

COL774
Machine Learning
Assignment 2

Atif Anwer
2020EE10479

4 October 2022

Contents

1	Text Classification	3
1.1	Naive Bayes	3
1.2	Random Guessing and Constant Prediction	4
1.3	Confusion Matrices	4
1.4	Stop-word removal and Stemming	5
1.5	Feature Engineering	6
1.6	Precision, Recall, F1-score	6
2	Binary Image Classification	7
2.1	Linear Kernel SVM	7
2.2	Gaussian Kernel SVM	10
2.3	Scikit-learn SVM model	12
3	Multi-Class Image Classification	14
3.1	One-vs-One SVM model with CVXOPT	14
3.2	Scikit-learn One-vs-One SVM model	14
3.3	Confusion matrices and Misclassifications	14
3.4	Cross-Validation and Regularization	17

Libraries Used

- numpy
- scipy
- matplotlib
- nltk
- sklearn
- sys
- sys
- random
- pandas
- os
- collections
- re
- wordcloud
- cvxopt
- shutill
- time

1 Text Classification

1.1 Naive Bayes

Naive Bayes was implemented as a *multinomial* event model.

$$P(x|y = 1; \phi_{x_1}, \phi_{x_2}, \dots, \phi_{x_{|V|}}) = \frac{n!}{f_{x_1}! f_{x_2}! \dots f_{x_n}!} \phi_{x_1}^{f_{x_1}} \phi_{x_2}^{f_{x_2}} \dots \phi_{x_n}^{f_{x_n}}$$

Where x is a collection of n distinct words x_i appearing with frequency f_{x_i} .

The model was trained after removal of html tags and punctuations

- Training Accuracy : 91.616%
- Test Accuracy : 82.92%
- F1-Score : 0.8657373

Data Visualization

The word clouds for the training data were plotted for both positive and negative classes

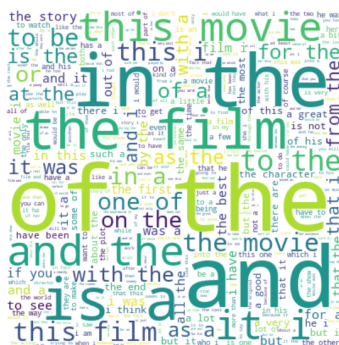


Figure 1: Word Cloud of Positive Reviews

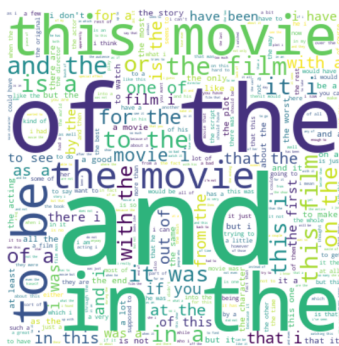


Figure 2: Word Cloud of Negative Reviews

1.2 Random Guessing and Constant Prediction

Random Guessing

Since the classification is binary, the probability of guessing the correct class is $1/2$. Hence, the expected accuracy is 50%. The observations were

- Test Accuracy : 50.0533333%
- F1-Score : 0.56957371

This is in compliance with the expected metrics and the F1-score quantizes it's poor performance.

Constant Prediction

If the prediction is set to constant, the test accuracy totally depends on the test dataset. The observations in our case were

- Test Accuracy : 66.66667%
- F1-Score : 0.8

The rise in accuracy with respect to the random model is because the test data is skewed towards the positive reviews. The Naive Bayes outperforms the above in accuracy by 16.25% and F1-Score by 0.0657373.

1.3 Confusion Matrices

Naive Bayes

$$\begin{pmatrix} 8260 & 822 \\ 1740 & 4178 \end{pmatrix}$$

Random Guessing

$$\begin{pmatrix} 4957 & 2449 \\ 5043 & 2551 \end{pmatrix}$$

Constant Prediction

$$\begin{pmatrix} 10000 & 5000 \\ 0 & 0 \end{pmatrix}$$

Observations

- It is observed that in the Naive Bayes model, classification split in each class ($TP/(TP+FN)$ and $TN/(TN+FP)$) are very close due to the same number of positive and negative examples in the training dataset.
- In case of Random Guessing, both classes are equally classified/misclassified bringing the accuracy close to 50%. The constant prediction model predicts positive for all test samples and hence misclassifies all samples with negative ground truth.
- All confusion matrices have high values in the first column, owing to the positive-skewed test dataset.

1.4 Stop-word removal and Stemming

The training and test data were cleaned off of Stop-words and stemming was performed using *PorterStemmer* from *nlTK*

The resulting word clouds for positive and negative classes are

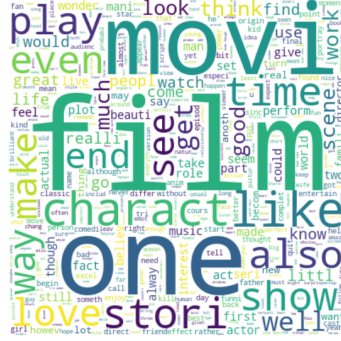


Figure 3: Word Cloud of Positive Reviews



Figure 4: Word Cloud of Negative Reviews

Learning the Naive Bayes model on the transformed data, the observations were

- Training Accuracy : 91.672
- Test Accuracy : 82.213
- F1-Score : 0.857555

The marginal drop in accuracy may be because the text reviews are not long enough to be affected by stop-word cleansing or due to the over-stemming nature of *nlTK*, resulting in losing information of different forms of a word expressing a different emotion.

1.5 Feature Engineering

Bigrams

Inclusion of Bigrams to the vocabulary as features with stemming and training the Naive Bayes model, results in the following

- Training Accuracy : 99.756%
- Test Accuracy : 84.78%
- F1-Score : 0.87864774

Trigrams + Bigrams

As a new feature, Trigrams were added to the vocabulary as features and training the Naive Bayes model, results in the following

- Training Accuracy : 99.988%
- Test Accuracy : 84.247%
- F1-Score : 0.87618548

Observations

- The addition of both Bigram and Trigram features have resulted in better performance on the test data. This is because of the information gained on the relative position of the words in the text review.
- Judging by the performance on the training and test data, it can be clearly seen that the inclusion of Trigrams results in overfitting and hence the drop in test-accuracy.

1.6 Precision, Recall, F1-score

The best model out of the above was Naive Bayes on the data with bigram feature engineering of stemmed words and exclusion of stop-words. The test data performance metrics were:

- Precision = 0.937819
- Recall = 0.8265
- F1-Score = 0.87864774

Considering the skewed nature of the test dataset, F1-score tends to a better measure of performance and hence, selecting a model by F1-score should be appropriate.

2 Binary Image Classification

Binary classification was performed for the classes 0 and 4 from the gives dataset.

2.1 Linear Kernel SVM

To utilize the *CVXOPT* package, the dual problem to be reduced to the following problem

$$\begin{aligned} \min & \left(\frac{1}{2} \alpha^T P \alpha + q^T \alpha \right) \\ & G \alpha \preceq H \\ & A \alpha = b \end{aligned}$$

The dual problem of the linear kernel SVM model is

$$\begin{aligned} & \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^i y^j (x^i)^T x^j - \sum_{i=1}^m \alpha_i \\ & 0 \leq \alpha_i \leq C \quad \forall i \in [1, m] \quad \text{and} \quad \sum_{i=1}^m \alpha_i y^i = 0 \end{aligned}$$

Defining the matrix Z such that $Z_{ij} = x_j^i y^i$, we can say

$$\alpha^T P \alpha = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j P_{ij}$$

Comparing $\frac{1}{2} \alpha^T P \alpha$ with the first term of the objective function,

$$\begin{aligned} P_{ij} &= (x_j^i y^i)^T (x_j^i y^i) \\ P &= Z Z^T \end{aligned}$$

Hence, $q^T \alpha$ represents the second term of the objective function, evaluating q to a m -size vector with $q_i = -1$ for all i

The inequality constraints can be represented as $2 * m$ inequalities

$$\begin{aligned} -\alpha_i &\leq 0 \quad \forall i \in [1, m] \\ \alpha_i &\leq C \quad \forall i \in [1, m] \end{aligned}$$

To represent the above, G and H matrices will evaluate as $2m \times m$ and $2m \times 1$ matrices respectively.

$$G = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \\ -1 & 0 & \dots & 0 \\ 0 & -1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & -1 \end{pmatrix}$$

$$H = \begin{pmatrix} c \\ c \\ \vdots \\ c \\ 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

The equality constraint $\sum_{i=1}^m \alpha_i y^i = 0$ can be easily represented using $A\alpha = b$. This give matrices A and b

$$A = (y^1 \quad y^2 \quad \dots \quad y^m)$$

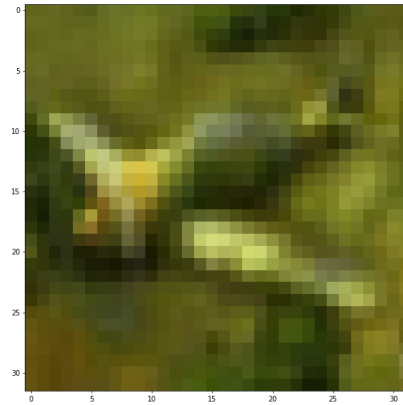
$$b = 0$$

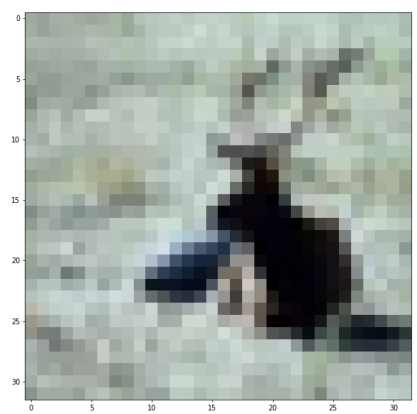
Feeding the matrices to the convex optimization algorithm for $C = 1.0$, the results were obtained as follows

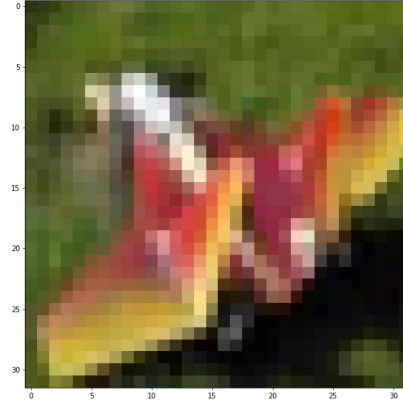
- No. of Support Vectors = 1505, which constitutes 37.625% of training sample
- $w = [0.40311351 \ 0.09715971 \ 0.99773958 \ \dots \ 0.48383299 \ -0.04012729 \ 0.53011341]$
- $b = 0.12942423232064204$
- Test Accuracy : 70.8%

Visualizing the Support Vectors

The top-5 support vectors(by value of α) were reshaped and plotted as an RGB image.







2.2 Gaussian Kernel SVM

The new dual problem with Gaussian Kernels will be

$$\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^i y^j \phi(x^i)^T \phi(x^j) - \sum_{i=1}^m \alpha_i \quad (1)$$

Where,

$$\phi(x^i)^T \phi(x^j) = e^{-\gamma \|x^i - x^j\|^2}$$

With the constraint being the same, change is only observed in the P matrix.

$$P_{ij} = y^i y^j e^{-\gamma \|x^i - x^j\|^2}$$

$$P_{ij} = y^i y^j e^{-\gamma (\|x^i\|^2 + \|x^j\|^2 - 2x^i x^j T)}$$

The above expression can hence be vectorized and calculated efficiently. Since we can't explicitly store w since it is an infinite-dimensional vector, we use the entire set of support vectors to make predictions on the test dataset.

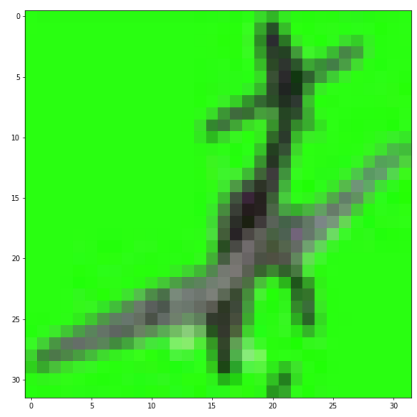
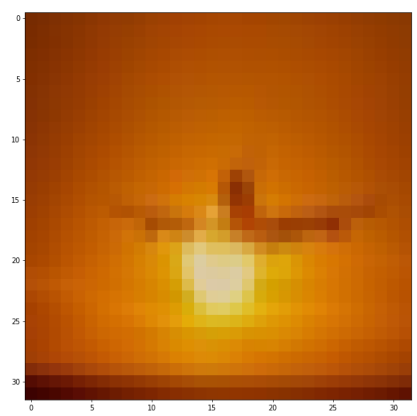
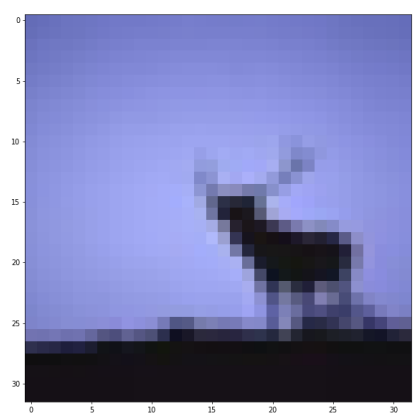
Feeding the matrices to the convex optimization algorithm with $C = 1.0$ and $\gamma = 0.001$, the results were obtained as follows

- No. of Support Vectors = 1754, which constitutes 43.85% of training samples
- Test Accuracy : 85.6%

The Gaussian Kernel obviously performs better than a linear kernel considering the it's capability to learn a lot more features than a linear kernel(which is limited to the number of features in a training sample).

Visualizing the Support Vectors

The top-5 support vectors(by value of α) were reshaped and plotted as an RGB image.





2.3 Scikit-learn SVM model

The SVM model from Scikit-learn was used to perform the binary classification task as above with Linear and Gaussian Kernels with the same hyperparameters.

Linear Kernel

- No. of Support Vectors = 1495, which constitutes 37.375% of training samples
- Test Accuracy : 79.100%
- $w = [0.40331326 \ 0.09731363 \ 0.99793737 \ \dots \ 0.48403465 \ -0.04035304 \ 0.53001174]$
- $b = 0.12816922455336766$

Observations

- The number of support vectors dropped marginally when compare to the model learnt using CVXOPT(nSV = 1505) with 1495 support vectors common between the two models (Hence,

the support vectors here are a subset of the ones in CVXOPT)

- It can be noted that the values of w and b closely match with the ones obtained with CVXOPT (The sign may reverse if the assignment of value 1 and -1 were switched for classes in CVXOPT.)

Gaussian Kernel

- No. of Support Vectors = 1743, which constitutes 43.575% of training samples
- Test Accuracy : 85.8%

The number of support vectors dropped marginally when compare to the model learnt using CVXOPT($nSV = 1754$) with 1743 support vectors common between the two models (Hence, the support vectors here are a subset of the ones in CVXOPT)

Computational Cost (Training)

SVM Learning Algorithm	Linear Kernel	Gaussian Kernel
<i>CVXOPT quadratic solver</i>	23.250s	29.385s
<i>Scikit-learn</i>	19.499s	10.851s

Judging by the above data, it is clear that the optimization algorithm used by Scikit-learn outperforms the *CVXOPT*'s quadratic solver. Also, for the give training data, CVXOPT seems to take more time in optimizing with a gaussian kernel than the linear kernel, which is the opposite with Scikit-learn.

3 Multi-Class Image Classification

The data(A subset of CIFAR-10) consisted of 5 classes, Airplane(0), Automobile(1), Bird(2), Cat(3) and Deer(4).

3.1 One-vs-One SVM model with CVXOPT

For the given training data with 5 classes, ${}^5C_2 = 10$ binary classifiers were trained. For testing, majority voting was considered among the 10 classifiers.

- Test Accuracy : 56.8799%
- Training time : 241.4825s

3.2 Scikit-learn One-vs-One SVM model

Using the one-vs-one SVM model by Scikit-learn, the results were

- Test Accuracy : 59.3%
- Training time : 157.6044s

Observations

- The test accuracy has improved slightly with Scikit-learn and hence outperforms the CVXOPT's quadratic solver optimization algorithm.
- The computational cost in training on the given training data has reduced drastically in case of Scikit-learn's SVM model.

3.3 Confusion matrices and Misclassifications

One-vs-One SVM model with CVXOPT

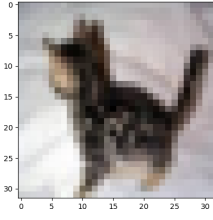
The confusion matrix on the test data was observed

$$\begin{pmatrix} 567 & 24 & 38 & 20 & 47 \\ 90 & 459 & 134 & 287 & 105 \\ 175 & 207 & 661 & 192 & 62 \\ 29 & 165 & 69 & 418 & 47 \\ 139 & 145 & 97 & 83 & 739 \end{pmatrix}$$

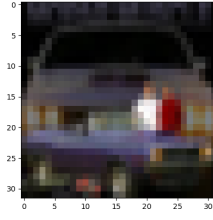
The misclassifications by the model were randomly picked and reshaped for visualization.

Observations

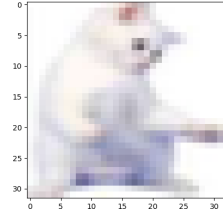
- Correct classifications can be read along the left diagonal of the confusion matrix, with the highest value being for the Deer(4) class.
- The highest misclassifications are observed in Automobile(1) and Cat(3) classes. It can also be seen that the model misclassifies too many Automobiles(1) as Bird(2) or Cat(3) and too many Cats(3) as Automobile(1) or Bird(2). The Misclassifications in the figure below are hence justified by the confusion matrix.



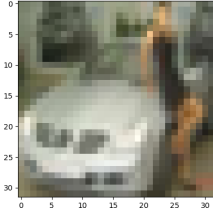
(a) Predicted : Bird(2)



(b) Predicted : Cat(3)



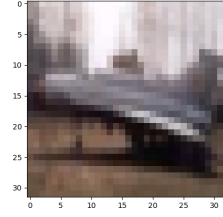
(c) Predicted : Airplane(0)



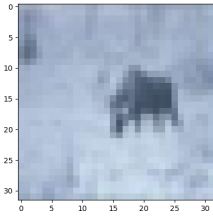
(d) Predicted : Cat(3)



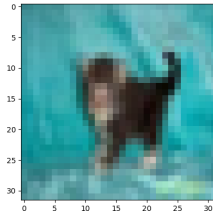
(e) Predicted : Airplane(0)



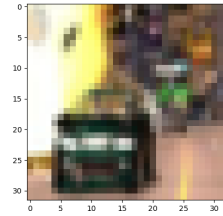
(f) Predicted : Automobile(1)



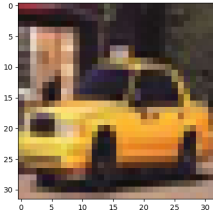
(g) Predicted : Airplane(0)



(h) Predicted : Bird(2)



(i) Predicted : Cat(3)



(j) Predicted : Cat(3)

Figure 5: Misclassifications by the model

One-vs-One SVM model with Scikit-Learn

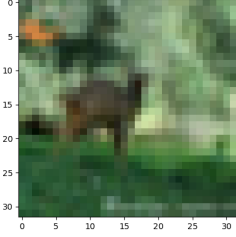
The confusion matrix on the test data was observed

$$\begin{pmatrix} 731 & 61 & 97 & 48 & 83 \\ 42 & 409 & 123 & 212 & 78 \\ 92 & 133 & 572 & 118 & 61 \\ 35 & 252 & 126 & 524 & 49 \\ 100 & 145 & 82 & 98 & 729 \end{pmatrix}$$

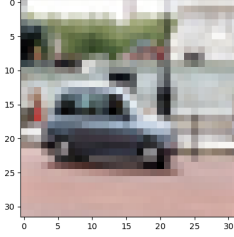
The misclassifications by the model were randomly picked and reshaped for visualization.

Observations

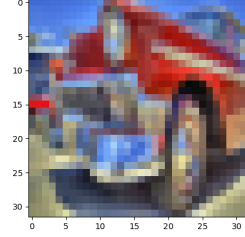
- Correct classifications can be read along the left diagonal of the confusion matrix, with the highest value being for the Airplane(0) and Deer(4) classes.
- The highest misclassifications are observed in Automobile(1) class. It can also be seen that the model misclassifies too many Automobiles(1) as Bird(2), Cat(3) or Deer(4). The Misclassifications in the figure below are hence justified by the confusion matrix.



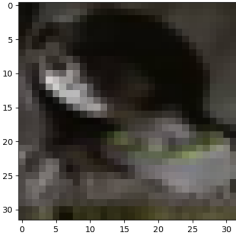
(a) Predicted : Bird(2)



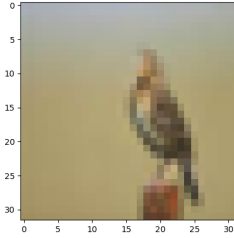
(b) Predicted : Airplane(0)



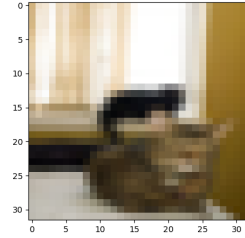
(c) Predicted : Cat(3)



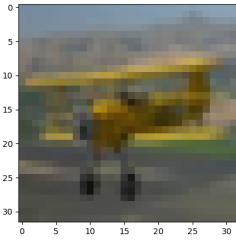
(d) Predicted : Deer(4)



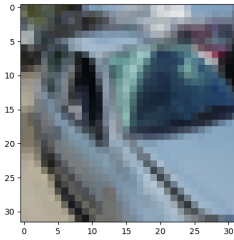
(e) Predicted : Deer(4)



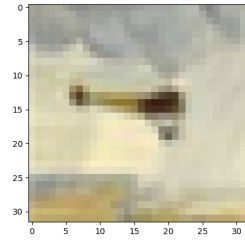
(f) Predicted : Automobile(1)



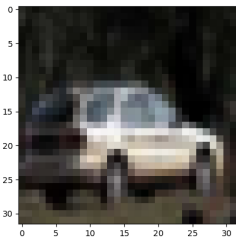
(g) Predicted : Deer(4)



(h) Predicted : Bird(2)



(i) Predicted : Bird(2)

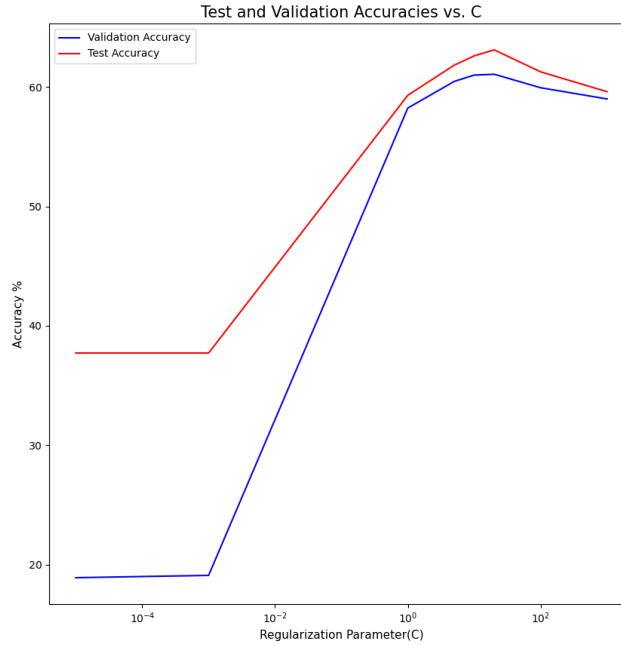


(j) Predicted : Cat(3)

3.4 Cross-Validation and Regularization

5-Fold Cross Validation was implemented with Sckit-learn's SVM model and Validation and Test Accuracies were noted for different values of the regularization parameter C .

C	Validation Accuracy	Test Accuracy
10^{-5}	18.91%	37.72
10^{-3}	19.10%	37.72%
1	58.23%	59.30%
5	60.46%	61.82%
10	60.99%	62.60%
20	61.06%	63.10%
100	59.93%	61.26%
1000	58.99%	59.60%



- The validation accuracy peaks at $C = 20$, which also gives the highest accuracy in the test data. This tells us that the cross-validation method is a good estimate to determine the hyperparameters for the model that could possibly perform the best in testing.
- From the primal problem of the SVM, we can deduce that for lower values of C , the model tends to have a high relaxation on ξ_i and hence under-fits on the training data, resulting in low validation and test accuracies. On the other side, for high values of C , the model has a high penalty on the values of ξ_i and results in a model that over-fits on the training data, and again dropping in accuracy. The model peaks in performance between under-fitting and over-fitting regions. This is all in compliance with the above experimental observations.