

EECS445 Project 1

Dong Zhengyuan (dongzy)

February 10, 2019

2 Feature Extraction

(c) The number of unique words is 2850, The number of average non-zero features is 15.624

3 Hyperparameter Feature Selection

3.1 Hyperparameter Selection for a Linear-Kernel SVM

(a) In an iteration of cross-validation, we take one fold as the validation set and the rest folds as the training set. Therefore, if we maintain class proportions across folds, we will have a balanced training set in each iteration of the cross-validation and do not need to compensate for class imbalances.

(c) The optimal parameter under different metrics is shown below:

Performance Metric	C	Performance
Accuracy	0.1	0.839
F1-Score	0.1	0.838
AUROC	0.1	0.920
Precision	10	0.841
Sensitivity	0.001	0.864
Specificity	10	0.844

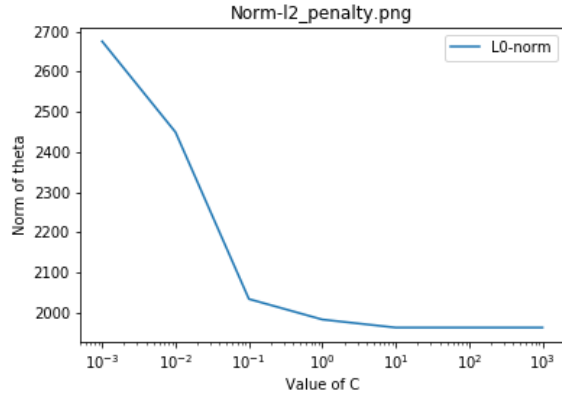
Generally, as C increases, CV performance increases to the peak, then gently decreases. This trend is shared among metrics, except for sensitivity, where the peak is reached at the beginning when $C = 0.001$.

If I have to train a final model, I will optimize for accuracy when choosing C , since the recognition of both positive and negative classes are important, and accuracy is the most direct metric to address that importance.

(d) Using $C = 0.1$ which maximizes accuracy, the performance of the SVM under different metrics is shown below:

Performance Metric	Performance
Accuracy	0.833
F1-Score	0.830
AUROC	0.921
Precision	0.845
Sensitivity	0.815
Specificity	0.850

(e) The plot is shown below:



We can see that the L0-norm of $\bar{\theta}$ decreases as C increases, but gradually converges to 2000.

(f) The most positive / negative words are shown below:

Positive Coefficient	Word	Negative Coefficient	Word
0.969	thanks	-0.615	hours
0.901	thank	-0.549	delayed
0.765	great	-0.521	due
0.596	good	-0.507	good

3.2 Hyperparameter Selection for a Quadratic-Kernel SVM

(a) The optimal parameters under different metrics are shown below:

Grid search:

Performance Metric	C	r	Performance
Accuracy	10	1000	0.840
F1-Score	10	1000	0.839
AUROC	1000	0.1	0.918
Precision	10	1000	0.846
Sensitivity	0.001	0.001	0.948
Specificity	10	100	0.848

Random search:

Performance Metric	C	r	Performance
Accuracy	136.61	5.077	0.848
F1-Score	136.61	5.077	0.846
AUROC	136.61	5.077	0.916
Precision	136.61	5.077	0.852
Sensitivity	0.0014	0.0033	0.918
Specificity	136.61	5.077	0.854

(b) The AUROC results are shown below:

Tuning Scheme	C	r	AUROC
Grid Search	1000	0.1	0.918
Random Search	136.61	5.077	0.916

Generally, for a given C , the performance increases as r increases, for a given r , the performance increases and then decreases as C increases. The use of random search is generally better than grid search, as

it visits more distinct values of C and r . If C and r do not contribute equally to CV-performance (which is usually the case), random search is more likely to get closer to the optimal value of the predominant hyperparameter. Random search also gives more consistent results among metrics (5 of the 6 metrics gives the same result)

3.3 Learning Non-linear Classifiers with a Linear-Kernel SVM

(a) The quadratic kernel is $K(\bar{x}, \bar{x}') = (\bar{x} \cdot \bar{x}' + r)^2$. Suppose $\bar{x}, \bar{x}' \in \mathbb{R}^d$, we can expand it as

$$\begin{aligned} K(\bar{x}, \bar{x}') &= \left(\sum_{i=1}^d x_i x'_i + r \right)^2 = \left(\sum_{i=1}^d x_i x'_i \right)^2 + 2 \left(\sum_{i=1}^d x_i x'_i \right) r + r^2 \\ &= \sum_{i=1}^d (x_i x'_i)^2 + 2 \sum_{i=1}^{d-1} \sum_{j=i+1}^d (x_i x'_i)(x_j x'_j) + 2 \left(\sum_{i=1}^d x_i x'_i \right) r + r^2 \\ &= \sum_{i=1}^d x_i^2 x_i'^2 + 2 \sum_{i=1}^{d-1} \sum_{j=i+1}^d (x_i x_j)(x'_i x'_j) + 2r \left(\sum_{i=1}^d x_i x'_i \right) + r^2 \\ &= \phi(\bar{x}) \phi(\bar{x}') \end{aligned}$$

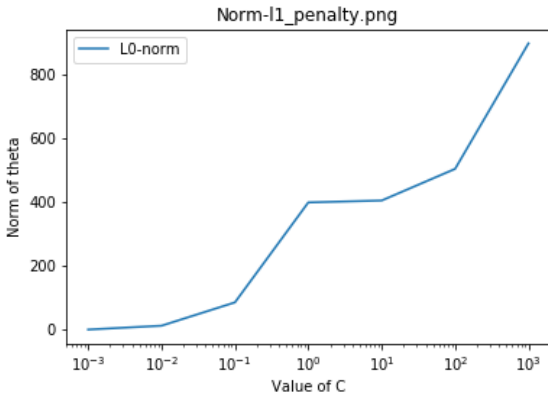
Therefore $\phi(x) = [(x_i^2)_{i=1..d}, (\sqrt{2}x_i x_j)_{i=1..d-1, j=i+1..d}, (\sqrt{2r}x_i)_{i=1..d}, r]^T$

(b) Explicit feature mappings are more flexible and always produces a valid kernel when multiplied, while kernel functions must satisfy Mercer's Theorem. However, explicit feature mappings can be difficult, sometimes impossible (rbf kernel) to compute, while kernel functions can be computed much faster.

3.4 Linear-Kernel SVM with L1 Penalty and Squared Hinge Loss

(a) The default max_iter=1000 is not enough for convergence, and the result is undeterministic. Among about 10 trials, AUROC varies between 0.91 and 0.92. In most cases $C = 1000$ is optimal, but sometimes $C = 100$. After modifying max_iter to 10000, the algorithm converges, but the result is still undeterministic.

(b) The plot is shown below:



(c) The L1 penalty significantly decreases the L0-norm of the learned parameter. Also, under L1 norm, the L0-norm of the learned parameter increases as C increases.

(d) If we use Hinge Loss instead of Squared Hinge Loss, the penalty on the misclassified data points will be lighter. More data points will become support vectors and have non-zero parameter. Therefore, the optimal solution will have a greater L0-norm under Hinge Loss.

4 Asymmetric Cost Functions and Class Imbalance

4.1 Arbitrary class weights

(a) If W_n is much greater than W_p , misclassified negative data points will be much more severely penalized than misclassified positive data points. Our model will be trained to emphasize more on correctly classifying negative data points.

(b) The performance of the modified SVM is shown below:

Performance Metric	Performance
Accuracy	0.563
F1-Score	0.222
AUROC	0.905
Precision	1.0
Sensitivity	0.125
Specificity	1.0

(c) Compared to 3.1(d), F1-score and sensitivity are affected the most. As stated in 4.1(a), the new class weights make the model emphasize more on classifying negative data points and lowers sensitivity, which is the metric for classifying positive data points. Since F1-score is a function of sensitivity, it gets affected as well.

4.2 Imbalanced data

(a) The performance of the SVM with $C = 0.01$, $W_n = W_p = 1$ is shown below:

Performance Metric	Performance
Accuracy	0.384
F1-Score	0.374
AUROC	0.911
Precision	1.0
Sensitivity	0.23
Specificity	1.0

(b) In the imbalanced data set in this question, there are more negative points than positive points. If we train the SVM with balanced class weights, we end up penalizing more on misclassified negative points and make the model better at classifying negative points, as indicated by specificity=1. However, the model performed poorly at classifying positive points, as indicated by a low sensitivity.

4.3 Choosing appropriate class weights

(a) Since there are more negative points than positive points, metrics that evaluates the classification of both classes tend to be dominated by the majority class. The F1-score is more informative in such situation since it focuses on classifying positive points and maintains a balance between sensitivity and precision, which helps mitigating the situation in 4.2.

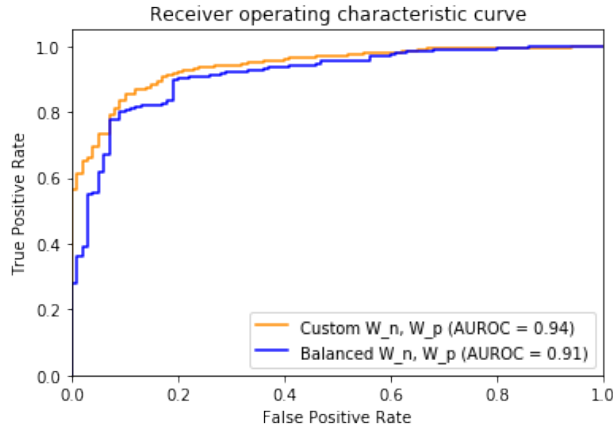
A coarse grid search is performed first to determine a reasonable, smaller range of class weights. After finding a smaller range, another random search is performed. Considering that in the imbalanced dataset here, the proportion between negative and positive points is about 3:1, $W_n : W_p$ is fixed inverse-proportionally to 1:3 to further reduce the effort wasted.

This approach gives its best results as $W_n = 0.212$, $W_p = 0.636$ ($C = 1$), which achieves F1-Score=0.777 in cross-validation. The corresponding test performances are shown below:

Performance Metric	Performance
Accuracy	0.858
F1-Score	0.905
AUROC	0.938
Precision	0.971
Sensitivity	0.848
Specificity	0.9

4.4 The ROC curve

(a) The ROC curves are shown below:



5 Challenge

For this challenge, I attempted various approaches to train the model, including those used in previous sections and some new approaches. Since we have three classes and we are equally interested in identifying each one, accuracy is a decent metric to evaluate the performance. Since the ground truth of the test set is unknown, I chose to optimize hyperparameters based on the accuracy score of 5-fold cross-validation.

Approach 1: Quadratic-kernel SVC, search for C and r

To start with, I used a quadratic-kernel SVC that is previously used. I first applied a grid search within range $\{10^{-3}, 10^{-2}, \dots, 10^2, 10^3\}$ for both C and r , and the best performance 0.715 was achieved when $C = 1$, $r = 1000$. Around that optimal value I conducted a second random search of 25 points, with $\lg C \in [-1, 1]$ and $\lg r \in [2, 4]$ uniformly distributed. The best performance 0.718 was achieved when $C = 0.206$, $r = 908.7$.

Approach 2: Linear SVC with L1 penalty, grid search for C

The quadratic SVC hyperparameter selection was time consuming, so I switched to linear SVC with L1 penalty. I searched C within range $\{10^{-3}, 10^{-2}, \dots, 10^2, 10^3\}$. The best performance 0.7076 was achieved when $C = 1$.

Approach 3: Linear SVC with L2 penalty, grid search for C

Still using a linear SVC, I switched to L2 penalty and searched for C within the same range. This time, the best performance 0.726 is achieved when $C = 0.1$.

Approach 4: Finer grid search with word counting and normalization

The results from linear SVC with L2 penalty seems better, so I decided to stick to that model and explore

the effect of feature engineering on the results. I modified the `generate_feature_matrix` and added two options: Counting the number of times a word occurs in a review and normalizing each row of the feature matrix. I also made a finer grid search by having 41 values evenly distributed between $\lg C = -2$ and $\lg C = 2$. It turns out that the two feature generating options does not improve the result of the model. The four combinations of these two options being On/Off gives the same best performance 0.731 when $C = 0.040$, the improvement mainly due to a finer grid. I think the reason it did not work is that compared to the dictionary size, the number of duplicate words in a comment is small, let alone these duplicate words are often neutral words without emotional inclination like "I", "and", "to", "at".

Approach 5: Quadratic-kernel SVC under OvO mode

Finally, I investigated the difference between OvO and OvR in a quadratic kernel SVC. In Approach 1 the `decision_function_shape` is set to OvR by default. Here I changed it to OvO and searched for C and r within the same range as used in Approach 1. The best performace 0.712 was achieved when $C = 1$, $r = 1000$, which is not very different from the result of approach 1.

Conclusion After these approaches, I decided to choose a linear SVC with L2 penalty as my final model, with $C = 0.040$

Code

```
#!/usr/bin/env python
```
