

# **Applied Machine Learning**

## **Kernelized Support Vector Machines**

**Kevyn Collins-Thompson**

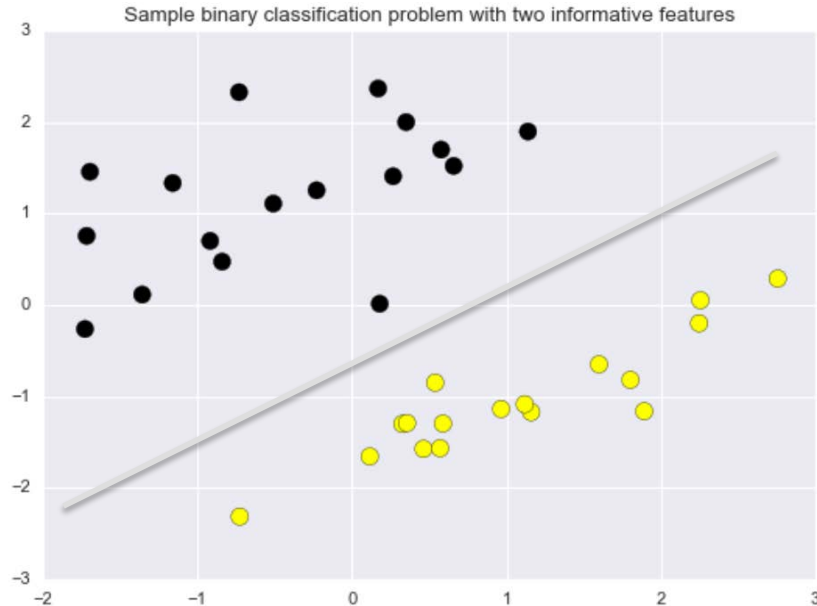
**Associate Professor of Information & Computer Science  
University of Michigan**

**We saw how linear support vector classifiers could effectively find a decision boundary with maximum margin**

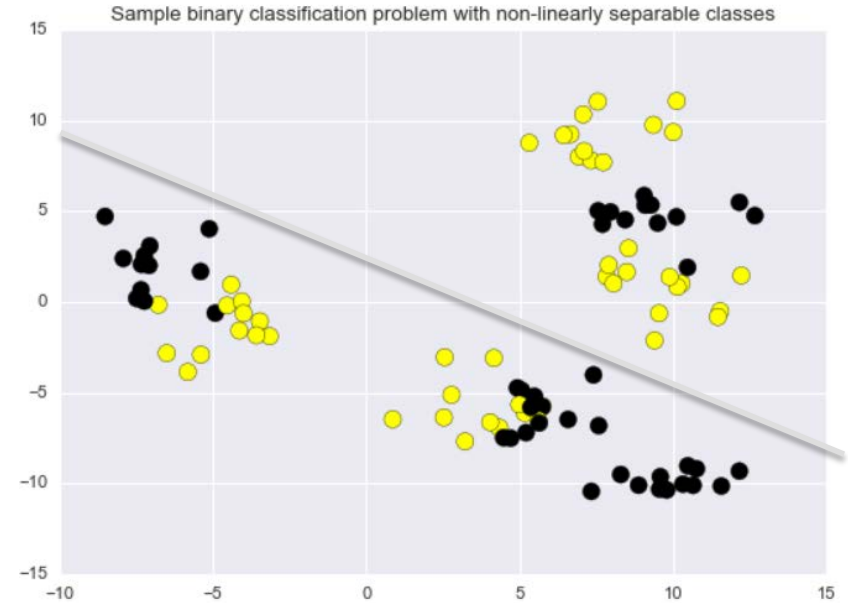


**Easy for a linear classifier**

## But what about more complex binary classification problems?

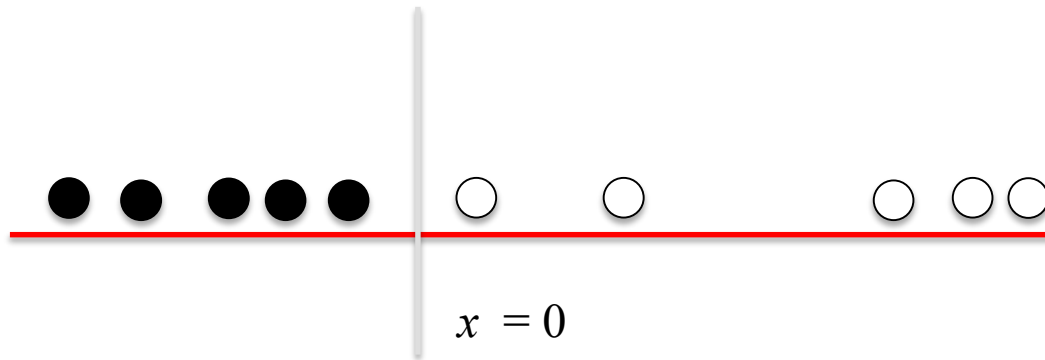


**Easy for a linear classifier**



**Difficult/impossible for a linear classifier**

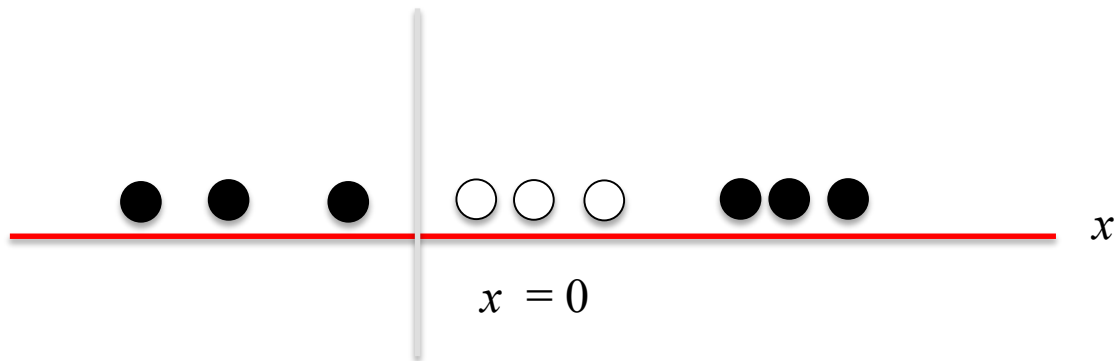
# A simple 1-dimensional classification problem for a linear classifier



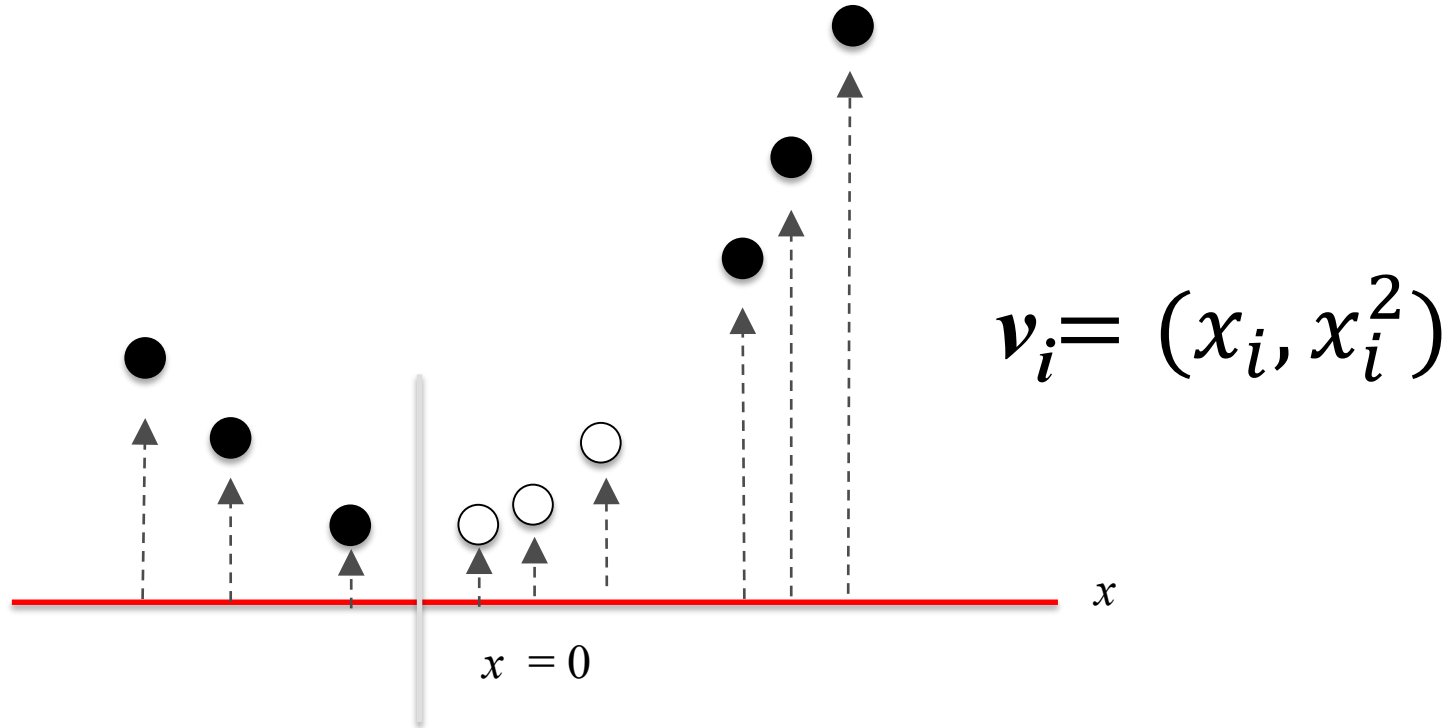
# A more perplexing 1-d classification problem for a linear classifier



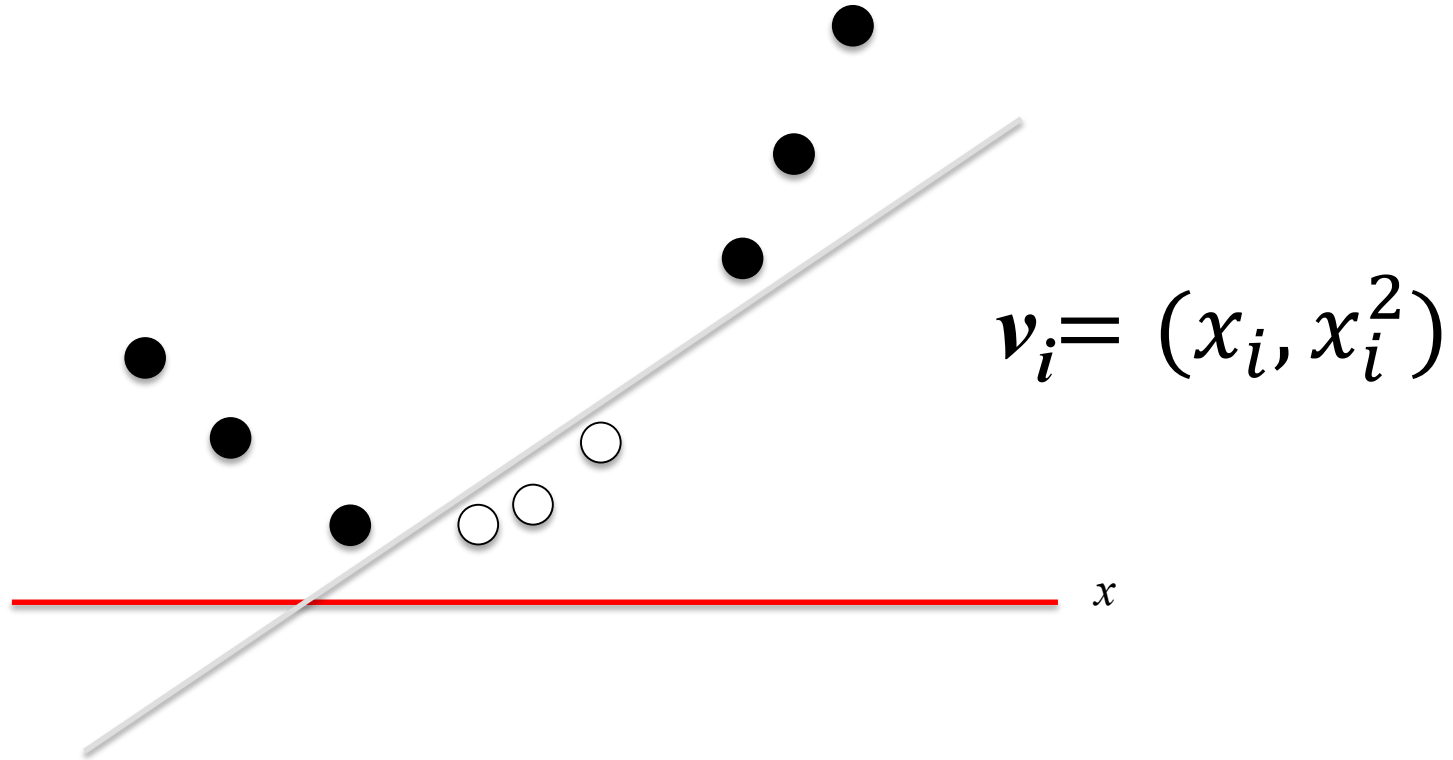
# A more perplexing 1-d classification problem for a linear classifier



Let's transform the data by adding a second dimension/feature  
(set to the squared value of the first feature)



**The data transformation makes it possible to solve this with a linear classifier**

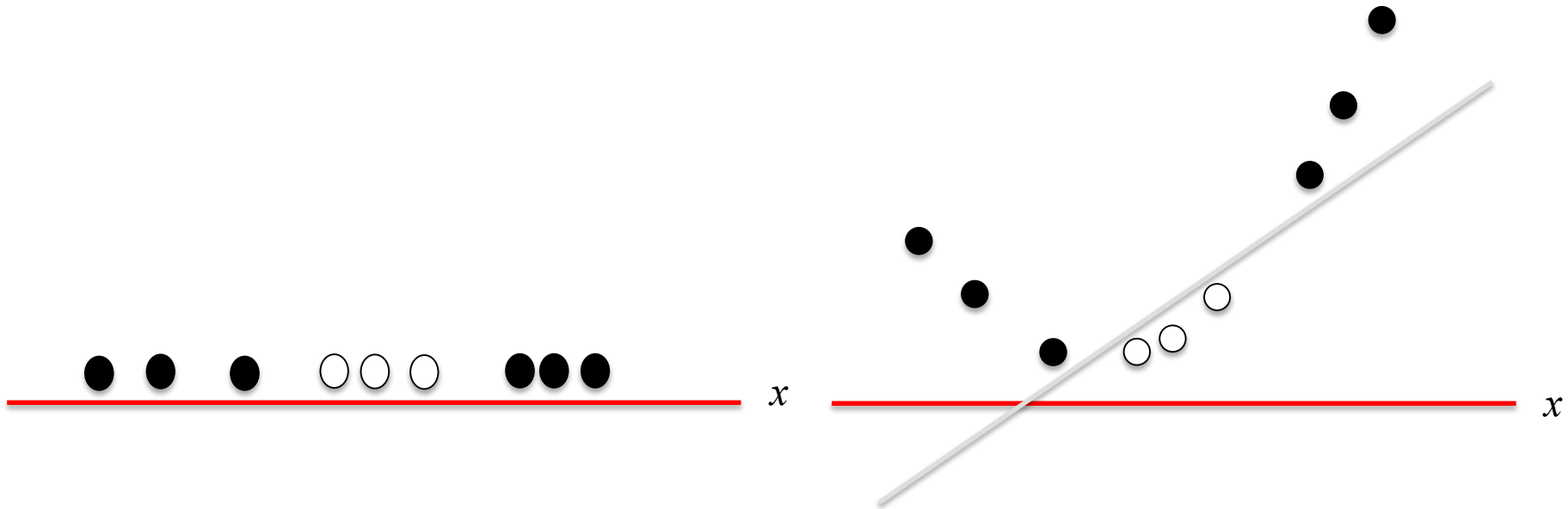




# What does the linear decision boundary in feature space correspond to in the original input space?

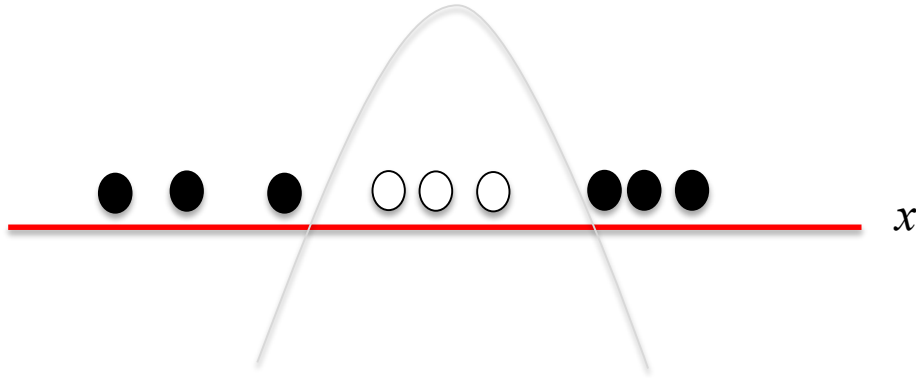
Original input space

Feature space

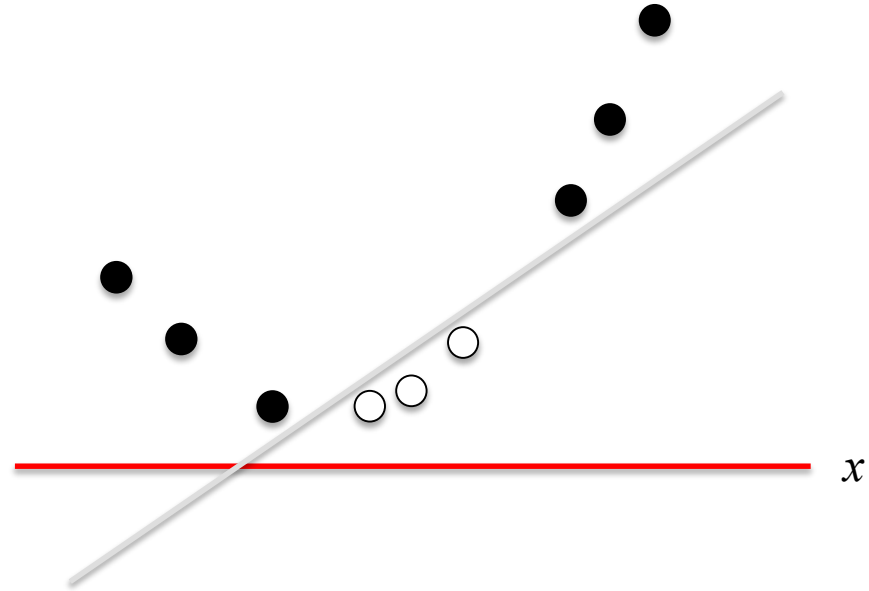


# What does the linear decision boundary correspond to in the original input space?

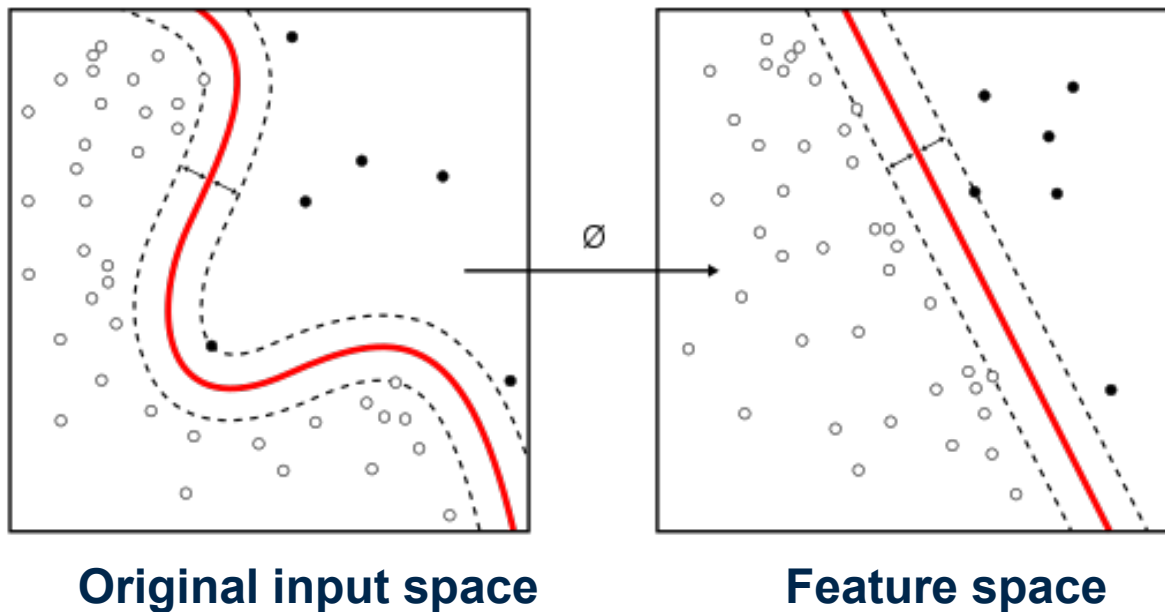
Original input space



Feature space



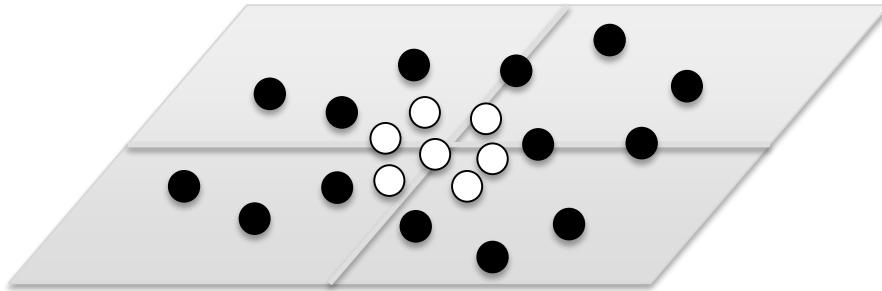
# Transforming the data can make it much easier for a linear classifier.



Source: Wikipedia "Kernel Machine" article.  
<https://commons.wikimedia.org/w/index.php?curid=47868867>

# Example of mapping a 2D classification problem to a 3D feature space to make it linearly separable

Original input space

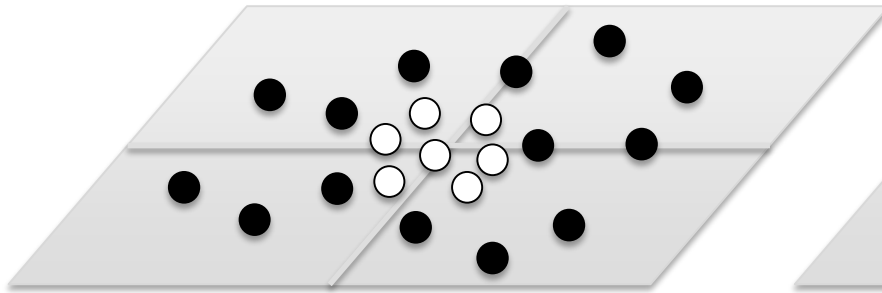


$$x_i = (x_0, x_1)$$

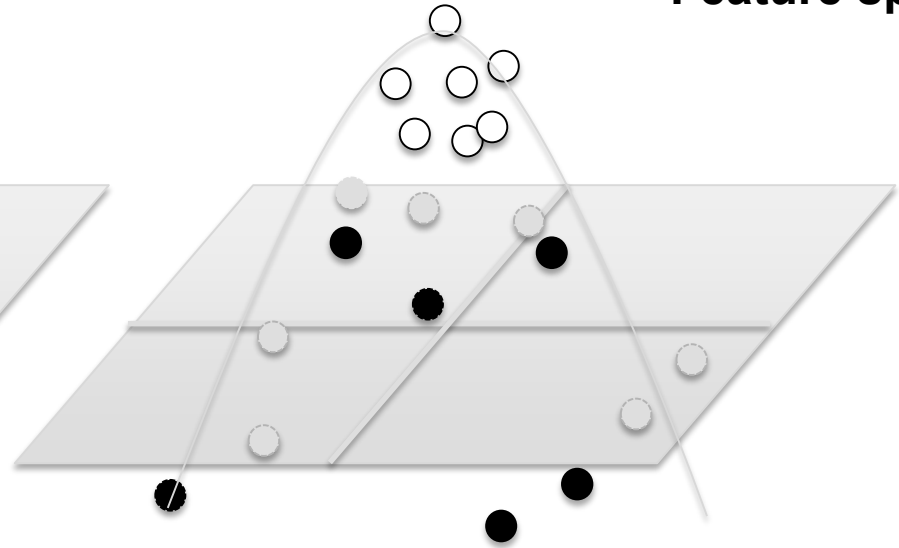
# Example of mapping a 2D classification problem to a 3D feature space to make it linearly separable

Feature space

Original input space

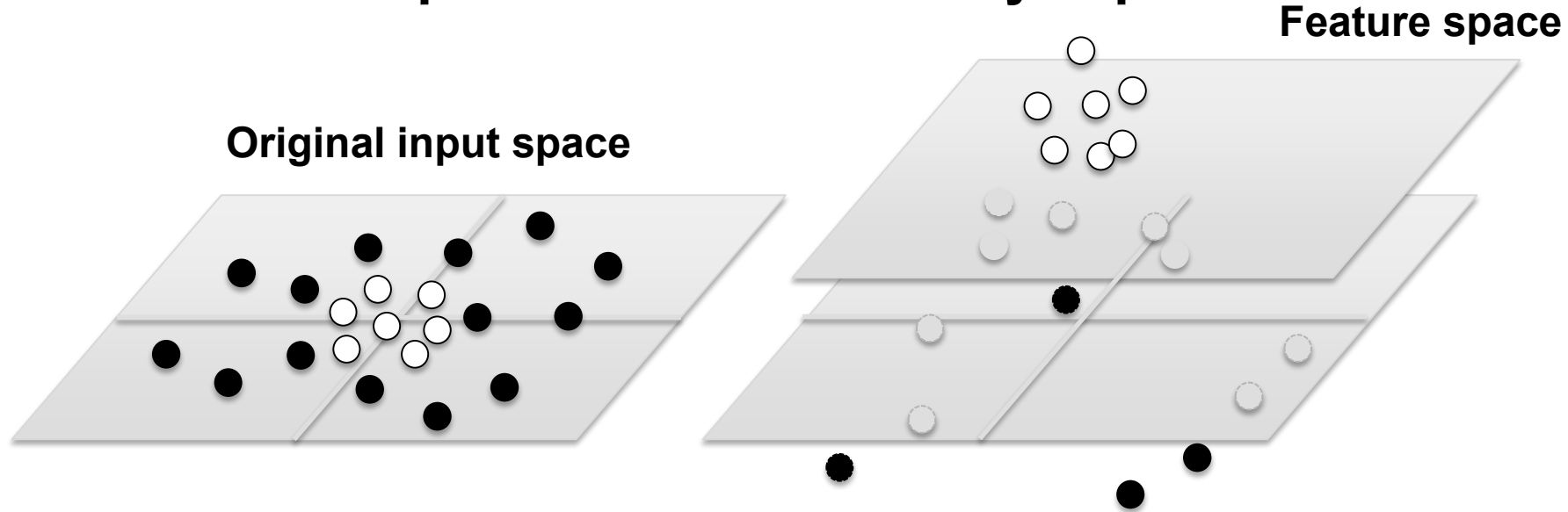


$$\mathbf{x}_i = (x_0, x_1)$$



$$\mathbf{v}_i = (x_0, x_1, 1 - (x_0^2 + x_1^2))$$

# Example of mapping a 2D classification problem to a 3D feature space to make it linearly separable

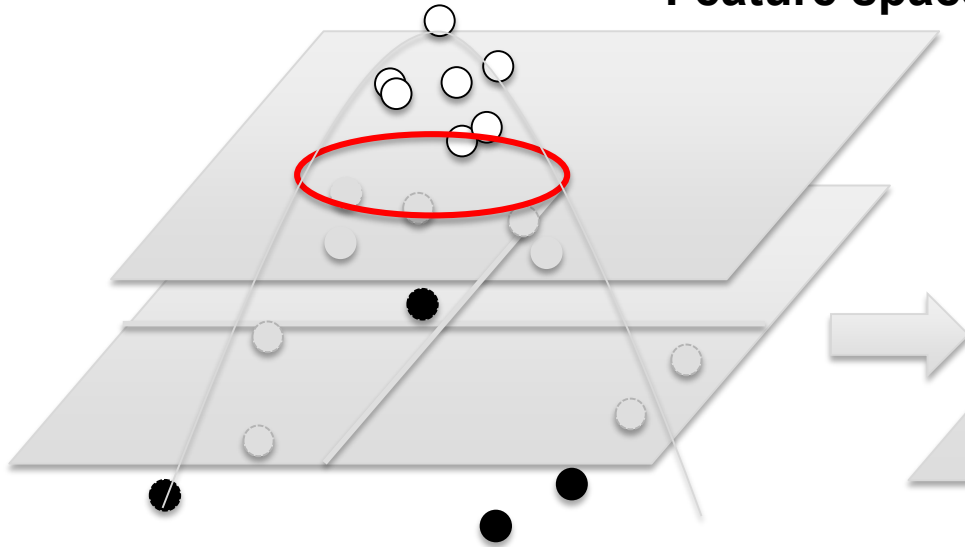


$$\mathbf{x}_i = (x_0, x_1)$$

$$\mathbf{v}_i = (x_0, x_1, 1 - (x_0^2 + x_1^2))$$

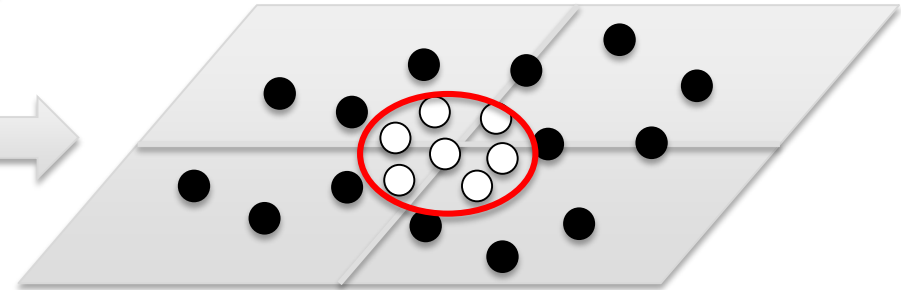
# Example of mapping a 2D classification problem to a 3D feature space to make it linearly separable

Feature space



$$\mathbf{v}_i = (x_0, x_1, 1 - (x_0^2 + x_1^2))$$

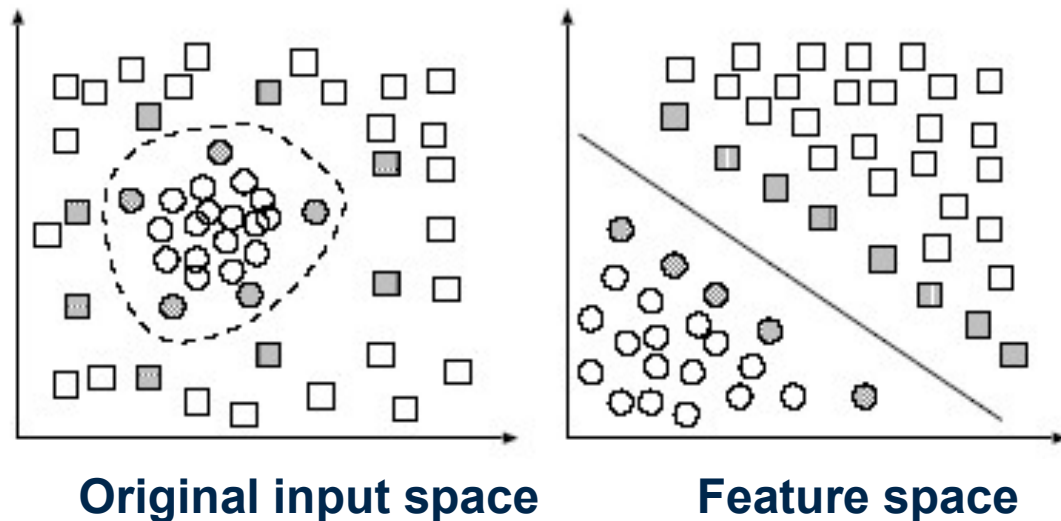
Original input space



$$\mathbf{x}_i = (x_0, x_1)$$

# Radial Basis Function Kernel

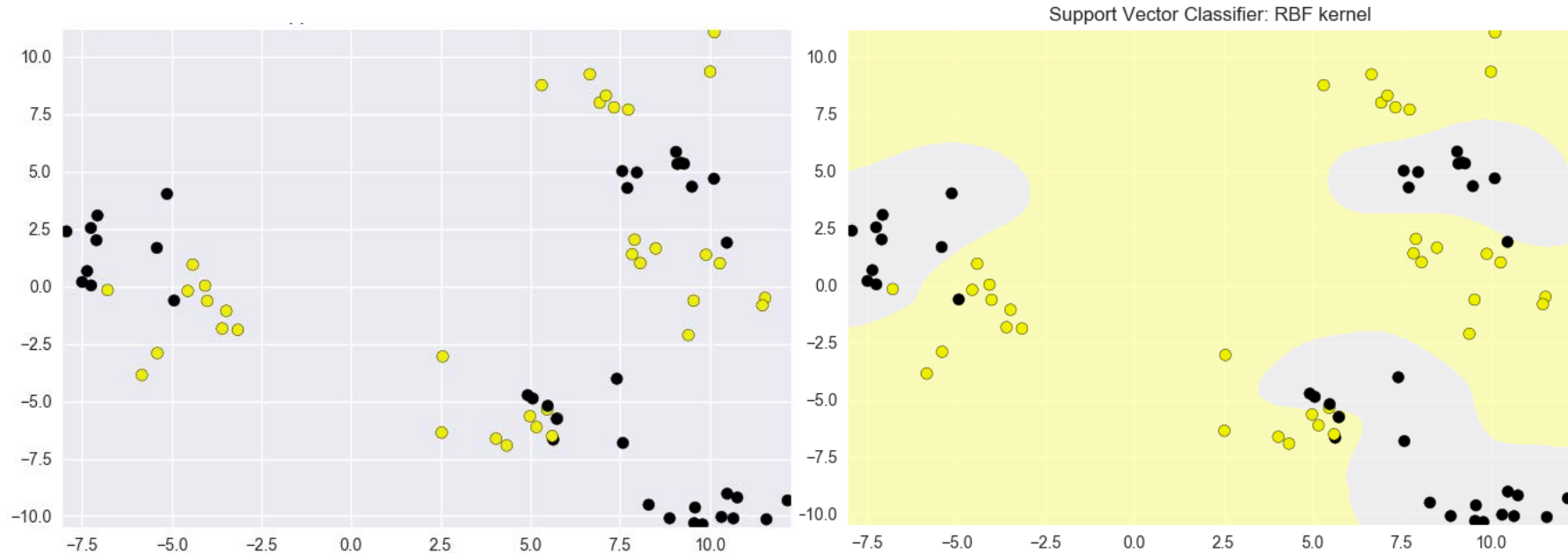
$$K(x, x') = \exp [-\gamma \cdot \|x - x'\|^2]$$



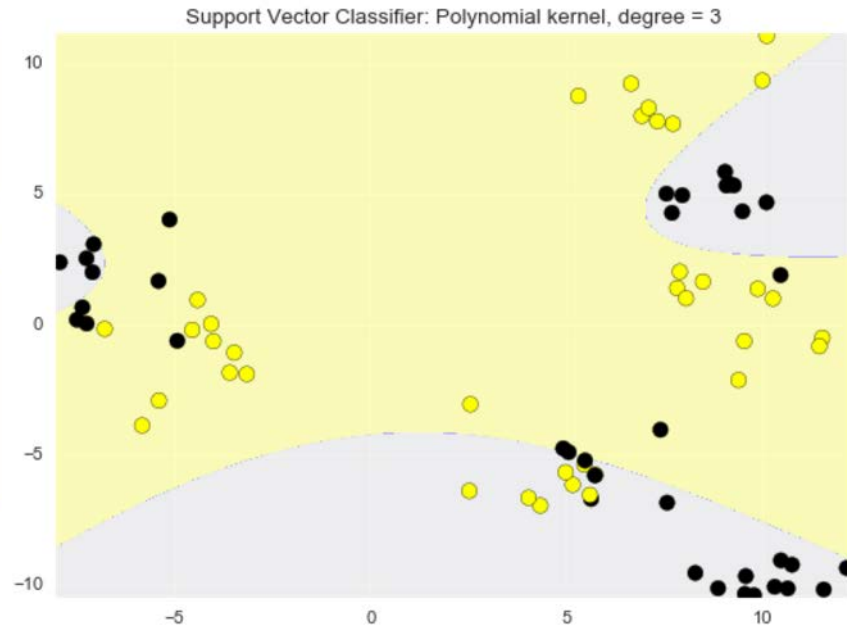
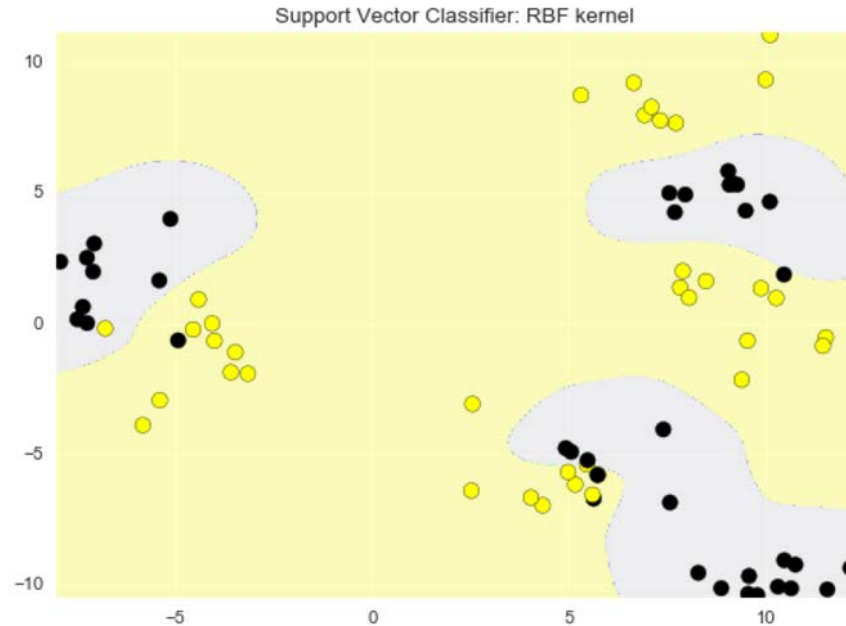
**A kernel is a similarity measure (modified dot product) between data points**



# Applying the SVM with RBF kernel



# Radial Basis Kernel vs Polynomial Kernel



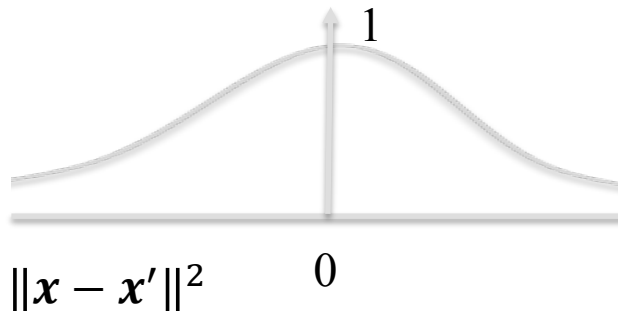
# Radial Basis Function kernel: Gamma Parameter

$$K(\mathbf{x}, \mathbf{x}') = \exp [-\gamma \cdot \|\mathbf{x} - \mathbf{x}'\|^2]$$



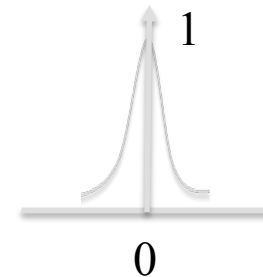
gamma ( $\gamma$ ): kernel width  
parameter

small gamma (0.01)

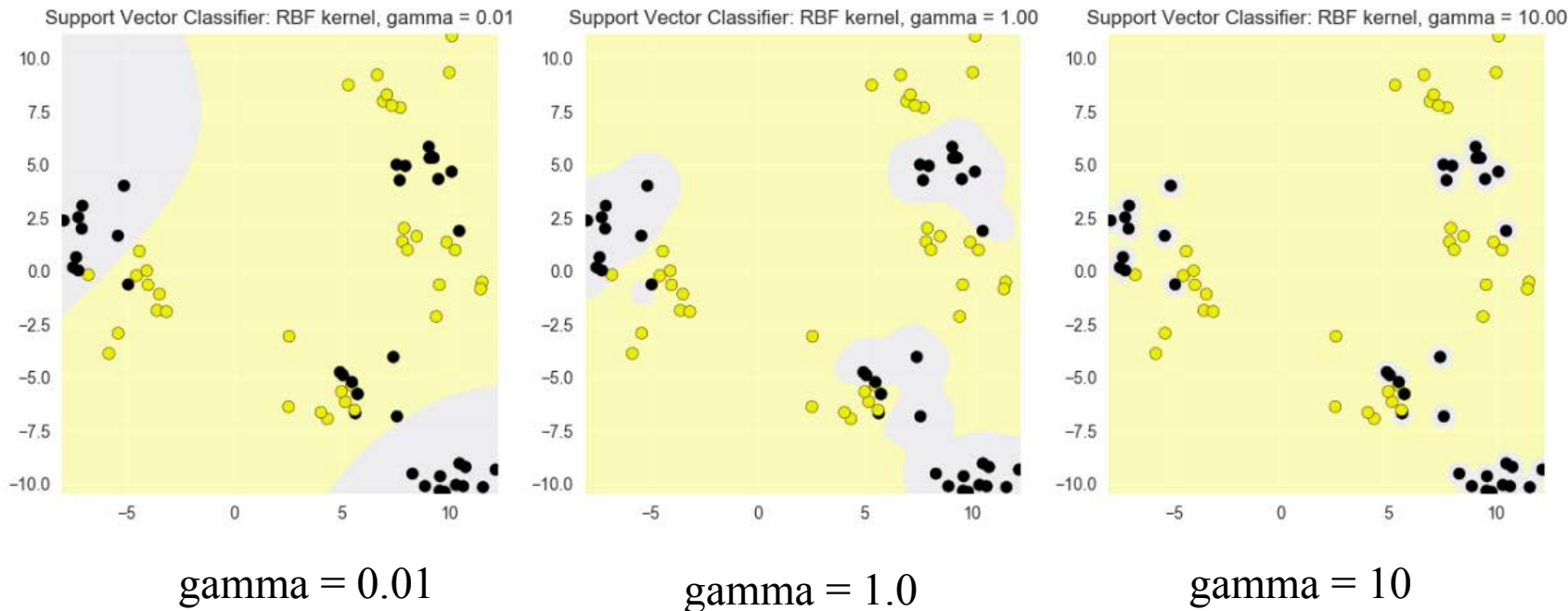


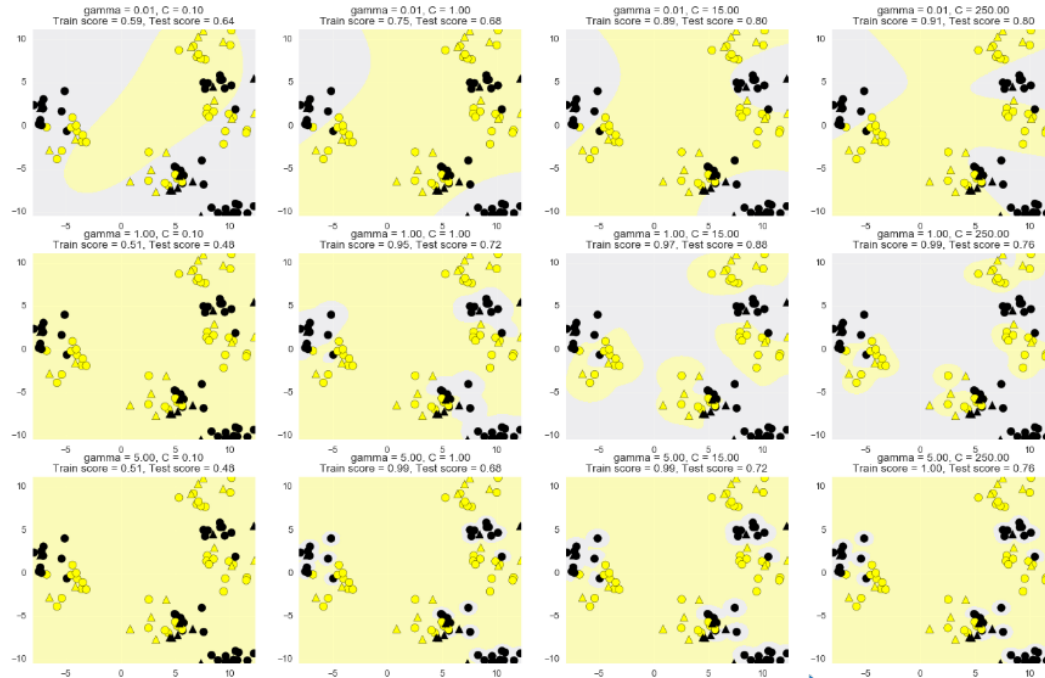
Squared distance  
between  $\mathbf{x}$  and  $\mathbf{x}'$

large gamma (10)



# The effect of the RBF gamma parameter on decision boundaries





Increasing gamma

Increasing C

## Reminder: Using a scaler object: fit and transform methods

```
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
scaler.fit(X_train)
X_train_scaled = scaler.transform(X_train)
X_test_scaled = scaler.transform(X_test)
clf = SVC().fit(X_train_scaled, y_train)
accuracy = clf.score(X_test_scaled, y_test)
```

**Tip:** It can be more efficient to do fitting and transforming together on the training set using the `fit_transform` method.

```
scaler = MinMaxScaler()
X_train_scaled = scaler.fit_transform(X_train)
```

# Kernelized Support Vector Machines: pros and cons

## Pros:

- Can perform well on a range of datasets.
- Versatile: different kernel functions can be specified, or custom kernels can be defined for specific data types.
- Works well for both low- and high-dimensional data.

## Cons:

- Efficiency (runtime speed and memory usage) decreases as training set size increases (e.g. over 50000 samples).
- Needs careful normalization of input data and parameter tuning.
- Does not provide direct probability estimates (but can be estimated using e.g. Platt scaling).
- Difficult to interpret why a prediction was made.

# Kernelized Support Vector Machines (SVC): Important parameters

## Model complexity

- **kernel:** Type of kernel function to be used
  - Default = 'rbf' for radial basis function
  - Other types include 'polynomial'
- **kernel parameters**
  - `gamma` ( $\gamma$ ): RBF kernel width
- **C: regularization parameter**
- Typically C and `gamma` are tuned at the same time.