# OTDM Constrained Optimization Lab Assignment: Support Vector Machines in AMPL

Mariem Gandouz, Egon Ferri

21th of NOVEMBER 2019

# Contents

# 1    Introduction

This report presents the Support Vector Machine (SVM) learning algorithm. We will need to first talk about margins and the idea of separating data with a large gap. Next, we will talk about the optimal margin classifier, which will lead us to talk about the duality for the SVM. We will also implement SVMs in AMPL on random generated data with primal and dual formulations, and finally, we will close off our report with a classic iris dataset example analyse.

# 2    Modelling SVMs

## 2.1    Non-linear separable data, Soft-Margin SVM

In general, the training data can rarely be separated linearly because they are often noisy. An outlier in the training data, caused by noise, for example an observation with a wrong class label, could clearly influence the position of the hyperplane. To avoid this problem, it is possible to use a slightly modified variant of the Hard-Margin SVM, that can tolerate a certain number of outliers. The goal is to find a hyperplane that allows some misclassifications. For this purpose, the constraint in the linearly separable case $y_i(w^\top x_i + \gamma) \geq 1 (*)$[1] is eased, that's why the resulting SVM is referred as Soft-Margin SVM. Furthermore, the slacks $s_i \geq 0$ are inserted for each observation $x_i$. Each slack indicates the deviation of an observation $x_i$ from the separating hyperplane and can be regarded as a penalty for this observation if $x_i$ does not satisfy the previous constraint (*). Specifically, these are given by:

1) $s_i = 0$, for correctly classified observations, i.e. if

$$y_i(w^\top x_i + \gamma) \geq 1$$

2) $0 < s_i \leq 1$, for observations within the separation area of the hyperplane that are on the correct side of the hyperplane,i.e. if

$$0 \leq y_i(w^\top x_i + \gamma) < 1$$

3) $s_i > 1$, for misclassified observations, i.e. if the observations are on the wrong side of the hyperplane and thus
$$y_i(w^\top x_i + \gamma) < 0$$

Based on the above considerations, the constraints of the modified formular result as follows:

$$y_i(w^\top x_i + \gamma) \geq 1 - s_i, s_i \geq 0 \text{ for } i = 1, \ldots, m$$

---

[1]Given $m$ pairs $(x_i, y_i) \in \mathbb{R}^n \times \{+1, -1\}$ and the hyperplane defined by $(w, \gamma) \in \mathbb{R}^{n+1}, i = 1, \ldots, m$

Thus the total deviation error is given by $\sum_{i=1}^{m} s_i$. Since the observations are not classified correctly with $s_i > 1$, so $\sum_{i=1}^{m} s_i$ is an upper bound for the number of misclassified observations. A hyperparameter $\nu > 0$ is introduced to regulate the influence of the slacks. This hyperparameter determines the tradeoff between two opposing goals, namely minimizing the number of missclassified points and maximizing the width of the margin. In other words the aim of $\nu$ is to achieve the best ratio between the widest possible margin and the lowest possible number of misclassifications. Accordingly, the optimization problem

$$
\begin{aligned}
\min_{(w,\gamma)\in\mathbb{R}^{n+1}} \quad & \tfrac{1}{2} w^\top w \\
\text{s. to} \quad & y_i(w^\top x_i + \gamma) \geq 1 && i = 1, \ldots, m
\end{aligned}
$$

is adjusted as follows:

$$
\begin{aligned}
\min_{(w,\gamma,s)\in\mathbb{R}^{n+1+m}} \quad & \frac{1}{2} w^\top w + \nu \sum_{i=1}^{m} s_i \\
\text{s. to} \quad & y_i \left( w^\top x_i + \gamma \right) + s_i \geq 1 \quad i = 1, \ldots, m && i = 1, \ldots, m \qquad (**) \\
& s_i \geq 0
\end{aligned}
$$

Now, all that remains is to talk about the duality for SVM.


## 2.2 Dual problem formulation

Similar to the procedure in the last section of our report, the SVM optimization problem is solved in the form of the dual problem instead of the primary problem (**).
It can be thus shown that, the dual problem can be formulated as follows:

$$
\begin{aligned}
\max \quad & \sum_{i=1}^{m} \lambda_i - \tfrac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\
\text{s. to} \quad & \sum_{i=1}^{m} \lambda_i y_i = 0 \\
& 0 \leq \lambda_i \leq \nu \quad i = 1, \ldots, m
\end{aligned}
$$

The dual problem is a quadratic convex optimization problem. The slacks and the Lagrange multipliers $\mu_i$ are not present in this dual problem.

Based on KKT conditions the point $w^*$ can be determined equivalent to the linearly separable case and is given by

$$
w^* = \sum_{i=1}^{m} \lambda_i^* y_i x_i \tag{1}
$$

For the determination of the parameter $\gamma^*$ we consider any support vector $x_s$ (with the class label $y_s$) with $0 < \lambda_s < \nu$. Thus

$$y_s\left(\langle \mathbf{w}^*, \mathbf{x}_s\rangle + \gamma^*\right) = 1$$

Then we compute $\gamma^*$ as

$$\gamma^* = \frac{1}{y_s} - \langle \mathbf{w}^*, \mathbf{x}_s\rangle \qquad y_s = \pm 1 \tag{2}$$

Let's move on. In the next section, we will exploit these ideas to our classification problem.

# 3  Generated data

## 3.1  Our generator

Since we had several problems with the given generator, we decide to write one of our own in Python, and a function that, taken the random generated data, is capable to spit out data in a format ready for AMPL. The function that we used can be found in the functions.py python script. The concept of the generator is the same of the given one: It randomly generates points n points $x \in \mathbb{R}^4$: if $\sum_{i=1}^{4} \geq 2$ the point belongs to class $+1$; otherwise, to class $-1$. We added also a probability of $0.05$ to misclassificate, in order to have not-linearly separable data.

## 3.2  The choice of $\nu$

$\nu \geq 0$ is a parameter that we have to fix a priori.
The goal is now to find the best value for $\nu$, so we tried different $\nu$'s and these are the results:

|  | accuracy | sensitivity | specifity |
|---|---|---|---|
| 0.001 | 0.54 | 1.000000 | 0.000000 |
| 0.01 | 0.71 | 0.470588 | 0.959184 |
| 1 | 0.84 | 0.844444 | 0.836364 |
| 10 | 0.91 | 0.907407 | 0.913043 |
| 100 | 0.92 | 0.894737 | 0.935484 |
| 1000 | 0.91 | 0.875000 | 0.933333 |

Larger $\nu$ means we assign more importance to reducing number of misclassified points, on the other hand, smaller $\nu$ means we assign more importance to having a large margin.
We discovered that with small $\nu$ results are really bad in terms of accuracy, but once we arrive to reasonably large , the accuracy is stable around the 90%, with a little bit of noise due to the randomicity of the training and the testing data. From now on, we will consider as a given $\nu$ of 10.

## 3.3 Results obtained with SVM

In the primal formulation we obtained:

$$w^* = [4.54007, 5.06347, 5.98317, 5.45954]$$

$$\gamma^* = -10.2623$$

## 3.4 Primal Dual Relationship

We want to show, that for the dual, we have the same hyperplane as the primal. From (1) and (2) let's compute $w^*$ and $\gamma^*$:

For $w^*$, we could do this directly in AMPL with a "fake constraint":

```
subject to get-w{i in 1..m}:
    w[i]=sum{j in 1..n}lambda[j]*y[j]*x[j,i];
```

Instead for $\gamma^*$, we need at first to identify a point displaying $\lambda$ and *slacks* in AMPL.



Figure 1: Slacks          Figure 2: Lambdas

We see that the point 19 is the first one that respects our constraints. We can use this point to retrive $\gamma^*$, with the help of python:

```
gamma_retrieved=1 - np.matmul(W,x_19)
```

that again returns our beloved $\gamma^*$:

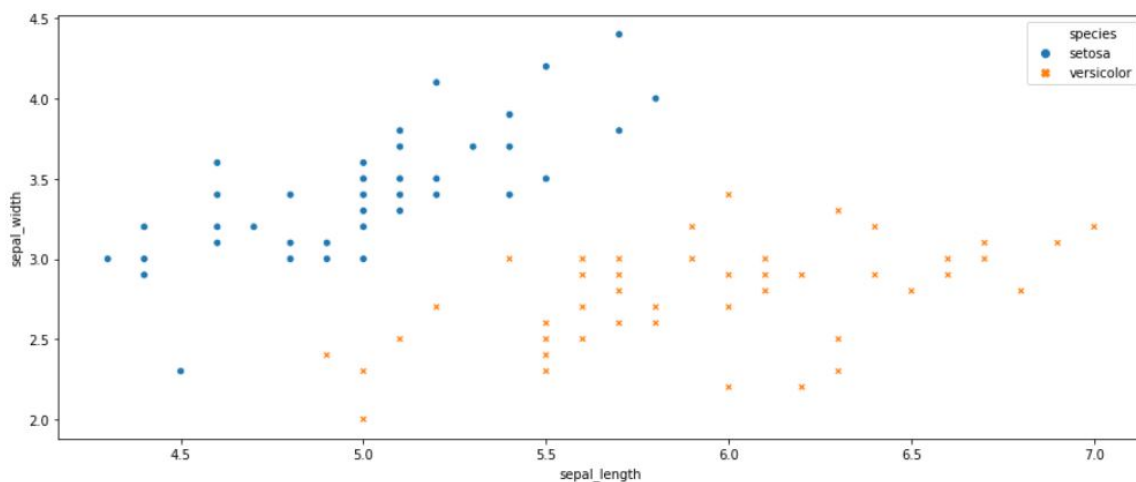$$\gamma^* = -10.2623$$

# 4    Real data



Figure 3: Iris setosa



Figure 4: Iris versicolor

## 4.1    The iris dataset

To test our SVM we decided to use a simple and famous dataset; the Iris dataset. Since we wanted only two classes (to feed the SVM) and two features (to have the possibility to have an easy and immediate graphical representation), we filtered the data set in order to have only two species; setosa and versicolor, and two attributes; the length and the width of the sepal.

Figure 5: Filtred iris dataset



We split our data into two subsets: training data (taking at random 80% of the data) and testing data (the rest).
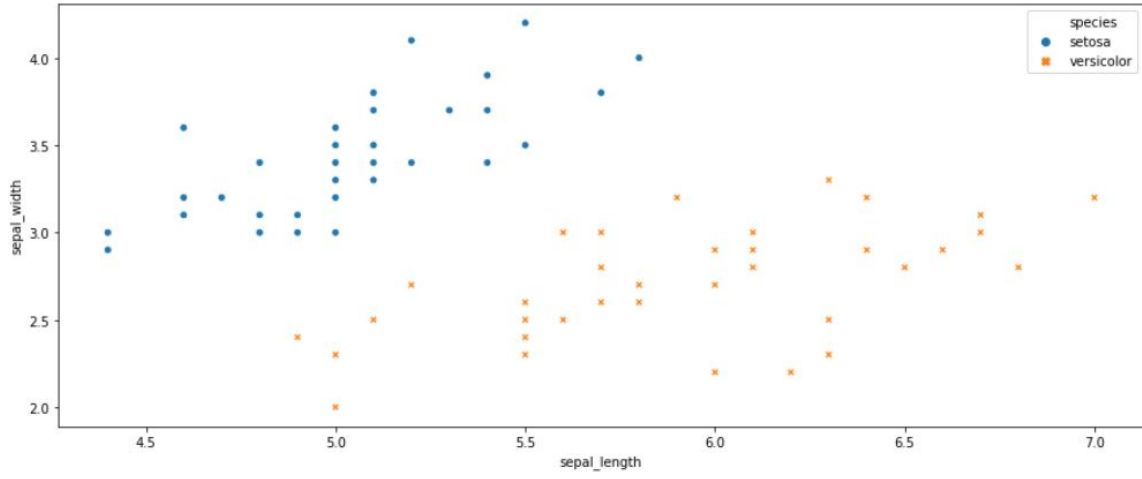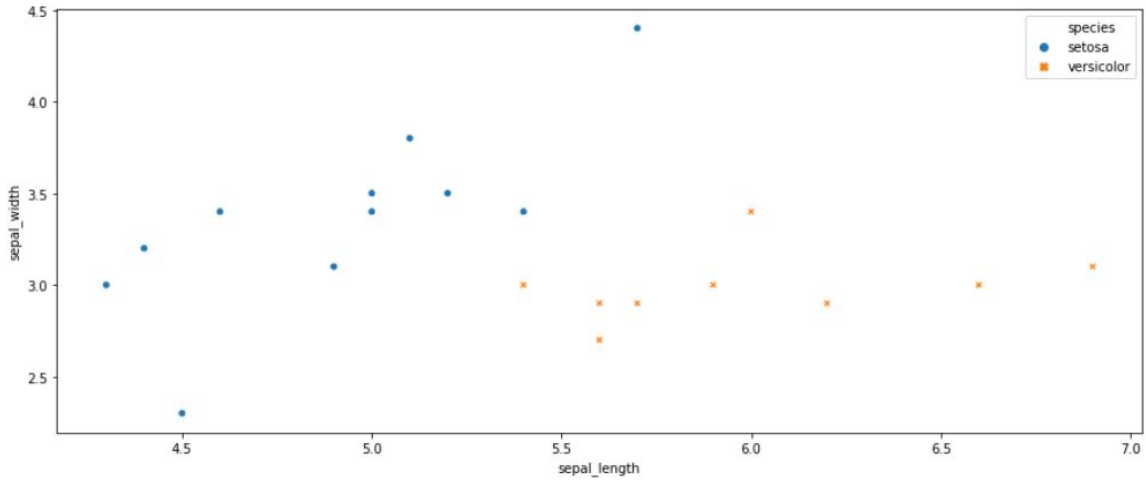
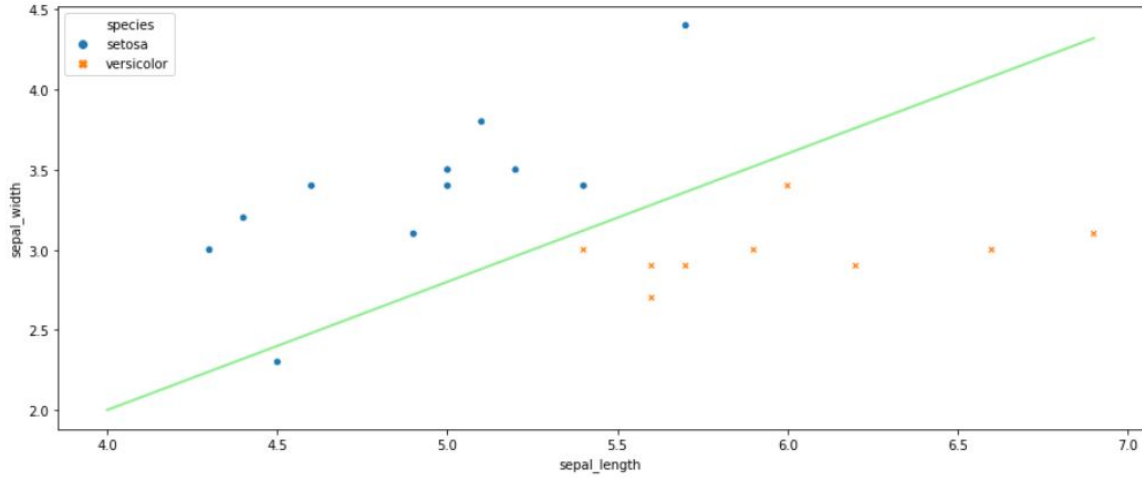Figure 6: Testing data



Figure 7: Training data



## 4.2 Results obtained with SVM

We feed the training data to our AMPL SVM, and we used the parameters $w_1, w_2, \gamma$ to recollect the optimal margin classifier with the equation:

$$y = -(\frac{w_1}{w_2}x + \frac{\gamma}{w_2})$$

applying what we found training our model to the test data, we found these results:

Figure 8: test classification



Our SVM seems to work very well. It only fails on one point, but taking a glance at our test (or full) dataset we see that our single wrong classified point is in some sense an outlier of his class, so its misclassification could be expected.

# 5    Conclusion

In this project, we have seen how to implement a SVM with primal and dual formulations in AMPL that is going to be trained to solve our classification problem. We also have tested the influence of the hyperparameter $\nu$ and how to come up with a way of efficiently finding this hyperparameter that minimize the mistakes and keep the margin as wide as possible. We found that the best accuracy is obtained with $\nu = 10$. Setting this hyperparameter is very important to train our SVM. Since SVM is a convex optimisation problem, we found as expected that the dual problem has the same optimum solution as our original, primal problem.

# 6    Resources

- AMPL model:

  svm_basic.mod $--$ $>$ Primal formulation of SVM

  svm_dual.mod $--$ $>$ Dual formulation of SVM

- Datasets:

  train.dat $--$ $>$ The main training data that we used ($\nu = 10$)

  train_001.dat $--$ $>$ Dataset with $\nu = 0.01$

  train_01.dat  $--$ $>$ Dataset with $\nu = 0.1$

  train_1.dat    $--$ $>$ Dataset with $\nu = 1$

  train_100.dat  $--$ $>$ Dataset with $\nu = 100$

  train_1000.dat $--$ $>$ Dataset with $\nu = 1000$

  iris.csv $--$ $>$ Whole iris dataset

  iris_train_ampl.dat $--$ $>$ iris dataset modified to feed AMPL

  train_iris $--$ $>$ Trainset of iris

  test_iris $--$ $>$ Testset of iris

- Python codes:

  functions.py $--$ $>$ Function that we used in the main notebook stored here for tidiness

  SVM.ipynb $--$ $>$ Main notebook that contains the Python part of the analysis