

Homework 2: Support Vector Machines

Anthony Andriano

University of Colorado Colorado Springs
aandrian@uccs.edu

Planning Relax Dataset —

<https://archive.ics.uci.edu/dataset/230/planning+relax>

Skin Segmentation Dataset —

<https://archive.ics.uci.edu/dataset/229/skin+segmentation>

Introduction

Support Vector Machines (SVMs) are supervised learning models that can be used for binary classification, multiclass classification, and regression. The fundamental principle of an SVM is to find a decision boundary between two or more classes, such that the width of the decision boundary is as wide as possible. The decision boundary is sometimes referred to as a hyperplane, which separates classes in as many dimensions as the dataset requires based on the number of features. It is possible to separate the features using non-linear methods using kernel functions with non-linear constructions. However, kernel functions and non-linear classification are beyond the scope of this paper.

A key term during the discussion and use of an SVM is support vector. Support vectors are the points that lie on or near the decision boundaries. Ideally, SVMs have linearly-separable data points as support vectors, but that is not always the case, which will be discussed in later sections.

Question 1: Linear Separability

If a linear function, which is called a hyperplane in higher dimensions, exists such that the examples of a dataset can be separated by the linear function with any errors, the examples are said to be linearly separable. In a classification problem, that means all data points from one class are on a different side of the linear function than the data points from the other class. The linear function is a line when the feature space has a cardinality of 2, and the linear function is a plane or hyperplane when the cardinality of the feature space is 3 or more, respectively.

To determine if the data points are in one class or another, a discriminator function is used. A linear discriminator is a linear function, which is used to predict the class of a label according to the following definitions:

$$h_{\vec{\theta}}(x) = \theta_0 + \theta_1 x \quad (1)$$

The bias term, θ_0 , is typically included in the parameter

vector, which means the equation becomes:

$$h_{\vec{\theta}}(x) = \vec{\theta} \cdot \vec{x} \quad (2)$$

where:

- $x \in R^n$: the feature vector (input example)
- $\theta \in R^n$: the parameter or weight vector
- $y \in \{0, 1\}$: the binary class label

For a dataset to be linearly separable, there must exist some $\vec{\theta}$ such that:

- $\vec{\theta} \cdot \vec{x} \geq 1$ for data points in class 1
- $\vec{\theta} \cdot \vec{x} \leq -1$ for data points in class 0

Question 2: Margin of Separation

When the classes are linearly separable, a margin of separation is said to exist between the data points in each class. The boundaries defined by the discriminator are called the upper and lower margins. The goal is to find the widest distance between the upper and lower margins to maximize the classification and prediction capabilities of the SVM. Variations in the parameter vector are minimized when the margin is the widest it can be. The data points from each class on or closest to the margins are called support vectors.

The linear discriminator in a binary classification problem is sometimes called a binary discriminator. Equation 2 defines the decision boundary when the value is 0, or $\vec{w} \cdot \vec{x} = 0$. The upper boundary is defined as $\vec{\theta} \cdot \vec{x} = 0$. The upper boundary is set to $\vec{\theta} \cdot \vec{x} \geq 1$ while the lower boundary is set to $\vec{\theta} \cdot \vec{x} \leq -1$.

The discriminator and both margins form the margin of separation, which determines the classification of the data points that are evaluated using the SVM.

Width of the Margin

The width of the margin is the perpendicular distance between the upper and lower margins. The discriminator is halfway between each margin. To maximize the width of the margin, the minimum of $\|\theta\|$ must be calculated. For mathematical convenience, it is often better to minimize $\|\theta^2\|$ instead.

Question 3: Regularized Loss

It is ideal to have linearly separable data, but that is not always possible. To accurately classify data that is not perfectly linearly separable, an objective function with a loss component must be used.

The loss function must be simple enough to implement with a reasonable amount of resources and run time, while also fitting well to the data. In class, we primarily discussed hinge loss, which is a type of loss function that is piecewise differentiable. The non-differentiable point creates problems with the mathematical analysis of the model, but that issue can be worked around by picking a value for the function at the non-differentiable point. It is not as polished as the loss functions that are continuously differentiable, but it is easier to implement and the results are generally considered to be good enough for robust classification.

The first component of the loss function is the regularization component: $\frac{1}{2}\|\theta\|^2$. This component creates the widest possible margin between the classes. The second component of the loss function is the hinge loss equation, which is $\max(0, 1 - y_i(\vec{\theta} \cdot \vec{x}_i))$. The hinge loss equation must be summed over the dataset. The point of the hinge loss function is to penalize the output if a value exists within the margin or if the value has been misclassified. To control the effect of both components, a hyperparameter called C is often used. The hyperparameter balances between the width of the margin and the penalties due to misclassifications.

The combined equation is as follow:

$$\frac{1}{2}\|\theta\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i(\vec{\theta} \cdot \vec{x}_i)) \quad (3)$$

The motivation for the combined loss function is to allow data from imperfect datasets to be correctly classified with as much confidence as possible. Real datasets are often noisy with linearly inseparable classes, but the majority of the data points in each class could still be linearly separable, which would make an SVM with a linear classifier useful.

Question 4: Stochastic Gradient Descent

To solve the SVM classification problem, we are using Stochastic Gradient Descent (SGD). The goal of SGD is to minimize equation 3, which is the equation that combines regularization and loss. To compute the gradient of the equation, partial derivatives must be calculated.

Using equation 2, which is the linear function that is used as the discriminator, determine the gradient in the continuous regions. When $x < 1$, the gradient is -1. When $x > 1$, the gradient is 0. When $x = 1$, the gradient is $[-1, 0]$.

The algorithm for SGD includes a learning rate, which is another hyperparameter that controls the weight the parameter vector. First assign random initial values to the parameter vector, and then repeat the following calculation for each data point in the dataset:

$$\vec{\theta} \leftarrow \vec{\theta} - \alpha \nabla \mathcal{L} \quad (4)$$

That process is continued until the output converges, which means the value of consecutive calculations becomes very small, or the number of iterations has exceeded a

Question 5: Datasets

The **Planning Relax** and **Skin Segmentation** datasets were downloaded from the UCI repository for use in this assignment. The datasets contain all numeric features, which makes binary classification straightforward to perform.

Question 6: Custom SVM Model

A custom SVM was written in Python for this assignment. The custom SVM uses the SGD algorithm to optimize the parameter vector by changing the hyperparameter, C , and the learning rate. To measure the performance of the SVM, the classification accuracy was measured as a percentage. The runtime of each iteration was also recorded for information purposes only. Each of the datasets from Question 5 were analyzed with the custom SVM as seen in the following subsections.

Relax Dataset

The results produced by the custom SVM model when using the Relax dataset were minimally affected by changes to the hyperparameter, C , the learning parameter, l . Modifying the learning rate resulted in trivial changes to the accuracy score, F1 score, precision score, and recall score of the model. The hyperparameter does not appear to modify the outputs, which means the combination of the model and dataset is not sensitive to errors in the margin. Visualizing the data shows the separation, which means the margin is well defined and linearly separable. Notably, the custom SVM produced different accuracies as the library-based SVMs, which will be detailed in the next sections.

The configurations and results for the Relax dataset are as follows:

Acc %	F1 %	Prec %	Rec %	C,LR
71.429	83.33	71.43	100.00	1E-4,1E-1
71.270	82.54	70.27	100.00	1E-4,1E-5
71.429	83.33	71.43	100.00	1E-5,1E-1
71.270	82.54	70.27	100.00	1E-5,1E-5

Skin Dataset

The results produced by the custom SVM model when using the Skin dataset were primarily affected by changes to the hyperparameter, C . Modifying the learning rate resulted in trivial changes to the accuracy score, F1 score, precision score, and recall score of the model. The trade between precision and recall is more pronounced when the hyperparameter value is increased, which shows the need to tune the model to suppress false positives and negatives. In this case, recall is more important because false negatives are more detrimental to the diagnosis of cancer. Notably, the custom SVM produced different accuracies as the library-based SVMs, which will be detailed in the next sections.

The configurations and results for the Skin dataset are as follows:

Acc %	F1 %	Prec %	Rec %	C,LR
93.810	95.999	98.518	93.605	1E-4,1E-2
94.411	94.411	95.071	93.761	1E-4,1E-1
88.476	88.476	79.333	100.00	3E0,1E-2
88.476	87.521	77.215	100.00	3E9,1E-1

Question 7: Library-Based SVM Models

For this part of the assignment, the following libraries were used:

- SVC from sklearn.svm
- LinearSVC from sklearn.svm
- SGDClassifier from sklearn.linear_model

To measure the performance of the library-based SVMs, the classification accuracy, f1 score, precision, and recall score were measured as percentages. The runtime of each iteration was also recorded for information purposes only. Each of the datasets from Question 5 were analyzed with the library-based SVMs as seen in the following subsections. The same library implementations were also used for question 8.

SVC

Relax Dataset The SVC library was optimized using the hyperparameter, C. As seen in the following table, the results approximately match the results of the custom SVM model from question 6 when using the Relax dataset. As an additional experiment, the learning rate parameter was also changed, but the effect was the same - no change in the output. The limitation to the accuracy and precision are likely related to the lack of total linear separability of the data.

Model	Acc %	F1 %	Prec %	Rec %
SVC	70.270	82.540	70.270	100.000

Skin Dataset The SVC library was optimized using the hyperparameter, C. As seen in the following table, the results approximately match the results of the custom SVM model from question 6 when using the Skin dataset. The minor increases in performance are likely attributable to the C-based implementation, which tends to be more accurate due to better control of the data structures and type conversions. Compared to the other library-based SVM models, the SVC model produced the same accuracy within 1%.

Model	Acc %	F1 %	Prec %	Rec %
SVC	92.873	95.380	98.436	92.508

LinearSVC

Relax Dataset The LinearSVC library was optimized using the hyperparameter, C. As seen in the following table, the classification accuracy matches the results of the custom SVM from question 6 when using the Relax dataset. As an additional experiment, the loss function parameter was also changed, but the effect was the same - no change in the output. Because changing C and the loss function appeared to have no impact on the accuracy, the limitation of the accuracy is likely due to the nature of the data.

Model	Acc %	F1 %	Prec %	Rec %
Linear SVC	70.270	82.540	70.270	100.000

Skin Dataset The LinearSVC library was optimized using the hyperparameter, C. As seen in the following table, the results approximately match the results of the custom SVM model from question 6 when using the Skin dataset.

As with the SVC model, the minor increases in performance are likely attributable to the C-based implementation, which tends to be more accurate due to better control of the data structures and type conversions. Compared to the other library-based SVM models, the LinearSVC model produced the same accuracy within 1%.

Model	Acc %	F1 %	Prec %	Rec %
Linear SVC	92.873	95.380	98.436	92.508

SGDClassifier

Relax Dataset The SGDClassifier library was optimized using the regularization parameter, alpha. The best results of the experiments with SGD are shown in the table below. When the regularization parameter was weak, the results were worse than the custom SVM model as well as the other library-based SVM models. As an additional experiment, the loss function parameter was also changed, but the effect was the same - no change in the output. Because changing the loss function appeared to have no impact on the results, the limitation of the maximum accuracy, precision, f1, and recall is likely due to the nature of the data. The accuracy being lower when the regularization parameter is reduced shows a need for regularization to prevent poor fit of the model to the data. Penalizing larger weights prevents overfitting, which seems to be required in this scenario. A recall score of 1.0 is a good trade-off in this case due to the nature of the situation in which the data was collected.

Model	Acc %	F1 %	Prec %	Rec %
SGD	72.973	83.871	72.222	100.000

Skin Dataset The SGDClassifier library was optimized using the regularization parameter, alpha. As seen in the following table, the classification accuracy exceeds the results of the custom SVM from question 6 when using the Skin dataset. Unlike the results from the Relax dataset, this dataset performs well regardless of the value of the regularization parameter, which suggests a difference in the data and associated weights. As an additional experiment, the loss function parameter was also changed, but there was no effect on the output. Because changing the loss function appeared to have no impact on the accuracy, the limitation of the results is likely due to the nature of the data. With this dataset, overfitting does not appear to be a problem. In this case, a maximum recall score of 1.0 was not achieved, which is worse than the custom SVM model's performance; a high recall score is advantageous for the diagnosis of skin cancer because false positives are acceptable, while false negatives are not.

Model	Acc %	F1 %	Prec %	Rec %
SGD	93.338	95.682	98.734	92.814

Summary

Relax Dataset The maximum accuracy of each SVM classification library when analyzing the Relax dataset was the same as the custom SVM: 70.27%. However, library-based implementations completed 80 to 100 times faster due to being implemented in a lower-level language than Python.

Skin Dataset The maximum accuracy of each SVM classification library when analyzing the Skin dataset was significantly higher than the custom SVM: approximately 92% versus approximately 20%. In addition, library-based implementations completed 80 to 100 times faster due to being implemented in a lower-level language than Python.

Question 8: Other Classifiers

The last question for the assignment requires comparisons of the custom SVM, library-based SVMs, and other classifiers.

Relax Dataset

The list of algorithms and the associated results that were compared using the Relax dataset is as follows:

Model	Acc %	F1 %	Prec %	Rec %
SGD	72.973	83.871	72.222	100.000
Dec Tree	72.973	83.871	72.222	100.000
Custom SVM	70.270	82.540	70.270	100.000
SVC	70.270	82.540	70.270	100.000
Linear SVC	70.270	82.540	70.270	100.000
Rndm Frst	70.270	82.540	70.270	100.000
AdaBoost	70.270	82.540	70.270	100.000
XGBoost	70.270	82.540	70.270	100.000
Gauss Pro	70.270	82.540	70.270	100.000
Gauss NB	70.270	82.540	70.270	100.000
LightGBM	70.270	82.540	70.270	100.000
CatBoost	70.270	82.540	70.270	100.000
MLP	70.270	82.540	70.270	100.000
KNN	67.568	80.645	69.444	96.154
Quad Dis	67.568	80.000	70.588	92.308

Skin Dataset

The list of algorithms and the associated results that were compared using the Relax dataset is as follows:

Model	Acc %	F1 %	Prec %	Rec %
XGBoost	99.949	79.481	79.481	79.481
KNN	99.947	99.967	99.992	99.941
Rndm Frst	97.688	98.533	99.446	97.637
Dec Tree	97.380	98.345	98.788	97.906
Quad Dis	96.493	97.790	98.005	97.575
Custom SVM	93.810	95.999	98.518	93.605
AdaBoost	93.826	96.114	96.215	96.013
SGD	93.338	95.682	98.734	92.814
SVC	92.873	95.380	98.436	92.508
Linear SVC	92.873	95.380	98.436	92.508
Gaus NB	92.483	95.365	93.571	97.229
LightGBM	99.943	99.964	99.982	99.946
CatBoost	99.916	99.947	99.992	99.903
MLP	99.839	99.899	99.987	99.810

Analysis

Using the Relax dataset, the custom SVM model performed as well as almost every library-based model. Experimentation with various parameters for each model showed almost no performance increase; there were performance decreases,

though, such as the SGD model having a minimum accuracy of 35% when the regularization parameter is too low. The highest accuracy was achieved with the Decision Tree model; I expected XGBoost or one of the other boosted models to achieve higher accuracy. It may be that the parameters were not configured correctly, which led to lower accuracies.

Using the Skin dataset, the custom SVM model performed significantly worse than every other model. I spent hours debugging and tweaking the model to no avail. A maximum accuracy of 20.475% with the custom SVM model compared to a maximum of 99.951% with LightGBM is an inexplicably larger difference. The lowest accuracy of the library-based models was AdaBoost at 89.849%, which is still more than four times higher than the accuracy of the custom SVM model. I am guessing there is a bug in the code that is yet to be identified.

Future Work

The custom SVM model could be improved by implementing an algorithm in C instead of Python, which would dramatically increase the iteration speed, and fixing the objective function to perform better with additional datasets.

References

- "Implementing SVM and Kernal SVM with Python's Scikit-Learn". 2023. Implementing SVM and Kernal SVM with Python's Scikit-Learn. Accessed: March 29, 2025.
- "Support Vector Machines". 2025. 1.4 Support Vector Machines. Accessed: March 30, 2025.
- "Support Vector Machines with Scikit-learn Tutorial". 2019. Support Vector Machines with Scikit-learn Tutorial. Accessed: March 30, 2025.
- "SVM Python Without Library". 2019. svm-python-without-library. Accessed: April 5, 2025.

Appendix

To support this effort, significant effort and time were invested into visualizing various datasets and algorithms. Shown below are examples results using the algorithms in the "Question 8: Other Classifiers" section of the paper. The rows are different configurations of input data. From top to bottom:

- make_moons - Data points that are in the shape of partial moons, which resemble crescents
- make_circles - Data points in circular patterns with varying diameters
- make_blobs - Data points in groups that do not overlap
- make_gaussian_quantiles - Data points in densely nested concentric spheres with increasing radius
- make_classification - Data points that are normally distributed about the vertices of an n-dimensional hypercube

Here is an example of the output of the script that is responsible for comparing randomized data using many sci-kit algorithms:

