

16. Support Vector Machine

M.A.Z. Chowdhury and M.A. Oehlschlaeger

Department of Mechanical, Aerospace and Nuclear Engineering
Rensselaer Polytechnic Institute
Troy, New York

chowdm@rpi.edu and oehlsm@rpi.edu

MANE 4962 and 6962

Regular announcement

- 👉 HW 5 is due March 20, 2023
- 👉 Quiz 5 on March 20, 2023
- 👉 Office hour on Friday, March 17, 3-5 PM
- 👉 HW 5 discussion

OR Gate and NOR Gate

input $\underline{x} = (x_1, x_2)$	output y	prediction
(0,0)	0	< 0.5
(0,1)	1	≥ 0.5
(1,0)	1	≥ 0.5
(1,1)	1	≥ 0.5

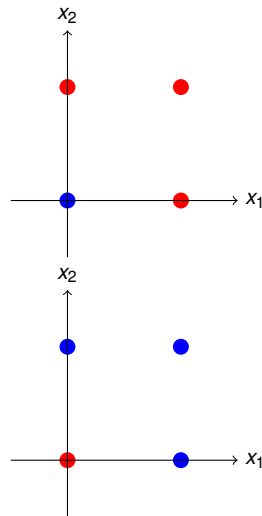
OR Gate: $w_0 = -1$ $w_1 = 2$ $w_2 = 2$

input $\underline{x} = (x_1, x_2)$	output y	prediction
(0,0)	1	≥ 0.5
(0,1)	0	< 0.5
(1,0)	0	< 0.5
(1,1)	0	< 0.5

NOR Gate: $w_0 = 1$ $w_1 = -2$ $w_2 = -2$

Note: w 's are not guaranteed to be unique

red-1, blue-0



AND Gate and NAND Gate

input $\underline{x} = (x_1, x_2)$	output y	prediction
(0,0)	0	< 0.5
(0,1)	0	< 0.5
(1,0)	0	< 0.5
(1,1)	1	≥ 0.5

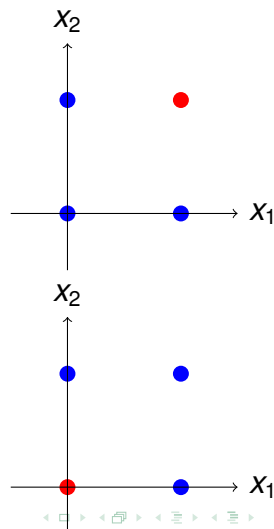
AND Gate: $w_0 = -3$ $w_1 = 2$ $w_2 = 2$

input $\underline{x} = (x_1, x_2)$	output y	prediction
(0,0)	1	≥ 0.5
(0,1)	0	< 0.5
(1,0)	0	< 0.5
(1,1)	0	< 0.5

NAND Gate: $w_0 = 3$ $w_1 = -2$ $w_2 = -2$

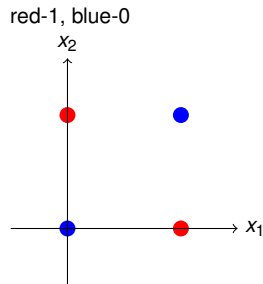
Note: w 's are not guaranteed to be unique

red-1, blue-0



XOR Gate

input $\underline{x} = (x_1, x_2)$	output y	prediction
(0,0)	0	×
(0,1)	1	×
(1,0)	1	×
(1,1)	0	×



XOR Gate: **No solution from logistic regression!**

- Use non-linear feature transformation. (The kernel trick) SVMs.
- Add extra layers. (Deep neural nets)

Support Vector Machines

A powerful and versatile machine learning model first proposed by Vladimir Vapnik. Check the very first notebook where we used SVM blindly to classify iris dataset which achieved very high classification accuracy even without data pre-processing.

- 📖 Linear or nonlinear classification
- 📖 Regression
- 📖 Outlier detection

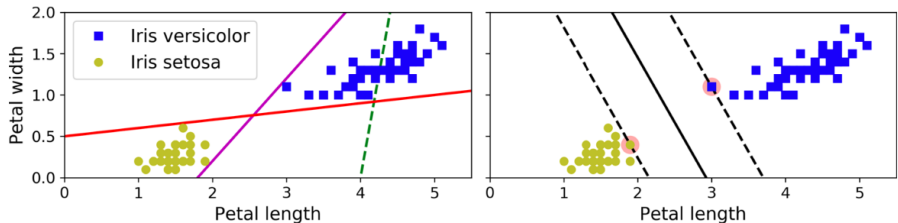
On solving XOR Gate

```

1  import numpy as np
2  from sklearn import svm
3  X = np.array([[0, 0],
4               [0, 1],
5               [1, 0],
6               [1, 1]])
7  y = np.array([0, 1, 1, 0])
8  model = svm.SVC(kernel='rbf', C=1, gamma='auto')
9  model.fit(X, y)
10 predictions = model.predict(X)
11 print("Predictions: ", predictions) # should classify all the data points correctly

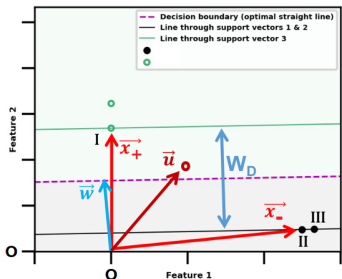
```

Linear Support Vector Machine Classifier



Adding more training data points within the margins of the classes will not change the classifier. The marginal data points are known as support vectors.

Support Vector Machine



Decision rule

$$\vec{w} \cdot \vec{x}_+ + b \geq 1, \quad y_i = +1$$

$$\vec{w} \cdot \vec{x}_- + b \leq -1, \quad y_i = -1$$

Width maximization

$$W_D = (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}$$

$$\max\left(\frac{2}{\|\vec{w}\|}\right) \rightarrow \min(\|\vec{w}\|) \rightarrow \min\left(\frac{1}{2}\|\vec{w}\|^2\right)$$

Constrained optimization problem

$$L = \frac{1}{2}\|\vec{w}\|^2 - \sum_{i=1}^m \alpha_i [y_i(\vec{w} \cdot \vec{x}_i + b) - 1]$$

(Primal form)

$$L = \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \Phi(x_i, x_j)$$

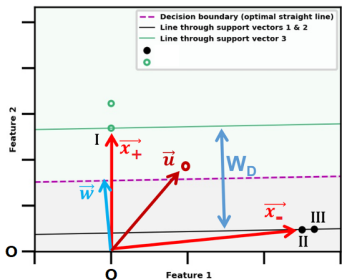
(Dual form)

solved under the constraints via the Lagrange multipliers

$$\alpha_i \geq 0 \quad \text{and} \quad \alpha_i \leq C$$

Larger values of C reduce misclassifications.

Support Vector Machine



Decision rule

$$\vec{w} \cdot \vec{x}_+ + b \geq 1, \quad y_i = +1$$

$$\vec{w} \cdot \vec{x}_- + b \leq -1, \quad y_i = -1$$

Width maximization

$$W_D = (\vec{x}_+ - \vec{x}_-) \cdot \frac{\vec{w}}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|}$$

$$\max\left(\frac{2}{\|\vec{w}\|}\right) \rightarrow \min(\|\vec{w}\|) \rightarrow \min\left(\frac{1}{2}\|\vec{w}\|^2\right)$$

Contd.

Kernel functions $\Phi(x_i, x_j)$

$$\Phi(x_i, x_j) = x_i \cdot x_j \text{ (linear)}$$

$$\Phi(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}} = e^{-\gamma\|x_i - x_j\|^2} \text{ (radial basis function)}$$

$$\Phi(x_i, x_j) = (x_i \cdot x_j + k)^d \text{ (polynomial)}$$

σ relates sensitivity to variance in the feature vectors

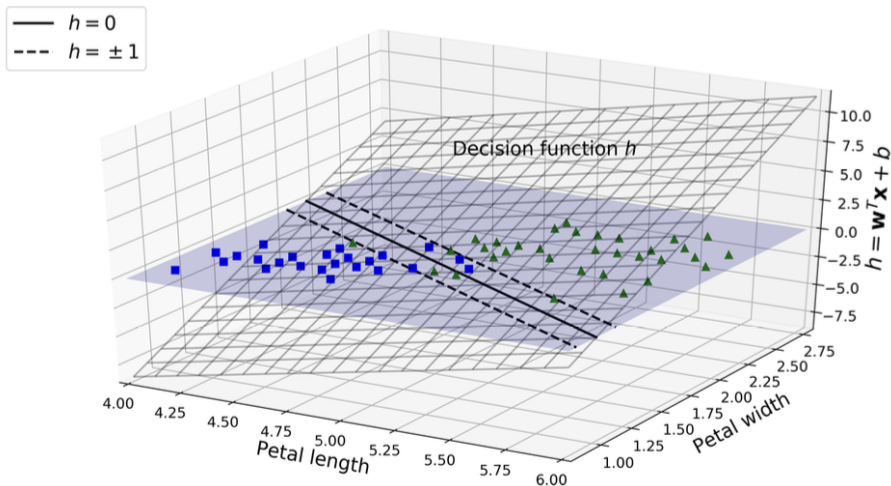
$$\gamma = \frac{1}{2\sigma^2} \geq 0$$

Final solution

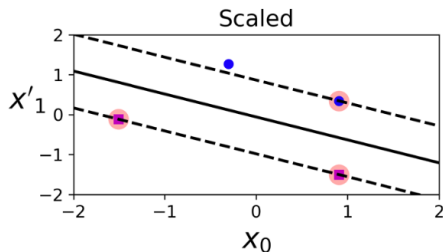
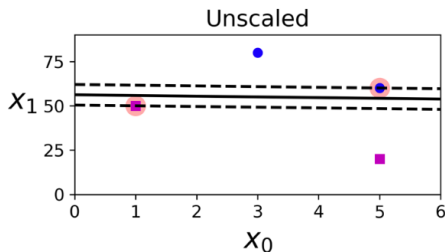
$$\vec{w} = \sum_{i=1}^m \alpha_i y_i \vec{u}$$

$$b = \frac{1}{N_S} \sum_{i \in y_i} y_i - \sum_{j \in y_j} \alpha_j y_j (x_i \cdot x_j)$$

Under the hood

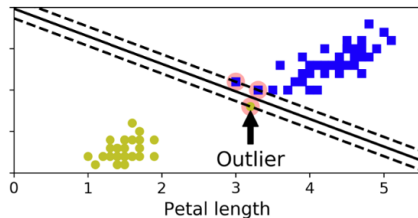
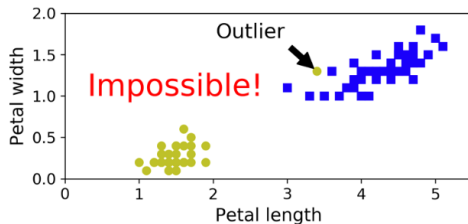


Feature scaling



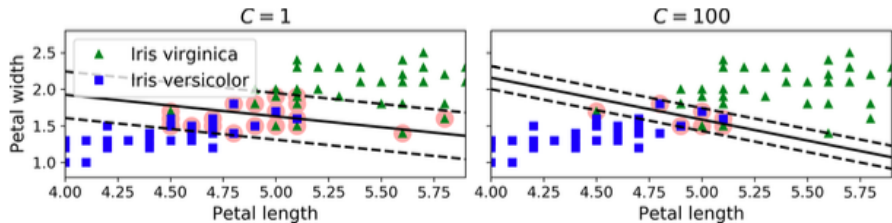
Feature scaling will help widen the gap between the margins.

Hard margin vs Soft margin



Hard margin works only for linearly separable data. It strictly separates the two classes. Soft margin allows misclassification on either side of the margin.

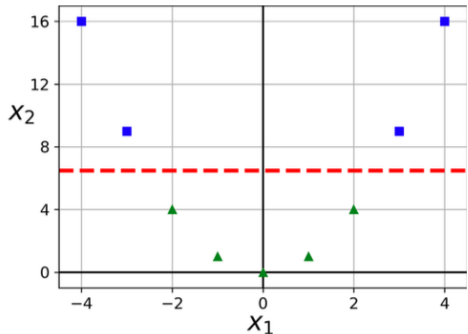
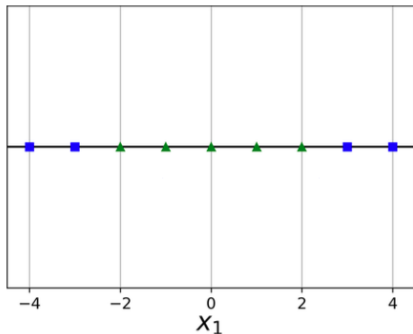
Hard margin vs Soft margin



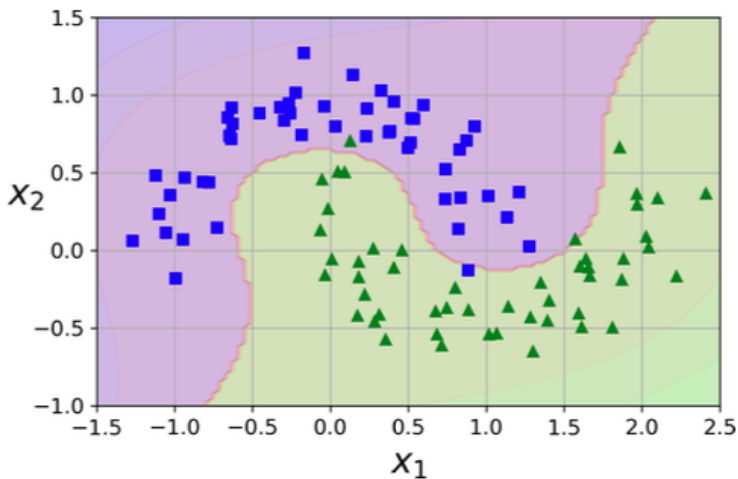
Margin violations are not desirable and thus minimized. Low C models generalize better. C is a regularization hyperparameter, called soft margin constant.

Non-linear SVM:

Adding polynomial features

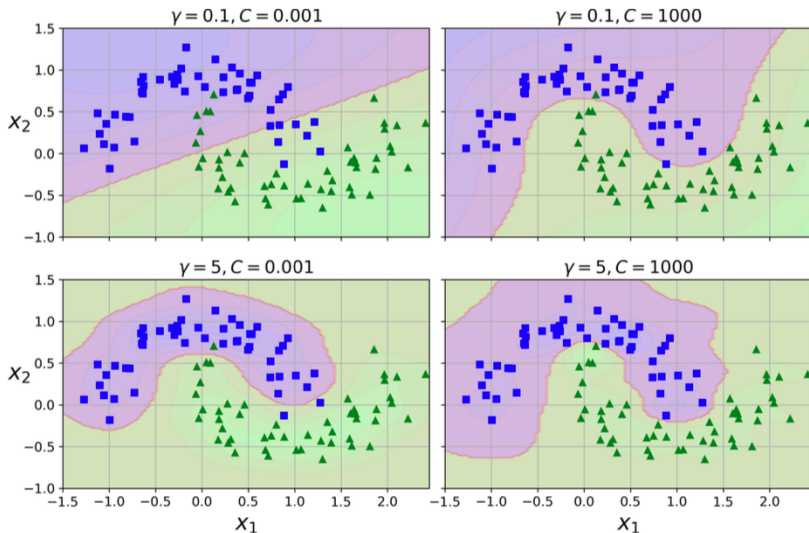


May make a problem linearly separable $x_2 = x_1^2$



Imagine classifying this dataset with just x_1 or x_2 alone.

Radial Basis Function Kernel



Heuristics for choosing the right kernel

Most commonly used kernels are linear, polynomial and radial basis function (Gaussian).

- Always try the linear kernel first, especially if the training set is very large or if it has plenty of features linear kernel will be faster.
- If the training set is not too large, you should also try the Gaussian RBF kernel; it works well in most cases.
- Experiment with other kernels, using cross-validation and grid search if additional computing resources and project time is available.

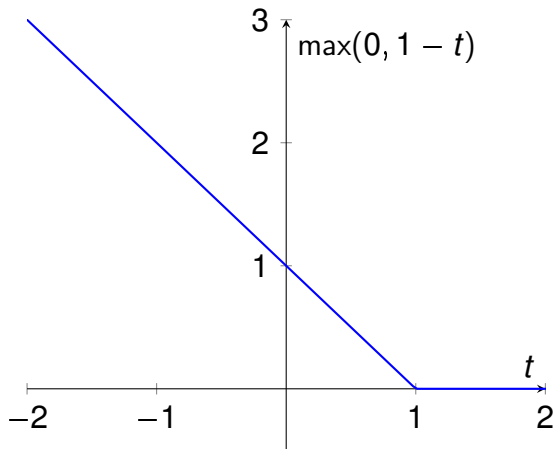
Scikit's implementation of SVM

Make sure to normalize the data. `StandardScaler()` is convenient for normalization. `LinearSVC` is faster but you must specify, the loss function to be hinge

```
1 from sklearn.pipeline import make_pipeline
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.svm import SVC
4 clf = make_pipeline(StandardScaler(), SVC(gamma='auto'))
5 clf.fit(X, y)
```

```
1 import numpy as np
2 from sklearn import datasets
3 from sklearn.pipeline import Pipeline
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.svm import LinearSVC
6 iris = datasets.load_iris()
7
8 X = iris["data"][:, (2, 3)] # petal length, petal width
9 y = (iris["target"] == 2).astype(np.float64) # Iris virginica
10 svm_clf = Pipeline([
11     ("scaler", StandardScaler()),
12     ("linear_svc", LinearSVC(C=1, loss="hinge")),
13 ])
14 svm_clf.fit(X, y)
```

Hinge loss

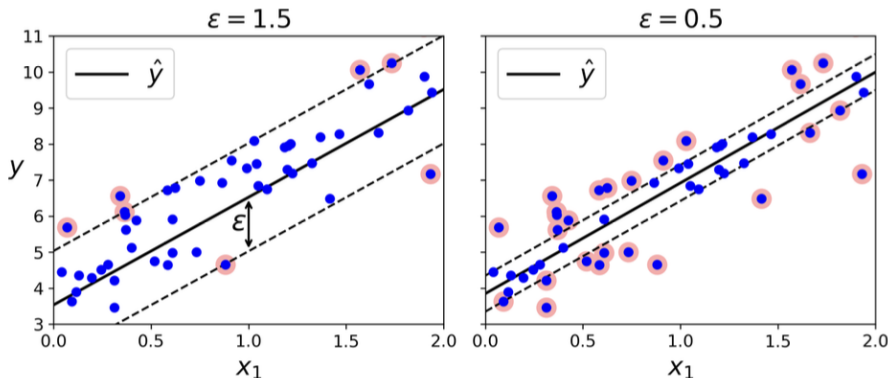


$$L = \max(0, 1 - y\hat{y}) \quad y = \pm 1$$

Computational Time

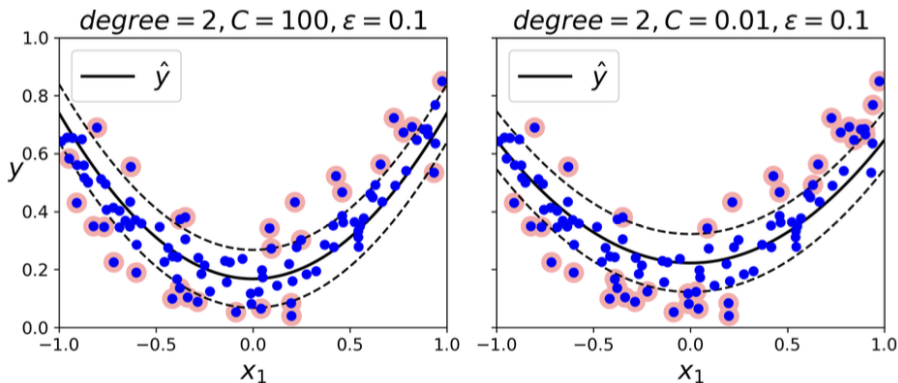
- Scikit's LinearSVC class is based on the LIBLINEAR library, which implements an optimized algorithm for linear SVMs.
- It does not support the kernel trick.
- Computation time scales almost linearly with the number of training data points and the number of features.
- The LIBLINEAR based linear SVM algorithm can be computationally expensive if you need very high precision. This is controlled by the tolerance hyperparameter ϵ (called tol in Scikit-Learn).
- SVC class is based on the LIBSVM library which supports the kernel trick.
- Computation time will scale as a product of either the square or the cubic order of the number of training data points and the number of features. So it will be excruciatingly slow for large datasets.
- Works well with sparse features (meaning most of the features have zero value). This makes the algorithm compute faster.

SVM Regression



SVM Regression attempts to fit as many data points as possible on the street while limiting margin violations. The width of the street is controlled by the epsilon hyperparameter. Results shown for linear kernel.

SVM Regression with polynomial kernel



```
1 from sklearn.svm import LinearSVR
2 svm_reg = LinearSVR(epsilon=1.5)
```

Advantages and disadvantages

Advantages:

- Works well for small to moderate datasets. It will work on larger dataset but computational time may be very high.
- Applicable to classification, regression, outlier detection.
- Needs only the support vectors to make decision.

Disadvantages:

- Can be slow and scaling issues may arise.
- Binary classifier by default, need to extend by training multiple classifiers
- For high-dimensional data (>3), interpretation is difficult.
- Unlike logistic regression, does not output estimated probabilities
- **Optimizing C is crucial for SVM performance.**