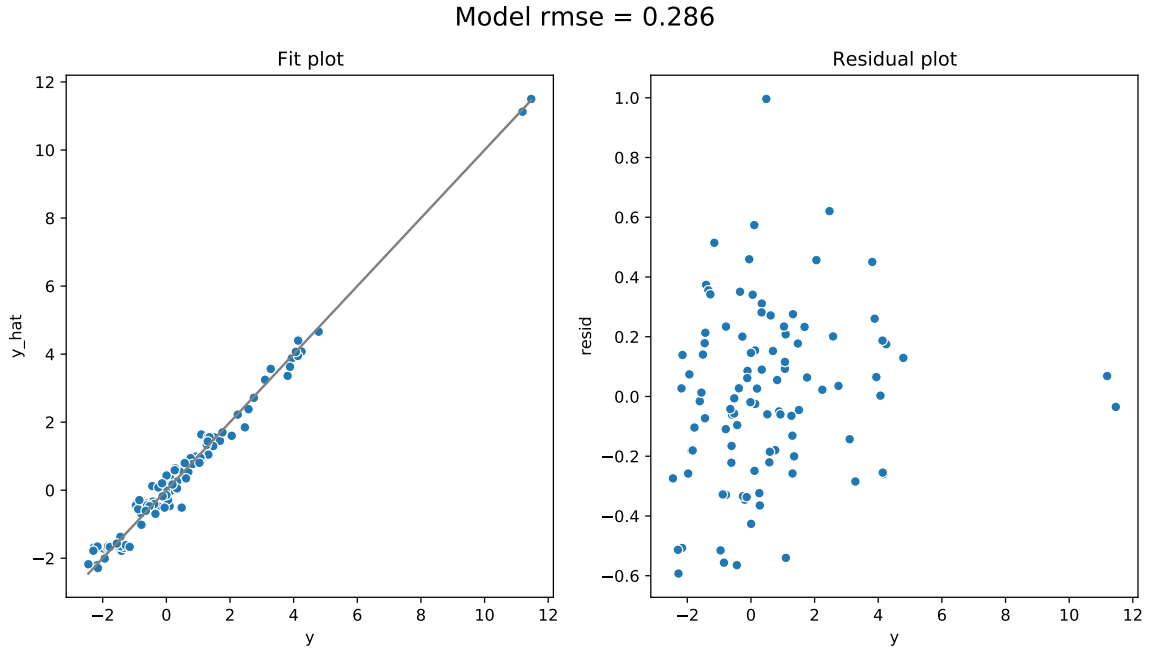
Figure 3: Summary of Ridge regression model with $\lambda = 1.2156$. Comparison of predicted y values with true y values (left) and residual plot (right)

Now we try to predict the value of $Y^*$ for $X^* \in \{-7, -6.9, -6.8, \ldots 6.8, 6.9, 7\}$ and $Z^* \in \{0, 1\}$. The equation we use to predict this is $Y^* = X^*\hat{\beta}$ where $\hat{\beta}$ are the model parameter from our fit above.

We can see the fit of these two models in figure 4. It is clear that our model performs well for both $z$ values. The fit seems to follow the trend of all points and even captures the few points which have higher values of $y$. This confirms a good model qualitatively on unseen data.
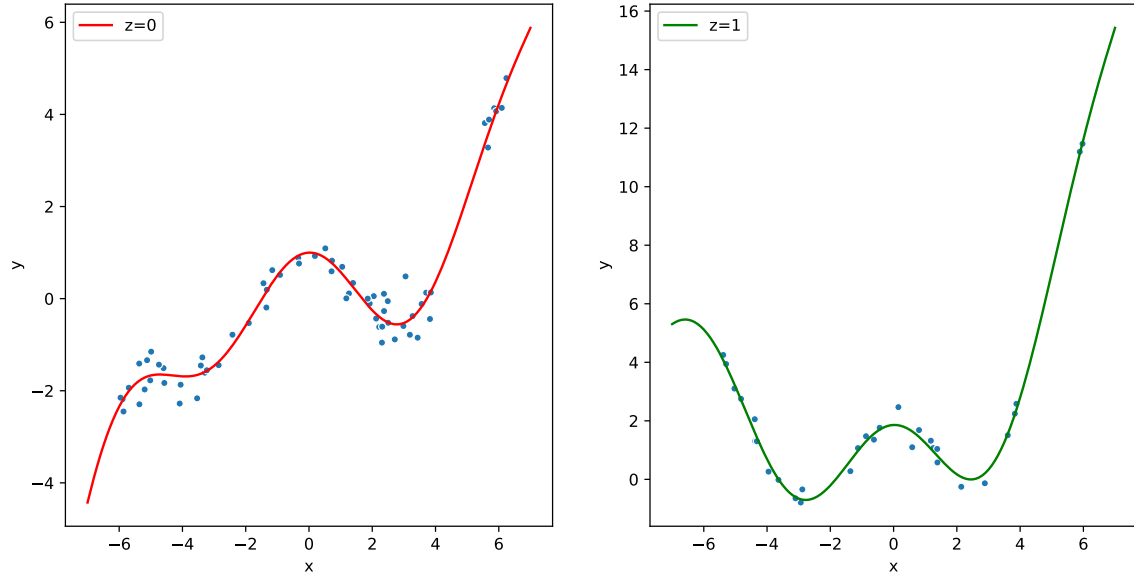


Figure 4: Fit of the model using Ridge regression on unseen data for $z = 0$ (left) and $z = 1$ (right).
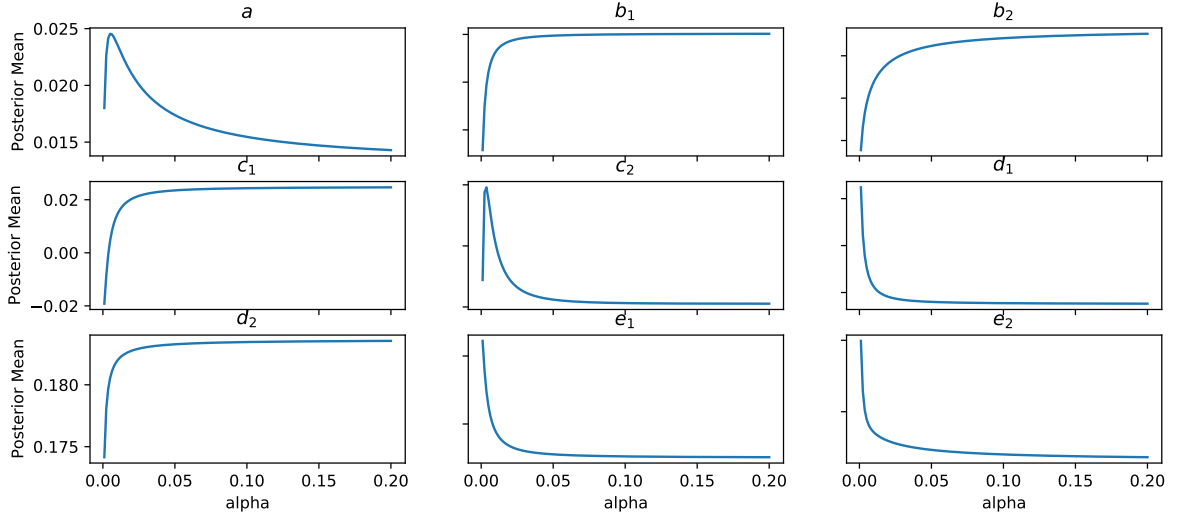
Figure 5: Posterior mean of 9 model parameters as a function of $\alpha$

## 2.

We now perform Bayesian inference to determine the same parameters of this problem. The fundamental idea is that these parameters cannot be considered as single values, instead they acquire probability distributions. We are therefore interested in the posterior distribution over the parameters $\boldsymbol{\beta}$,

$$p(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y}, \sigma) = \frac{p(\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{X}, \sigma)p(\boldsymbol{\beta})}{p(\mathbf{y} \mid \mathbf{X}, \sigma)} \tag{5}$$

where $\sigma^2 = 0.01$ and the marginal likelihood is written as $p(\mathbf{y} \mid \mathbf{X}, \sigma) = \int p(\mathbf{y} \mid \boldsymbol{\beta}, \mathbf{X}, \sigma)p(\boldsymbol{\beta})d\boldsymbol{\beta}$. We will consider a Gaussian prior of the form $p(\boldsymbol{\beta}) = \mathcal{N}_{\boldsymbol{\beta}}(\mathbf{0}, \Lambda)$, where $\Lambda = \alpha \boldsymbol{I}$, the multivariate normal centred at zero and with scaled identity covariance which encodes the belief that the parameters should be small which means that $\alpha$ is a parameter which we choose.

This is solved by [1]

$$\begin{aligned} p(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y}, \sigma) &= \frac{\mathcal{N}_{\mathbf{y}}(\mathbf{x}\beta, \sigma\mathbf{I})\mathcal{N}_{\beta}(\mathbf{0}, \boldsymbol{\Lambda})}{\int \mathcal{N}_{\mathbf{y}}(\mathbf{X}\beta, \sigma\mathbf{I})\mathcal{N}_{\boldsymbol{\beta}}(\mathbf{0}, \boldsymbol{\Lambda})d\boldsymbol{\beta}} \\ &= \mathcal{N}_{\boldsymbol{\beta}}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \end{aligned} \tag{6}$$

where

$$\boldsymbol{\mu} = \left(\mathbf{X}^T\mathbf{X} + \frac{\sigma^2}{\alpha}\mathbf{I}\right)^{-1}\mathbf{X}^T\mathbf{y} \quad, \quad \boldsymbol{\Sigma} = \sigma^2\left(\mathbf{X}^T\mathbf{X} + \frac{\sigma^2}{\alpha}\mathbf{I}\right)^{-1} \tag{7}$$

Calculating 6 in `python` we can find the probability distribution of each parameter. Plotting the posterior mean of the 9 model parameters as a function of $\alpha$ we obtain figure 5.

It is clear that as the value of $\alpha$ increases, all of the 9 parameters converge either to a maximum or a minimum value. This means that for any $\alpha > 0.2$ we can be confident that the posterior mean of the parameter is converged.

To calculate the predictive posterior distribution we use

$$\begin{aligned} p\left(\mathbf{y}_{\text{new}} \mid \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{y}, \sigma\right) &= \int p\left(\mathbf{y}_{\text{new}} \mid \mathbf{x}_{\text{new}}, \boldsymbol{\beta}, \sigma\right)p(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y}, \sigma)d\boldsymbol{\beta} \\ &= \mathcal{N}_{\mathbf{y}_{\text{new}}}\left(\mathbf{x}_{\text{new}}^T\boldsymbol{\mu}, \mathbf{x}_{\text{new}}^T\boldsymbol{\Sigma}\mathbf{x}_{\text{new}} + \sigma^2\right) \end{aligned} \tag{8}$$

as stated in the notes [1]. We again use the same new values as above, namely $X^* \in \{-7, -6.9, -6.8, \ldots 6.8, 6.9, 7\}$ and $Z^* \in \{0, 1\}$ to obtain the new matrix $\boldsymbol{x}_{\text{new}}$. We then calculate the posterior mean and its 95% credible interval. Credible intervals are analogous to the confidence intervals in frequentist statistics but for the bayesian interpretation. They determine the distribution of possible values of the posterior mean [2]. We calculate the
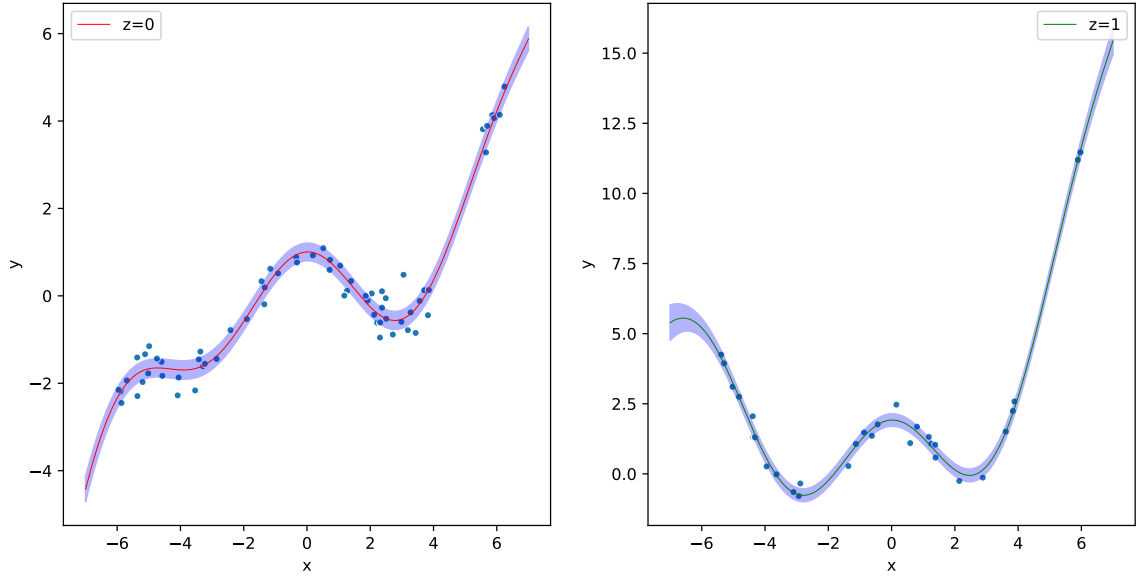
Figure 6: Mean of Posterior Predictive Distribution as a function of $X$ with $z = 0$ (left) and $z = 1$ (right). The shaded blue shows the 95% credible interval of the computation.

credible interval using the 2.5% and 97.5% quantiles of the Gaussian distribution with the parameters stated in equation 6. The plots for $z = 0$ and $z = 1$ are shown in figure 6.

The fit of the model seems to be very good for both values of $z$. The 95% credible interval shows that it covers most of the points in our data except some of the points that exhibit local maxima and local minima. If the range of the credible interval was increased to 90% then the interval would cover all points.

**3.**

Now we will make predictions on $(X, Z) \in \{-5, 7\} \times \{0, 1\}$. We will use 3 different values of $\alpha$ for our predictions. The first one consists of understanding the the regularised ridge least squares problem is exactly the maximum a posteriori estimate of the parameters $\beta$ with the regularisation parameter $\lambda = \sigma^2/\alpha$. So then using the optimum $\lambda$ found in part 1, $\alpha = 0.1/1.2156 = 0.0823$. The second $\alpha$ will be a value much smaller which would mean that the parameters in figure 5 would not have converged (here we use $\alpha = 0.001$) and the third value would be a much larger value that ensures that the posterior mean of the parameters has converged ($\alpha = 1$).
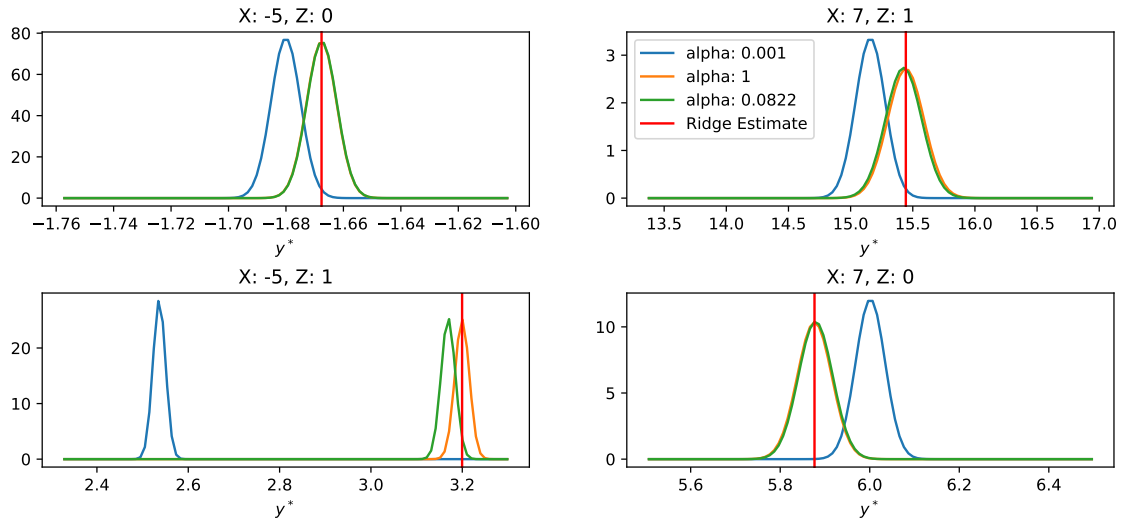


Figure 7: PDF of Predictive distributions for different observation points.

Looking at figure 7 we can observe that the posterior distributions for $\alpha = 0.0822, 1$ are very close for all 4 observations and this makes sense because as observed in figure 5 our parameters seemed to almost converge after 0.1. For $\alpha = 0.001$ as expected in all 4

5

observations the posterior distribution is far away from the other 2 values especially for $X = -5, Z = 1$ which shows that the prediction would not be accurate. Indeed, using the relationship with the $\lambda$ of the ridge regression we see that a very small value of $\alpha$ means there is heavy regularisation in the problem which was not the case when we searched for the optimum $\lambda$ using cross validation. The red vertical line shows that the predictions using Bayesian inference are almost equivalent to the Ridge regression problem, as expected.

## Question 2

We now consider a slightly modified version of the Million Song Dataset [3]. For each song in the dataset there are 91 attributes with the first one being the year it was published. We will create models that predict if a song was produced before or after 1980. This means we will change the first attribute to a binary variable for each song, namely it will be a 0 if it was produced before 1980 and 1 if it was produced after. We can see the distribution of songs for each year in figure 8.
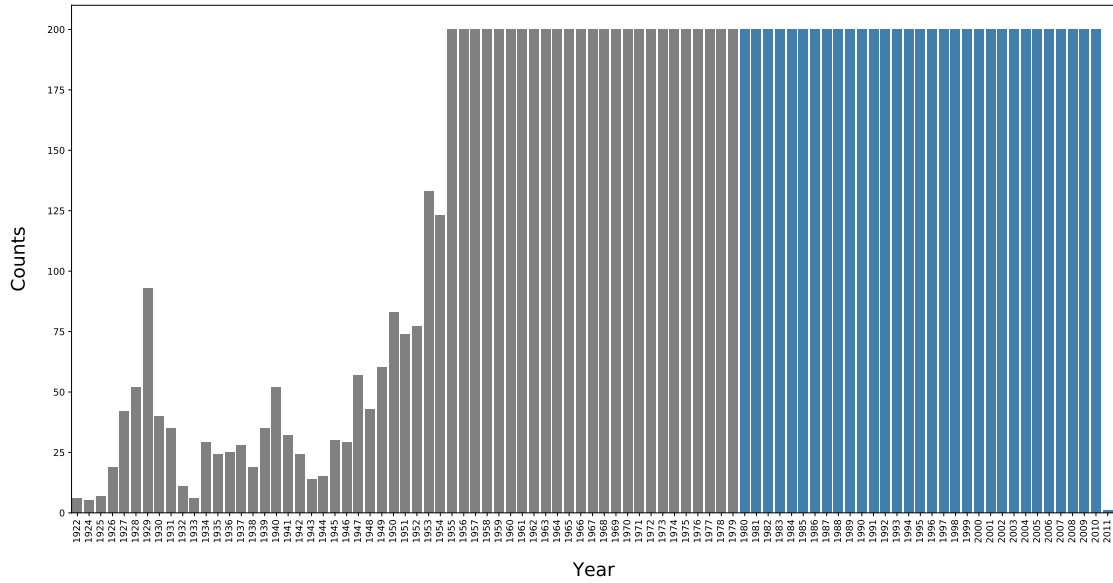
Figure 8: Distribution of songs produced each year. Grey are all the songs before 1980 and blue all the songs after.

We can observe that the number of songs after 1955 has been right censored to 200 counts. This was done in order to have roughly the same number of songs before and after 1980. Indeed when we count them, we have a dataset with 6201 songs after 1980 and 6322 songs before, a balanced dataset.

Since we have 90 attributes, we cannot visualise the attributes to explore some of their features. We note that we have no missing data, and therefore we can right away use classification. In order to be able to assess the efficiency of our classifiers, we split the data to training and testing using the function `train_test_split` from `sklearn` which shuffles the data and then randomly splits it in training and testing subsets. Since we have a large dataset we use a test size of 30% to reduce the computational cost in training. We then use the `StandardScaler` function of `sklern` to scale all of the attributes of the songs so that they have mean 0 and variance 1. We do that so that the effect of each attribute will be captured in our model and we make sure to do the scaling after the split of train and testing sets.

### LDA

The Linear Discriminant Analysis (LDA) classifier partitions our data into regions with the same class predictions via separating hyperplanes. This is a generative approach that models the class conditional densities $p(x \mid y = k)$ and class priors $p(y = k)$. Then using Bayes Theorem we compute the posterior probabilities

$$p(y = k \mid x) = \frac{p(x \mid y = k)p(y = k)}{p(x)} \tag{9}$$

where $p(x) = \sum_{j=1}^{K} p(x \mid y = j)p(y = j)$ [1].

Then it assumes that the class-conditional densities are Gaussian,

$$x \mid y = k \sim \mathcal{N}(\mu_k, \Sigma), \quad \text{for } 1 \le k \le K \tag{10}$$

where $\Sigma$ is the covariance matrix and $\mu_k$ are the means. Then to create the boundaries we use parameter estimation using Maximum Likelihood. These estimates are

$$
\begin{aligned}
\hat{\pi}_k &= \frac{N_k}{N} \\
\hat{\mu}_k &= \frac{1}{N_k} \sum_{i, y^{(i)} = k} x^{(i)} \\
\hat{\Sigma} &= \frac{1}{N} \sum_{k=1}^{K} \sum_{i, y^{(i)} = k} \left( x^{(i)} - \mu_k \right) \left( x^{(i)} - \mu_k \right)^T
\end{aligned}
\tag{11}
$$

where $N_k$ are the number of $i$s such that $y^{(i)} = k$, $N$ are the total number of points.

To calculate this we use the `LinearDiscriminantAnalysis` function of `sklearn`. It has a linear decision boundary which is generated by fitting Gaussian class conditional densities to the data assuming that all classes share the same covariance matrix [4]. We assume no shrinkage, and use the `svd` solver which performs singular value decomposition. We could also use the `lsqr` method which uses the dedicated least squares routine and it is the fastest solver but it produced less accurate results so we preferred `svd`.

We then fit this solver to our training data and make predictions on our testing data. A common metric for the performance of our classifier is a confusion matrix [5] which summarises the correct and false classifications.

|       | FALSE | TRUE |   |       | FALSE | TRUE |   |       | FALSE | TRUE |
|-------|-------|------|---|-------|-------|------|---|-------|-------|------|
| FALSE | 1470  | 428  |   | FALSE | 1478  | 459  |   | FALSE | 994   | 910  |
| TRUE  | 475   | 1384 |   | TRUE  | 479   | 1341 |   | TRUE  | 355   | 1498 |

Table 1: Confusion Matrix of Test Set (30%) using LDA (left), KNN (middle) and Naive Bayes (right)

As seen in table 1 our classifier was able to classify 76% of the songs correctly if they were before or after 1980 (Accuracy) with error rate 24%. Another useful metric for the performance of classifiers are the Receiver Operatic Characteristic Curves (ROC) [6]. They illustrate the diagnostic ability of a binary classifier as its discrimination threshold is varied. We plot the True Positive Rate

$$\text{TPR} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

against the False Negative Rate (FNR $= 1 - \text{TPR}$)[1]. The perfect classifier would exhibit a vertical line to TPR $= 1$ and then a horizontal line to FPR $= 1$. We can see the ROC Curve of our classifier in figure 10 and it clear that it follows the pattern we expect with the curve peaking relatively close to to the (0,1) point.

A quantitative metric that can judge the performance of our classifier is the AUC (Area Under the Curve) of an ROC curve. This area under the curve is equal to the probability that the classifier will rank a randomly chosen positive instance higher than a randomly chosen negative one [7]. The LDA classifier achieves an 84.6 % AUC score for the Training set and 83.3% for the Test set. These are very good scores for a dataset in its high dimensional nature.

---

[1]True Positive is when a test result correctly indicates the presence of a condition. False Negative is a test result which wrongly indicates that a particular condition is absent.

## KNN

K-nearest neighbours (KNN) is a discriminative approach which aims to model the conditional probability $p(y = k \mid x)$ directly. This method measures distances between each data point and the number of nearest neighbours. The pseudo algorithm for this method is found in Algorithm 1.

---

**Algorithm 1** KNN algorithm [8]

---

**Input:** $(X_i, C_i)$ where $i = 1, \ldots n$ be $n$ data points, $X_i$ denotes feture values and $C_i$ denotes labels for each $i$. Integer $\boldsymbol{k}$.

**for** *i=1,2,...* **do**
$\quad \mid$ Calculate $\boldsymbol{d}(x, x_i)$ where $\boldsymbol{d}$ is the Euclidean distance
**end**
Arrange $\boldsymbol{d}$ in increasing order
$\quad$ Take first $\boldsymbol{k}$ distances from the sorted list
$\quad$ Find those $\boldsymbol{k}$-points corresponding to these $\boldsymbol{k}$-distances
$\quad$ Let $\boldsymbol{k}_i$ be the number of points in the $i^{\text{th}}$ class among $\boldsymbol{k}$ points

**for** *i=1,2,...* **do**
$\quad$ **for** *j=1,2,...* **do**
$\quad\quad$ **if** $\boldsymbol{k}_i > \boldsymbol{k}_j$ *for all* $i \neq j$ **then**
$\quad\quad\quad \mid$ put x in class i
$\quad\quad$ **end**
$\quad$ **end**
**end**
**Return:** Data points in different classes

---

To compute the KNN we use the `KNeighborsClassifier` function from `sklearn` which implements learning based on the $k$ nearest neighbors of each query point. We choose the Euclidean distance as the metric since it provides the least error rate. The function uses uniform weights (all points i each neighbourhood are weighted equally). The algorithm used to compute the nearest neighbours was `ball_tree` with leaf size 30. This is a similar algorithm to 1, more about the algorithm can be found in [9]. The leaf size is 30 which is a good compromise between the computational time and amount of nodes created.

In order to tune the hyperparameter $k$, we try different $k$-values and for each value we calculate the mean accuracy rate (the times the predicted label matches the true label). To do this we again use cross validation with 5 folds as described in question 1. Using the function `GridSearchCV` we were able to plot the mean accuracy for each value of $k$ for $k \in \{5, 10, 15, 20, 25, 30, 35, 40, 45, 50\}$ as seen in figure 9. We see that the accuracy stabilizes after $k = 25$ so we choose this value to avoid underfitting (decision boundary is closer to being linear as $k$ increases). A smaller value would result in the decision boundary to be overly flexible and find patterns inn the data that do not correspond to the Bayes decision boundary (overfitting).

We can asses the success of this classifier again by looking at the confusion matrix, error rate, ROC curve and AUC score. The confusion matrix can be seen in table 1. The error rate is 25%, 1% larger than the LDA classifier. The ROC curve seems to be slightly better than the LDA, with an AUC of 81.9% for the test set. For a model with only 12 nearest neighbours, these are very good statistics for the classifier.

## Naive Bayes Classifier

The Naive Bayes classifier is a generative approach and it assumes all measured variables and features are independent given the label. The clasifer makes the assumption that the features are conditionally independent for the class label:

$$p(x \mid y = k) = \prod_{i=1}^{d} p(x_i \mid y = k, \beta_{ik}) \tag{12}$$
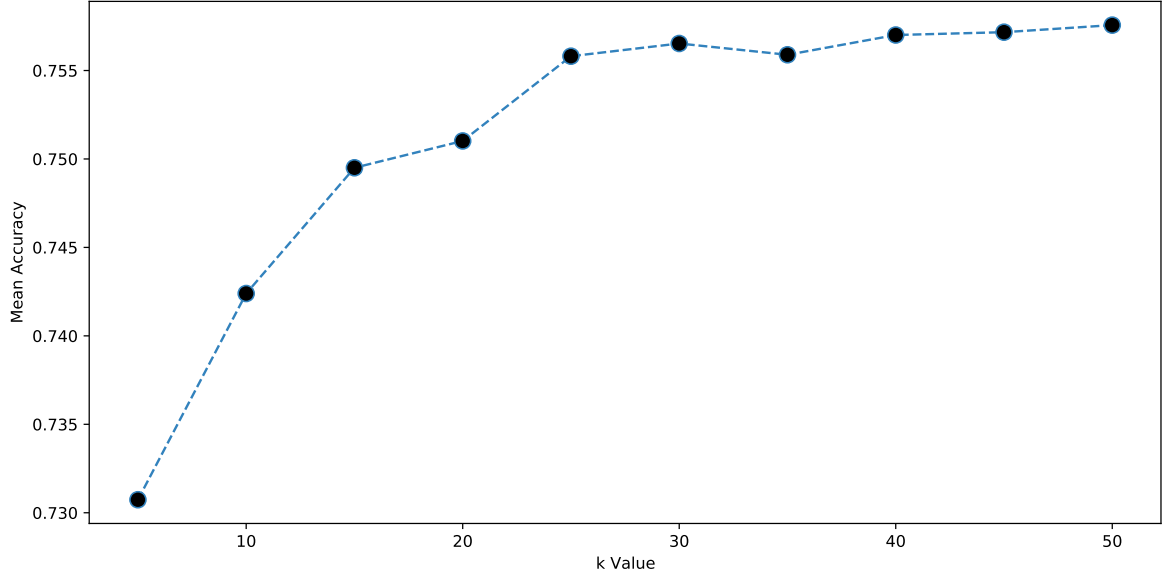
8

Figure 9: Mean Accurary against number of nearest neighbours for KNN model.

Since in this dataset we have real-valued features we decided to use the Gaussian distribution, ie:

$$p\left(x \mid y = k', \theta\right) = \prod_{j=1}^{d} \Phi\left(x_j \mid \mu_{jk}, \sigma_{jk}^2\right) \tag{13}$$

where $\Phi$ is the probability density function of the normal distribution.

We then use the Maximum Likelihood approach for parameter estimation, with the probability for a single data point $x^i, y^i$ given by:

$$p\left(x^{(i)}, y^{(i)} \mid \theta, \pi\right) = p\left(y^{(i)} \mid \pi\right) \prod_{j=1}^{d} p\left(x_j^{(i)} \mid \theta_j\right) = \prod_{k=1}^{K} \pi_k^{1\left(y^{(i)}=k\right)} \prod_{j=1}^{d} \prod_{k=1}^{K} p\left(x_j^{(i)} \mid \theta_{jk}\right)^{1\left(y^{(i)}=k\right)} \tag{14}$$

where $\pi_k$ is the prior probability that $y^{(i)} = k$.

To fit the Naive Bayes classifier we will use the `GaussianNB` function from `sklearn` which implements the Gaussian Naive Bayes algorithm mentioned above. Fitting on the training set and testing on the test set we observe that the confusion matrix is much worse than the other two classifiers, with an error rate of 33%. This is mostly due to the number of instances the true label was 0 (before 1980) and was misclassified by 1 (after 1980). Looking at the ROC curve in figure 10 we see a slight decrease in the peak with comparison of the other curves. The AUC score for the test set was just at 75.6%, not bad in general but worse than our other two classifiers.

## Conclusion

All three classifiers are able to classify songs before and after 1980 fairly good. By the nature of the problem, there is no hard decision boundary at 1980 that differentiates music in a fundamental way. That is to say, we can expect that songs of 1982 for example do not defer a lot with 1978. This suggests that an extension to the study would be to create softer margins, and group maybe decades together to see clear difference between the classes. Nevertheless, since there is a clear difference of the music style for songs of the 40s, 50s, 60s and 70s with everything after the 90s, we achieved very good classifiers that can predict the age of a song. Right-censoring the data so that we have roughly equal number of songs before and after the decision boundary really helped with the classification problem and meant that we can use the underlying assumptions of LDA and Naive Bayes. When trying to decide between the three classifiers we can refer to figure 11 to test that LDA seems to have the best ROC Curve with the highest AUC followed by KNN and then by Naive Bayes.
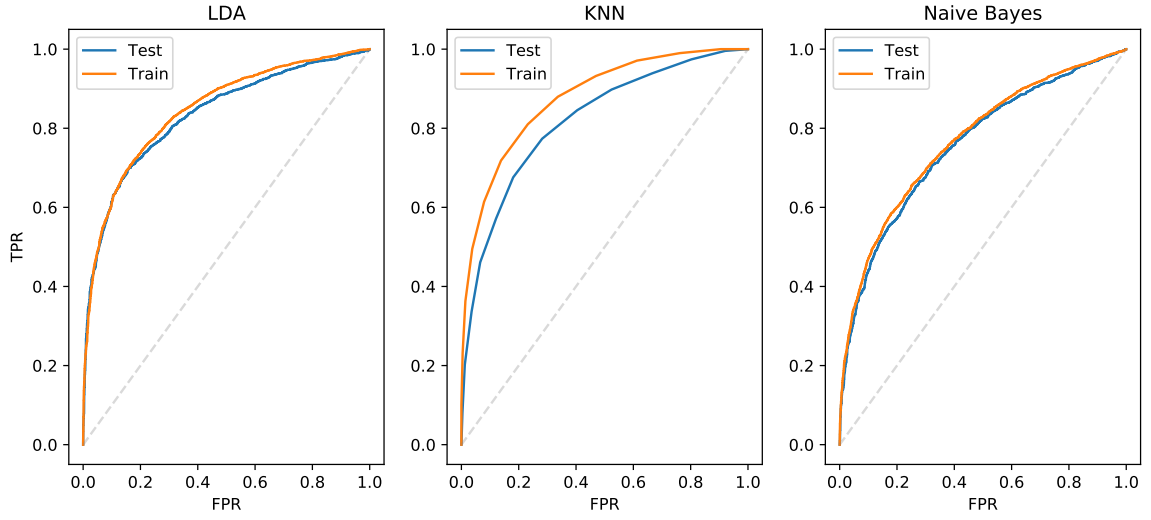
Figure 10: ROC curves of LDA (left), KNN (right) and Naive Bayes (right) for both Test and Training sets (30-70 % split)

| Classifier | LDA | KNN | NB |
|:---:|:---:|:---:|:---:|
| **Computation time** | 0.5s | 5.7s | 0.4s |

Table 2: Computational time of the fitting of all three classifiers for the same training subset of the data. The measurements were made on the same computer, with minimal background processes at the fitting procedure.

KNN was the only classifier which we could tune hyperparameters, which suggest that the other two models cannot be improved to a large extent, but we also believe we have found a good hyperparameter $k$ for the KNN that does not overfit our data. Comparing the computational times of the three classifiers in table 2 it is clear that the KNN is the slowest and LDA and Naive Bayes to be almost identical with Naive Bayes having the edge. This makes sense since Naive Bayes makes the (crude) assumption of independent features given a label. We did not challenge this assumption because of how many features we had, but computing the correlation matrix did not show any clear correlation with the labels.

The KNN iteratively goes through every point and assignns it to a classs which is computationally extensive. The ball tree algorithm does not work very well with high dimensions since the split of the tree at each node is based on all of the labels of our dataset. LDA works very well with high dimensional data.

In conclusion, the LDA seems to be the best classifier, satisfying both accuracy, computational cost, and clear interpretability of results.
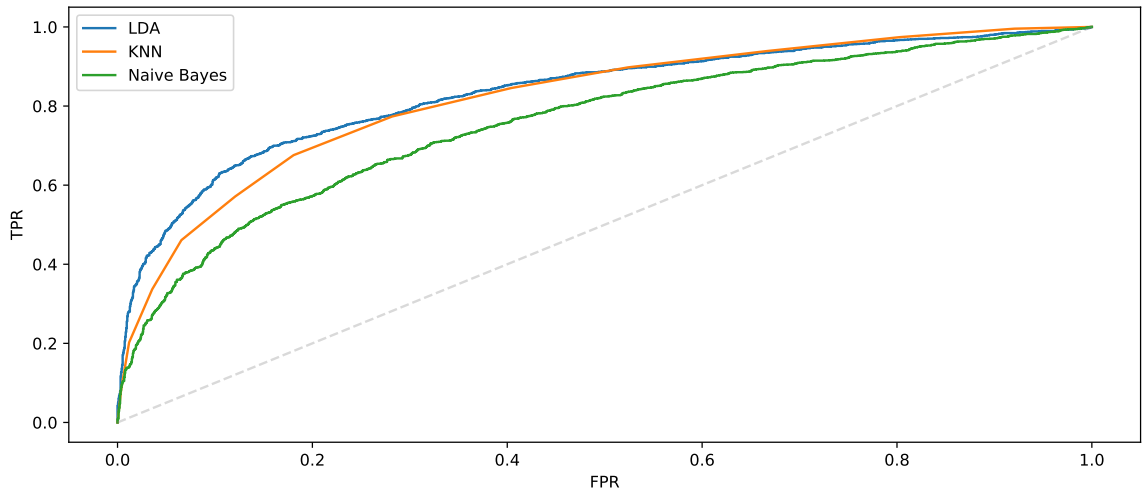


Figure 11: ROC curves of Test sets of all 3 classifiers (LDA, KNN, Naive Bayes)

# References

[1] Dr Sarah Filippi. Machine learning lecture notes. *Imperial College London*, pages 30–40, October 2021.

[2] Richard D Morey, Rink Hoekstra, Jeffrey N Rouder, Michael D Lee, and Eric-Jan Wagenmakers. The fallacy of placing confidence in confidence intervals. *Psychonomic bulletin & review*, 23(1):103–123, 02 2016. doi: 10.3758/s13423-015-0947-8. URL https://pubmed.ncbi.nlm.nih.gov/26450628.

[3] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.

[4] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

[5] Wikipedia. Confusion matrix — Wikipedia, the free encyclopedia. http://en.wikipedia.org/w/index.php?title=Confusion%20matrix&oldid=1058352752, 2022. [Online; accessed 29-January-2022].

[6] Wikipedia. Receiver operating characteristic — Wikipedia, the free encyclopedia. http://en.wikipedia.org/w/index.php?title=Receiver%20operating%20characteristic&oldid=1067917334, 2022. [Online; accessed 29-January-2022].

[7] Fawcett, Tom. An introduction to roc analysis. 2006. URL https://www.math.ucdavis.edu/~saito/data/roc/fawcett-roc.pdf.

[8] Rahul Saxena. Knn classifier, introduction to k-nearest neighbor algorithm. https://dataaspirant.com/k-nearest-neighbor-classifier-intro/, 2022. [Online; accessed 29-January-2022].

[9] Stephen M. Omohundro. Five balltree construction algorithms. Technical report, 1989.

# A   Code for Question 1

```python
# %%
import pandas as pd
import numpy as np
import os
import matplotlib.pyplot as plt
import seaborn as sns
# from sklearn.preprocessing import PolynomialFeatures
# from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.pipeline import make_pipeline
from sklearn.model_selection import GridSearchCV, KFold
# from sklearn.preprocessing import LabelEncoder   # binary encoding
from sklearn.model_selection import train_test_split
# from sklearn.preprocessing import StandardScaler
# from sklearn.cluster import KMeans
import sklearn.pipeline
from sklearn.linear_model import Ridge
from sklearn.linear_model import RidgeCV
from scipy.stats import norm
seed = 222
np.random.seed(seed)

plt.rcParams['figure.figsize'] = (12,5)
plt.rcParams['figure.dpi'] = 80
```

```python
25
26
27  # %%
28  d = pd.read_csv('data.csv')
29  d.describe()
30
31  # %%
32  sns.scatterplot(x='x', y='y', data = d, hue= 'z', s=20)
33  plt.savefig('Plots/my_data.pdf')
34  plt.show()
35
36  # %%
37  def create_design_matrix(x,z):
38
39      X = np.vstack((z, np.cos(x), z * np.cos(x), x, z * x, x * x,
40      z * x * x, x * x * x, z * x * x * x ))
41      X = X.T
42      return(X)
43
44  X =create_design_matrix(d.x, d.z)
45  y = d.y
46
47
48  # %%
49  def get_coefs(m):
50      if (isinstance(m, sklearn.pipeline.Pipeline)):
51          m = m.steps[-1][1]
52      if m.intercept_ is None:
53          return m.coef_
54      return np.concatenate([[m.intercept_], m.coef_])
55
56
57  # %%
58  def model_fit(m, X, y, plot = False):
59
60      y_hat = m.predict(X)
61      rmse = mean_squared_error(y, y_hat, squared=False)
62      res = pd.DataFrame(
63              data = {'y': y, 'y_hat': y_hat, 'resid': y - y_hat}
64      )
65      if plot:
66              plt.figure(figsize=(12, 6))
67              plt.subplot(121)
68              sns.lineplot(x='y', y='y_hat', color="grey", data = pd.DataFrame(dat
69              sns.scatterplot(x='y', y='y_hat', data=res).set_title("Fit plot")
70              plt.subplot(122)
71              sns.scatterplot(x='y', y='resid', data=res).set_title("Residual plot
72              plt.suptitle("Model rmse = " + str(round(rmse, 4)), fontsize=16)
73              plt.savefig('Plots/model_summary.pdf')
74              plt.show()
75      return(rmse)
76
77
78  # %%
79  r = make_pipeline(
80      Ridge(alpha=1)
81  ).fit(X, y)
82
83  # %%
84  model_fit(r, X, y, plot=True)
85
86  # %%
87  alphas = np.linspace(0.01, 5, num=150)
88
```

```
 89  gs = GridSearchCV(
 90      make_pipeline(
 91          Ridge()
 92          ),
 93          param_grid={'ridge__alpha': alphas},
 94          cv=KFold(5, shuffle=True, random_state=1234),
 95          scoring="neg_root_mean_squared_error"
 96  ).fit(X, y)
 97
 98  # %%
 99  print(gs.best_params_)
100  model_fit(gs.best_estimator_, X, y, plot=True)
101  best_alpha = gs.best_params_['ridge__alpha']
102
103  # %%
104  cv_res = pd.DataFrame(
105      data = gs.cv_results_
106  ).filter(
107      regex = '(split[0-9]+|mean)_test_score'
108  ).assign(
109  # Add the alphas as a column
110      alpha = alphas
111  )
112  cv_res.update(
113  # Convert negative rmses to positive
114  -1 * cv_res.filter(regex = '_test_score')
115  )
116  sns.lineplot(x='alpha', y='mean_test_score', data=cv_res)
117  plt.xlabel('$\lambda$')
118  plt.ylabel('Negative RMSE')
119  # plt.savefig('Plots/single_l.pdf')
120  plt.show()
121
122  # %%
123  d1 = cv_res.melt(
124      id_vars=('alpha','mean_test_score'),
125      var_name='fold',
126      value_name='rmse'
127  )
128
129  sns.lineplot(x='alpha', y='rmse', color='black', ci = None, data = d1) # Plot
130  sns.lineplot(x='alpha', y='rmse', hue='fold', data = d1) # Plot the curves for
131  plt.xlabel('$\lambda$')
132  # plt.savefig('Plots/splits.pdf')
133  plt.show()
134
135  # %%
136  m = make_pipeline(
137      Ridge(alpha=best_alpha, fit_intercept = False, solver='svd' )
138  ).fit(X, y)
139
140  # %%
141  model_fit(m, X, y, plot= True)
142
143  # %%
144  new_x = np.arange(-7, 7.1, step=0.1)
145
146  X_new1 = create_design_matrix(new_x, np.zeros(new_x.shape))
147  X_new2 = create_design_matrix(new_x, np.ones(new_x.shape))
148
149  # %%
150  y_hat1 = m.predict(X_new1)
151  y_hat2 = m.predict(X_new2)
152
```

```
153
154  # %%
155  plt.figure(figsize=(12, 6))
156  plt.subplot(121)
157  sns.scatterplot(x='x', y='y', data = d[d.z == 0], s=20)
158  plt.plot(new_x, y_hat1, color = 'red', label ='z=0')
159  plt.legend()
160  plt.subplot(122)
161  sns.scatterplot(x='x', y='y', data = d[d.z == 1], s=20)
162  plt.plot(new_x, y_hat2, color = 'green', label='z=1')
163  plt.legend()
164  # plt.savefig('Plots/two_models_ridge.pdf')
165  plt.show()
166
167  # %% [markdown]
168  # ## Part 2
169
170  # %%
171  sigma_sq = 0.1
172
173  # %%
174  # see equations page 31 of slides
175  def solve_posterior(
176      X: np.ndarray, y: np.ndarray, sigma_sq: np.ndarray, alpha: np.ndarray
177  ) -> np.ndarray:
178      n = X.shape[0]
179      p = X.shape[1]
180      mu = np.linalg.solve(X.T @ X + sigma_sq / alpha * np.eye(p), X.T @ y)
181      Sigma = sigma_sq * np.linalg.solve(
182          X.T @ X + (sigma_sq / alpha) * np.eye(p), np.eye(p)
183      )
184      return mu, Sigma
185
186  # %%
187  alphas = np.linspace(0.001, 0.2, 150)
188  mus = np.empty((alphas.shape[0], 9))
189  for i, a in enumerate(alphas):
190      mus[i] = solve_posterior(X, y, sigma_sq, a)[0]
191  mus.shape
192
193
194  # %%
195  my_lst = [(0,0), (0,1), (0,2), (1,0), (1, 1), (1,2), (2,0), (2,1), (2,2)]
196  all_par = ['$a$', '$b_1$', '$b_2$', '$c_1$', '$c_2$', '$d_1$', '$d_2$', '$e_1$',
197  fig, axs = plt.subplots(3, 3)
198  for i, (t, p) in enumerate(zip(my_lst, all_par)):
199      axs[t].plot(alphas, mus[:,i])
200      axs[t].set_title(p)
201
202  for ax in axs.flat:
203      ax.set(xlabel='alpha', ylabel='Posterior Mean')
204
205  for ax in axs.flat:
206      ax.label_outer()
207
208  # fig.savefig('Plots/all_par.pdf')
209
210  # %%
211  mu, sigma = solve_posterior(X, y, sigma_sq, 1)
212
213  # %%
214  pred_1 = X_new1 @ mu
215  pred_2 = X_new2 @ mu
216  pred_1_var = X_new1 @ sigma @ X_new1.T + sigma_sq
```

```
217  pred_1_var = X_new2 @ sigma @ X_new2.T + sigma_sq
218
219  # %%
220  plt.scatter(d.x,d.y)
221  plt.plot(new_x, pred_1, color = 'red')
222  plt.plot(new_x, pred_2, color = 'green')
223
224
225  # %%
226  cred_int_1 = np.empty((X_new1.shape[0],2))
227  for i in range(0, X_new1.shape[0]):
228      big_sigma = X_new1[i] @ sigma @ X_new1[i].T + sigma_sq
229      cred_int_1[i,0] = norm.ppf(0.025, pred_1[i], big_sigma)
230      cred_int_1[i,1] = norm.ppf(0.975, pred_1[i], big_sigma)
231
232  cred_int_2 = np.empty((X_new2.shape[0],2))
233  for i in range(0, X_new2.shape[0]):
234      big_sigma = X_new2[i] @ sigma @ X_new2[i].T + sigma_sq
235      cred_int_2[i,0] = norm.ppf(0.025, pred_2[i], big_sigma)
236      cred_int_2[i,1] = norm.ppf(0.975, pred_2[i], big_sigma)
237
238  # %%
239  plt.figure(figsize=(12, 6))
240  plt.subplot(121)
241  sns.scatterplot(x='x', y='y', data = d[d.z == 0], s=20)
242  plt.plot(new_x, pred_1, color = 'red', linewidth=0.5, label='z=0')
243  plt.fill_between(new_x, cred_int_1[:,0], cred_int_1[:,1], alpha=0.3, color='blue
244  plt.legend()
245  plt.subplot(122)
246  sns.scatterplot(x='x', y='y', data = d[d.z == 1], s=20)
247  plt.plot(new_x, pred_2, color = 'green', linewidth=0.5, label='z=1')
248  plt.fill_between(new_x, cred_int_2[:,0], cred_int_2[:,1], alpha=0.3, color='blue
249  plt.legend()
250  # plt.savefig('Plots/two_models_bayes.pdf')
251  plt.show()
252
253  # %%
254  X_trial = create_design_matrix(np.array([-5,7, -5, 7]), np.array([0,1, 1,0]))
255
256  # %%
257  pred_bayesian1 = X_trial @ mu
258
259  pred_ridge = m.predict(X_trial)
260
261  # %%
262  ridge_alpha = sigma_sq / best_alpha
263  ridge_alpha
264
265  # %%
266  def calc_all(i):
267      new_mu = solve_posterior(X,y,sigma_sq,0.001)[0]
268      new_sigma =  solve_posterior(X,y,sigma_sq,0.001)[1]
269      mean_small = X_trial[i] @ new_mu
270      sigma_small = X_trial[i] @ new_sigma @ X_trial[i].T
271      new_mu = solve_posterior(X,y,sigma_sq,1)[0]
272      new_sigma =  solve_posterior(X,y,sigma_sq,1)[1]
273      mean_medium = X_trial[i] @ new_mu
274      sigma_medium = X_trial[i] @ new_sigma @ X_trial[i].T
275      new_mu = solve_posterior(X,y,sigma_sq,ridge_alpha)[0]
276      new_sigma =  solve_posterior(X,y,sigma_sq,ridge_alpha)[1]
277      mean_ridge = X_trial[i] @ new_mu
278      sigma_ridge = X_trial[i] @ new_sigma @ X_trial[i].T
279      return(mean_small, sigma_small, mean_medium, sigma_medium, mean_ridge, sigma
280
```

```
281  # %%
282  plt.subplots_adjust(hspace=0.5)
283
284  plt.subplot(221)
285  mean_small, sigma_small, mean_medium, sigma_medium, mean_ridge, sigma_ridge = ca
286  t1 = np.linspace(mean_small - 15*sigma_small, mean_small + 15*sigma_small, 100)
287  plt.plot(t1, norm.pdf(t1, mean_small, sigma_small),label='alpha: 0.001' )
288  plt.plot(t1, norm.pdf(t1, mean_medium, sigma_medium), label='alpha: 1')
289  plt.plot(t1, norm.pdf(t1, mean_ridge, sigma_ridge),label='alpha: 0.0822' )
290  plt.axvline(x=pred_bayesian1[0], color='red', label ='Ridge Estimate')
291  plt.xlabel('$y^*$')
292  plt.title('X: -5, Z: 0')
293  # plt.legend()
294
295  plt.subplot(222)
296  mean_small, sigma_small, mean_medium, sigma_medium, mean_ridge, sigma_ridge = ca
297  t1 = np.linspace(mean_small - 15*sigma_small, mean_small + 15*sigma_small, 100)
298  plt.plot(t1, norm.pdf(t1, mean_small, sigma_small),label='alpha: 0.001' )
299  plt.plot(t1, norm.pdf(t1, mean_medium, sigma_medium), label='alpha: 1')
300  plt.plot(t1, norm.pdf(t1, mean_ridge, sigma_ridge),label='alpha: 0.0822' )
301  plt.axvline(x=pred_bayesian1[1], color='red', label ='Ridge Estimate')
302  plt.title('X: 7, Z: 1')
303  plt.xlabel('$y^*$')
304
305  plt.legend()
306
307  plt.subplot(223)
308  mean_small, sigma_small, mean_medium, sigma_medium, mean_ridge, sigma_ridge = ca
309  t1 = np.linspace(mean_small - 15*sigma_small, mean_small + 55*sigma_small, 100)
310  plt.plot(t1, norm.pdf(t1, mean_small, sigma_small),label='alpha: 0.001' )
311  plt.plot(t1, norm.pdf(t1, mean_medium, sigma_medium), label='alpha: 1')
312  plt.plot(t1, norm.pdf(t1, mean_ridge, sigma_ridge),label='alpha: 0.0822' )
313  plt.axvline(x=pred_bayesian1[2], color='red', label ='Ridge Estimate')
314  plt.xlabel('$y^*$')
315
316  plt.title('X: -5, Z: 1')
317  # plt.legend()
318
319  plt.subplot(224)
320  mean_small, sigma_small, mean_medium, sigma_medium, mean_ridge, sigma_ridge = ca
321  t1 = np.linspace(mean_small - 15*sigma_small, mean_small + 15*sigma_small, 100)
322  plt.plot(t1, norm.pdf(t1, mean_small, sigma_small),label='alpha: 0.001' )
323  plt.plot(t1, norm.pdf(t1, mean_medium, sigma_medium), label='alpha: 1')
324  plt.plot(t1, norm.pdf(t1, mean_ridge, sigma_ridge),label='alpha: 0.0822' )
325  plt.axvline(x=pred_bayesian1[3], color='red', label ='Ridge Estimate')
326  plt.title('X: 7, Z: 0')
327  plt.xlabel('$y^*$')
328  # plt.legend()
329
330  # plt.savefig('Plots/post_all.pdf')
331
332
333  # %%
```

## B  Code for Question 2

```
1  # %%
2  import pandas as pd
3  import numpy as np
4  import os
5  import matplotlib.pyplot as plt
6  import seaborn as sns
7  from sklearn.metrics import mean_squared_error
8  from sklearn.pipeline import make_pipeline
```

```
 9  from sklearn.model_selection import KFold
10  from sklearn.model_selection import train_test_split
11  from sklearn.preprocessing import StandardScaler
12  from sklearn.metrics import confusion_matrix
13  from sklearn.metrics import roc_curve, precision_recall_curve
14  from sklearn.metrics import roc_auc_score
15  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
16  from sklearn.naive_bayes import GaussianNB
17  from sklearn.neighbors import KNeighborsClassifier
18  from sklearn.model_selection import GridSearchCV, KFold
19
20  np.random.seed(222)
21
22  # %%
23  plt.rcParams['figure.figsize'] = (12,5)
24  plt.rcParams['figure.dpi'] = 80
25
26  # %%
27  d = pd.read_csv("dataQ2.csv")
28
29  # %%
30  df = d.groupby('V1').count()
31
32  clrs = ['grey' if (x < 1980) else '#3282bd' for x in df.index ]
33
34  plt.figure(figsize=(20,10))
35  plt.xticks(rotation=90)
36  # plot = sns.barplot(x=df.index, y = df.V3, color='#3282bd')
37  plot = sns.barplot(x=df.index, y = df.V3, palette=clrs)
38  plot.set_ylabel('Counts', fontsize=20, labelpad=20)
39  plot.set_xlabel('Year', fontsize=20, labelpad=20)
40  plt.savefig('Plots/song_year.pdf')
41  plt.show()
42
43  # %%
44  d['V1'] = np.where(d['V1'] >= 1980, 1,0)
45  d
46
47  # %% [markdown]
48  # ### LDA
49
50  # %%
51  y = d['V1']
52  X = d.loc[:,d.columns != 'V1']
53
54
55
56  # %%
57  TEST_SIZE = 0.3
58
59  # train/test split
60  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=TEST_SIZE)
61
62  # %%
63  sc = StandardScaler()
64  X_train = sc.fit_transform(X_train)
65  X_test = sc.transform(X_test)
66
67  # %%
68  lda = LDA(solver='svd', tol=1.e-4)
69  song_lda = lda.fit(X=X_train, y=y_train)
70
71  # %%
72  print(f"Prior probabilities of groups: \n {song_lda.priors_} \n")
```

```
73
74  # %%
75  lda_train = song_lda.predict(X_train)
76  lda_test = song_lda.predict(X_test)
77
78  # %%
79  cf_train = confusion_matrix(y_train, lda_train)
80  pd.DataFrame(cf_train, columns=["FALSE", "TRUE"], index=["FALSE", "TRUE"])
81
82  # %%
83  cf_test = confusion_matrix(y_test, lda_test)
84  pd.DataFrame(cf_test, columns=["FALSE", "TRUE"], index=["FALSE", "TRUE"])
85
86  # %%
87  from sklearn.metrics import roc_curve
88  fpr_1,tpr_1,thresh=roc_curve(y_test,lda.predict_proba(X_test)[:,-1])
89  auc_score=roc_auc_score(y_test,lda.predict_proba(X_test)[:,-1])
90
91  fpr_12,tpr_12,thresh2=roc_curve(y_train,lda.predict_proba(X_train)[:,-1])
92  auc_score2=roc_auc_score(y_train,lda.predict_proba(X_train)[:,-1])
93
94  print(f'AUC of Train: {auc_score2} \nAUC of Test: {auc_score}')
95
96  plt.plot(fpr_1,tpr_1, label='Test')
97  plt.plot(fpr_12, tpr_12, label='Train')
98  plt.xlabel('fpr1')
99  plt.ylabel('TPR')
100 plt.legend()
101
102 # %% [markdown]
103 # ## KNN
104
105 # %%
106 d2 = d.copy()
107 y = d2['V1']
108 X = d2.loc[:,d.columns != 'V1']
109
110 TEST_SIZE = 0.3
111
112 # train/test split
113 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=TEST_SIZE,
114
115 sc = StandardScaler()
116 X_train = sc.fit_transform(X_train)
117 X_test = sc.transform(X_test)
118
119
120 # X = sc.fit_transform(X)
121
122 # %%
123 KNeighborsClassifier().get_params().keys()
124
125 # %%
126 #trial:
127 ks = np.arange(5, 35)
128 ks = np.array([5,10,15,20,25,30,35,40,45,50])
129 gs = GridSearchCV(
130     make_pipeline(
131         StandardScaler(),
132         KNeighborsClassifier()
133         ),
134         param_grid={'kneighborsclassifier__n_neighbors': ks},
135         cv=KFold(5, shuffle=True, random_state=1234)
136 ).fit(X, y)
```

```
137
138
139  # %%
140  print(gs.best_params_)
141  test_score = gs.cv_results_['mean_test_score']
142
143  plt.figure(figsize=(12, 6))
144  plt.plot(ks, test_score, color='#3282bd', linestyle='dashed', marker='o',
145           markerfacecolor='black', markersize=10)
146  plt.xlabel('k Value')
147  plt.ylabel('Mean Accuracy')
148  plt.savefig('Plots/errorknn.pdf')
149  plt.show()
150
151  # %%
152  classifier = KNeighborsClassifier(n_neighbors=25)
153  classifier.fit(X_train, y_train)
154
155  y_pred = classifier.predict(X_test)
156
157
158  # %%
159  cf_test = confusion_matrix(y_test, y_pred)
160  pd.DataFrame(cf_test, columns=["FALSE", "TRUE"], index=["FALSE", "TRUE"])
161
162  # %%
163  (479+459)/(1478+459+479+1341)
164
165  # %%
166  fpr_2,tpr_2,thresh=roc_curve(y_test,classifier.predict_proba(X_test)[:,-1])
167  auc_score=roc_auc_score(y_test,classifier.predict_proba(X_test)[:,-1])
168
169  fpr_22,tpr_22,thresh2=roc_curve(y_train,classifier.predict_proba(X_train)[:,-1])
170  auc_score2=roc_auc_score(y_train,classifier.predict_proba(X_train)[:,-1])
171
172  print(f'AUC of Train: {auc_score2} \nAUC of Test: {auc_score}')
173
174  plt.plot(fpr_2,tpr_2, label='Test')
175  plt.plot(fpr_22, tpr_22, label='Train')
176  plt.legend()
177
178  # %% [markdown]
179  # ## Naive Bayes
180
181  # %%
182  d3 = d.copy()
183
184  y = d3['V1']
185  X = d3.loc[:,d3.columns != 'V1']
186
187  TEST_SIZE = 0.3 # smaller training set to reduce computational cost
188
189  # train/test split
190  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=TEST_SIZE)
191
192  sc = StandardScaler()
193  X_train = sc.fit_transform(X_train)
194  X_test = sc.transform(X_test)
195
196  # %%
197  from sklearn.naive_bayes import *
198
199  NB = GaussianNB()
200  # NB = MultinomialNB()
```

```
201  # NB = ComplementNB()
202
203  nb_fit = NB.fit(X_train, y_train)
204
205  y_pred = nb_fit.predict(X_test)
206
207  # %%
208  print("Number of mislabeled points out of a total %d points : %d"% (X_test.shape
209
210  # %%
211  cf_test = confusion_matrix(y_test, y_pred)
212  pd.DataFrame(cf_test, columns=["FALSE", "TRUE"], index=["FALSE", "TRUE"])
213
214  # %%
215  fpr_3,tpr_3,thresh=roc_curve(y_test,NB.predict_proba(X_test)[:,-1])
216  auc_score=roc_auc_score(y_test,NB.predict_proba(X_test)[:,-1])
217
218  fpr_32,tpr_32,thresh2=roc_curve(y_train,NB.predict_proba(X_train)[:,-1])
219  auc_score2=roc_auc_score(y_train,NB.predict_proba(X_train)[:,-1])
220
221  print(f'AUC of Train: {auc_score2} \nAUC of Test: {auc_score}')
222
223  plt.plot(fpr_3,tpr_3, label='Test')
224  plt.plot(fpr_32, tpr_32, label='Train')
225  plt.legend()
226
227  # %%
228  nb_fit.score(X_test, y_test)
229
230  # %%
231  t = np.linspace(0,1,50)
232  plt.plot(fpr_1,tpr_1, label='LDA')
233  # plt.plot(fpr_12, tpr_12, label='Train')
234  plt.plot(fpr_2,tpr_2, label='KNN')
235  # plt.plot(fpr_22, tpr_22, label='Train')
236  plt.plot(fpr_3,tpr_3, label='Naive Bayes')
237  # plt.plot(fpr_32, tpr_32, label='Train')
238  plt.plot(t, t, '--', alpha=0.3, color='grey')
239  plt.legend()
240  plt.xlabel('FPR')
241  plt.ylabel('TPR')
242  plt.savefig('Plots/together_plots.pdf')
243  plt.show()
244
245
246  # %%
247  t = np.linspace(0,1,50)
248  plt.subplot(131)
249  plt.plot(fpr_1,tpr_1, label='Test')
250  plt.plot(fpr_12, tpr_12, label='Train')
251  plt.plot(t, t, '--', alpha=0.3, color='grey')
252  plt.title('LDA')
253  plt.xlabel('FPR')
254  plt.ylabel('TPR')
255  plt.legend()
256  plt.subplot(132)
257  plt.plot(fpr_2,tpr_2, label='Test')
258  plt.plot(fpr_22, tpr_22, label='Train')
259  plt.plot(t, t, '--', alpha=0.3, color='grey')
260  plt.title('KNN')
261  plt.xlabel('FPR')
262  plt.legend()
263  plt.subplot(133)
264  plt.plot(fpr_3,tpr_3, label='Test')
```

```
265 plt.plot(fpr_32, tpr_32, label='Train')
266 plt.plot(t, t, '--', alpha=0.3, color='grey')
267 plt.title('Naive Bayes')
268 plt.xlabel('FPR')
269 plt.legend()
270 plt.savefig('Plots/allroc.pdf')
271 plt.show()
272
273 # %%
```