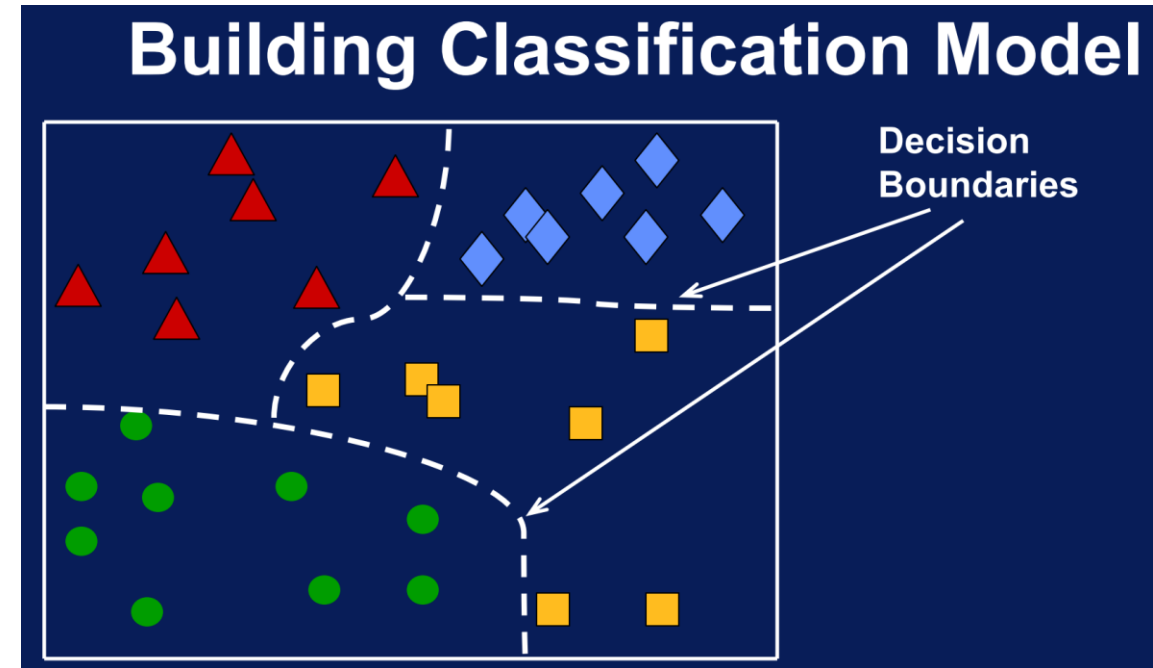
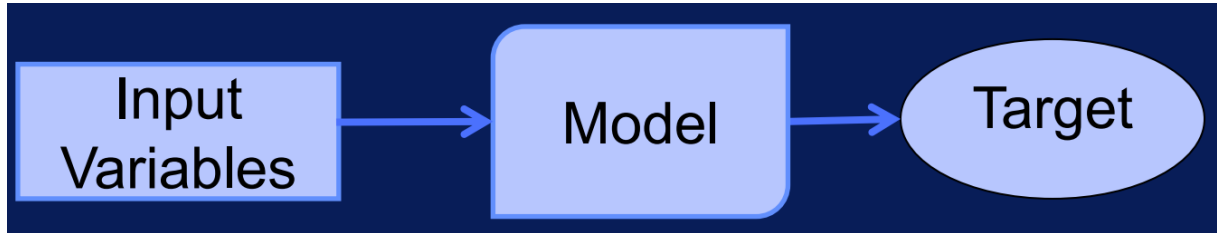


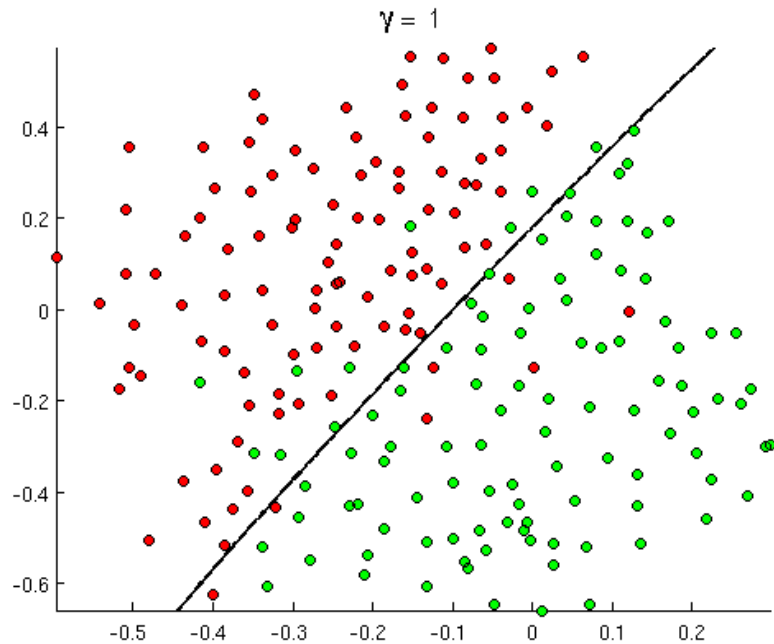
# Classification Model

# classification

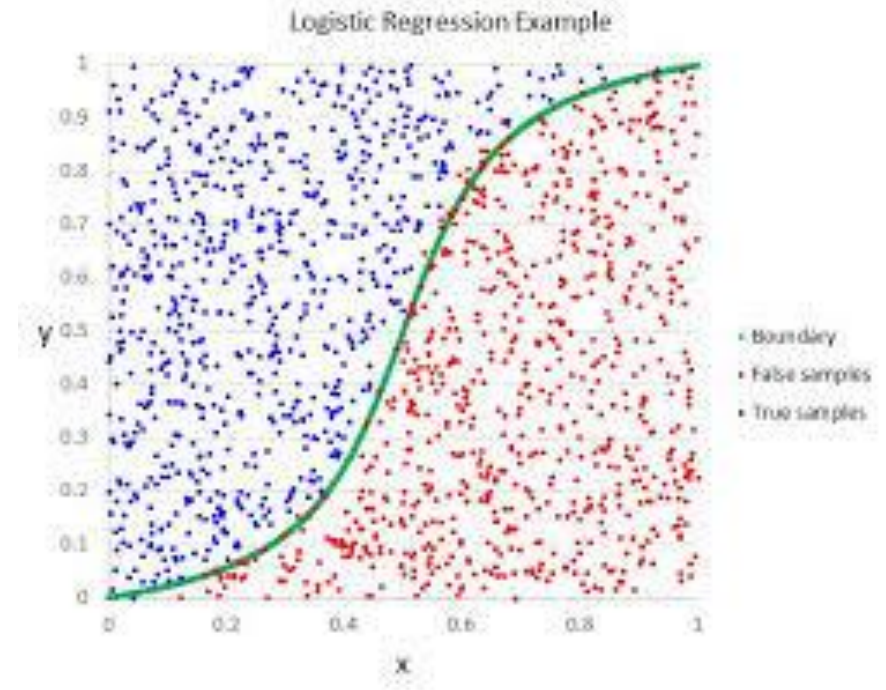
- Predict: Category from input variables
- Goal: Match model outputs to targets (desired outputs)



# Classification (single dimension Input)

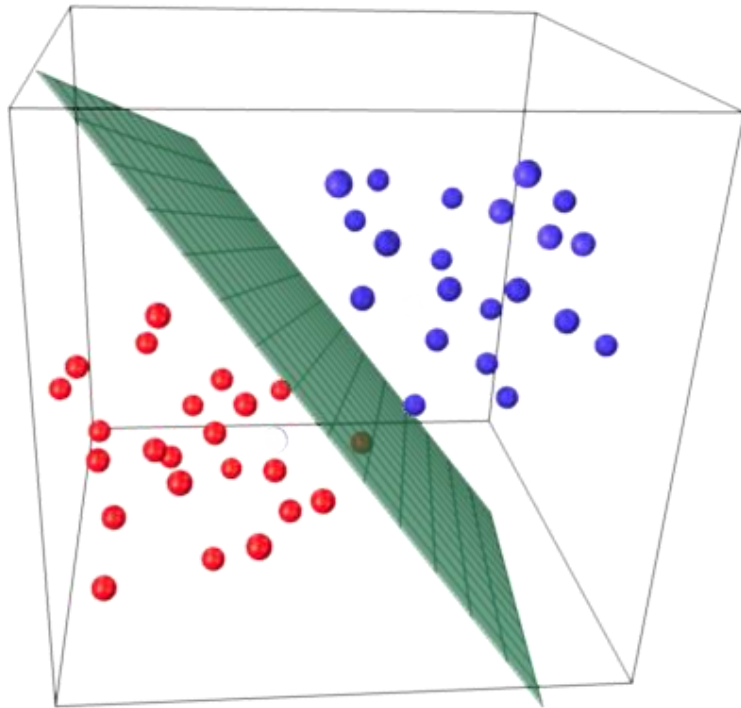


**Linear Classifier**

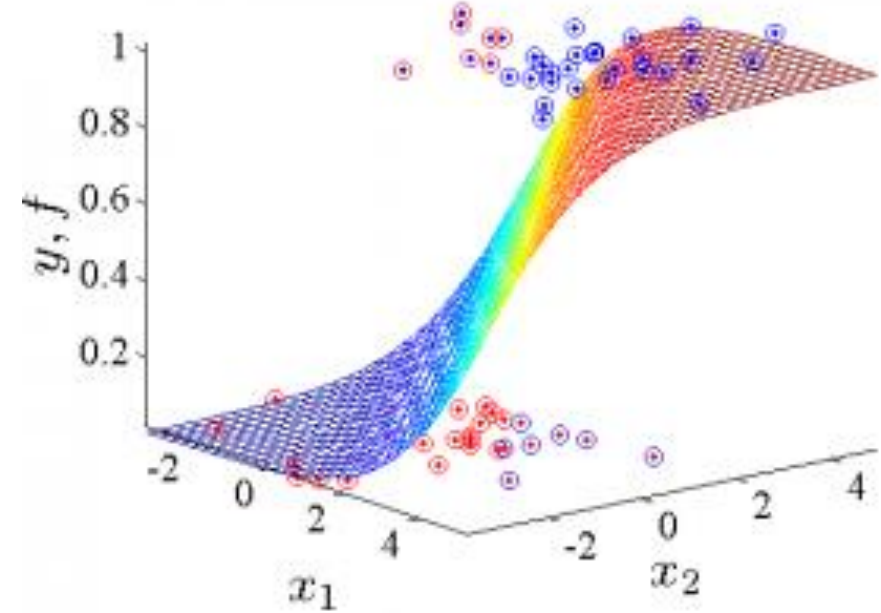


**Nonlinear Classifier  
(Logistic / Sigmoid)**

# Classification (Multi-dimension Inputs)



**Linear Classifier**



**Nonlinear Classifier  
(Logistic / Sigmoid)**

# Supervised vs Unsupervised Model

Supervised Classification	Unsupervised Clustering
<ul style="list-style-type: none"><li>• known number of classes</li><li>• based on a training set</li><li>• used to classify future observations</li></ul>	<ul style="list-style-type: none"><li>• unknown number of classes</li><li>• no prior knowledge</li><li>• used to understand (explore) data</li></ul>

# Supervised Classification

## Single Node MODEL

- **kNN (k-Nearest Neighbor)**
- **Logistic Regression**
- **Support Vector Machine**

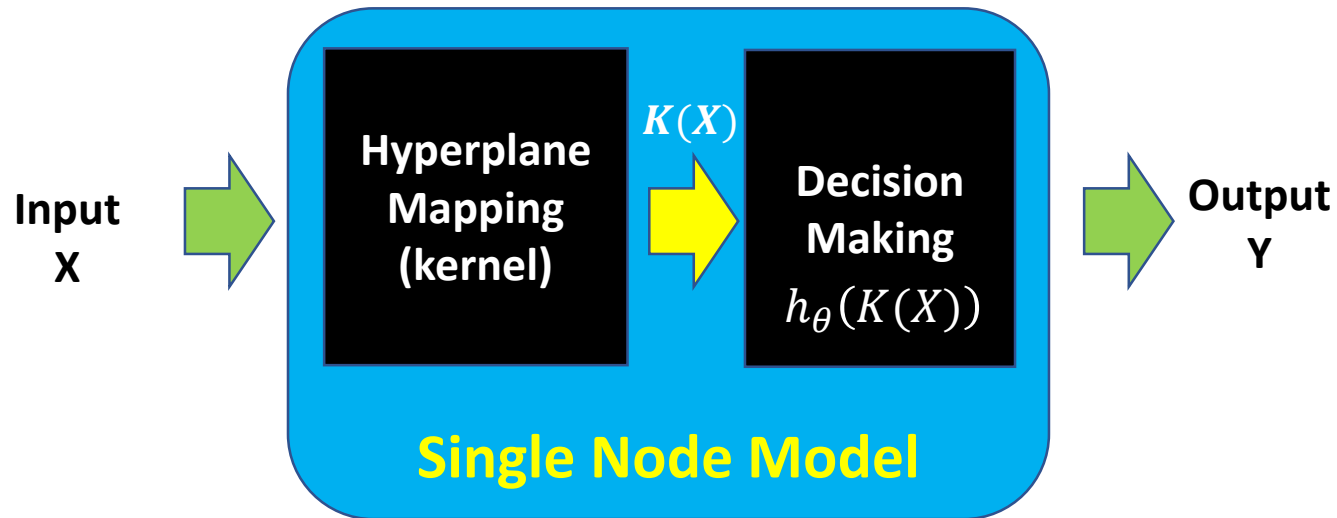
## Multi-Node Model

- **Neural Networks (Deep Networks)**
- **Convolutional Neural Network**
- **Recurrent Neural Network**
  - **Long Short Term Memory ( LSTM )**

# SINGLE NODE MODEL

**Logistic regression      vs      SVM**

# Single node Model Architecture

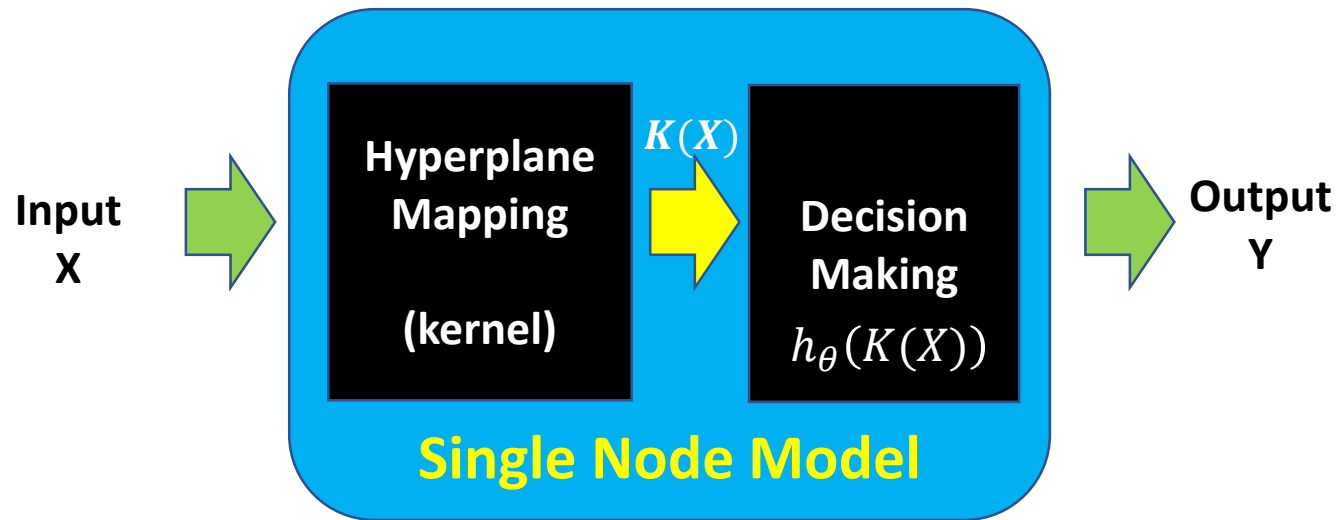


- **Model Architecture**

- Hyperplane Mapping through
  - Model Kernel
- Decision Making through
  - Hypothesis Function
  - Activation Function

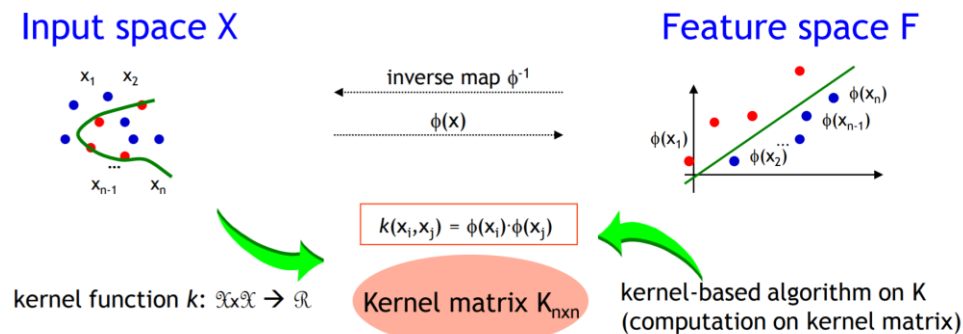


# Single node Model Architecture

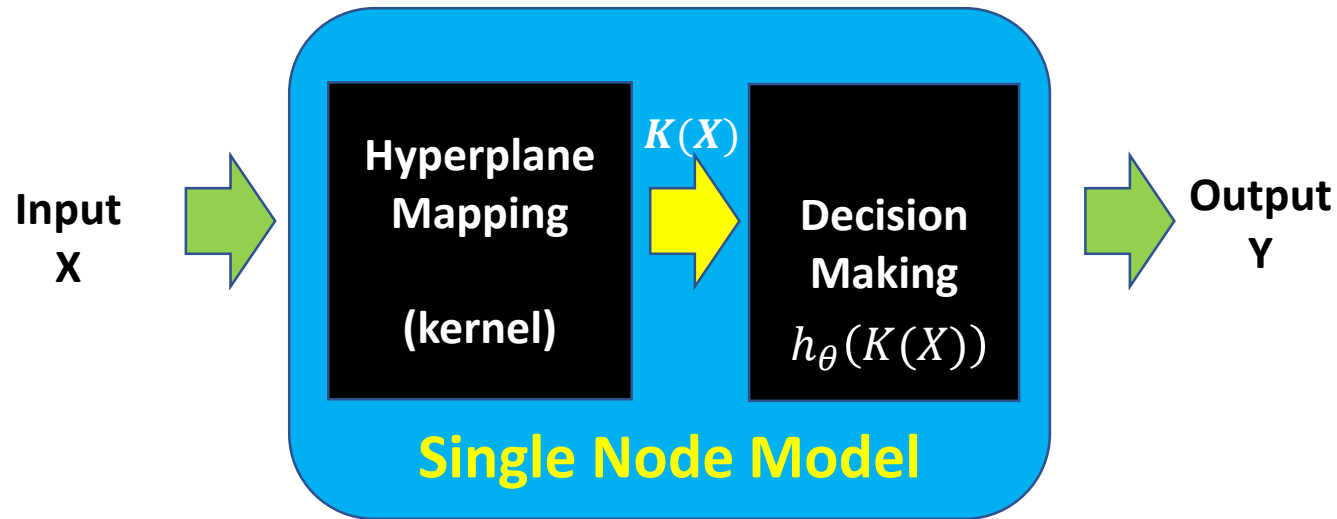


- **Model Architecture**

- Hyperplane Mapping through
  - Kernel (K)
    - Linear kernel
    - Nonlinear kernel
      - Gaussian
      - Radial Basis Function (RBF)
      - Polynomial

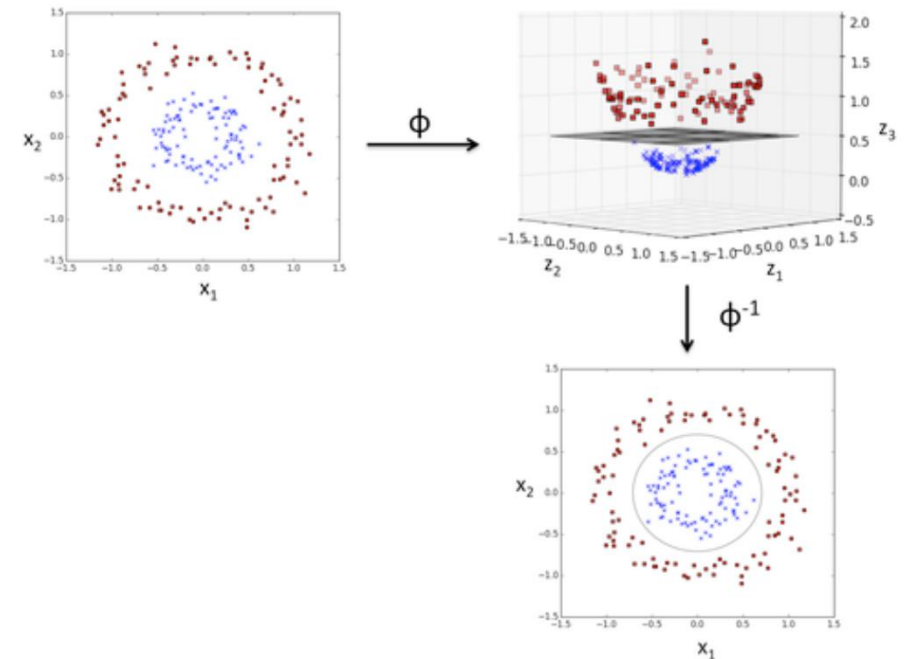


# Single node Model Architecture



- **Model Architecture**

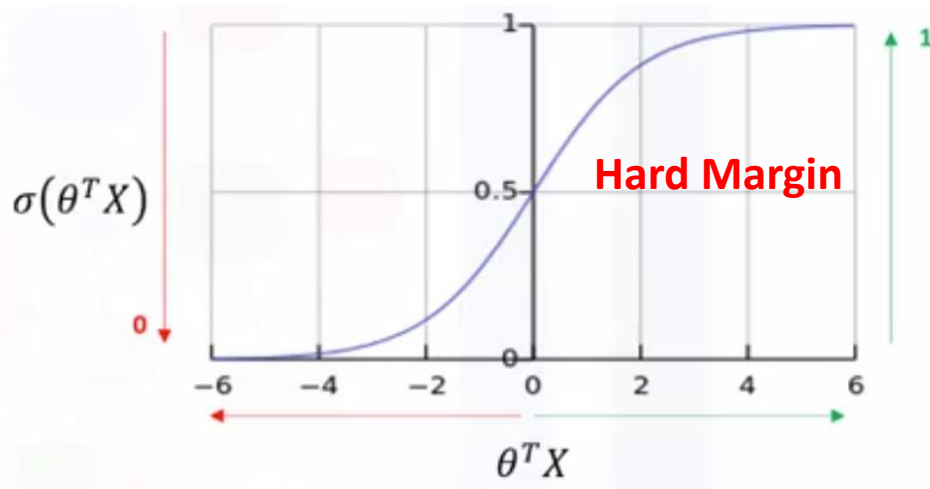
- Hyperplane Mapping through
  - Kernel (K)



# Model Kernel: Hyperplane Mapping

## Logistic Regression

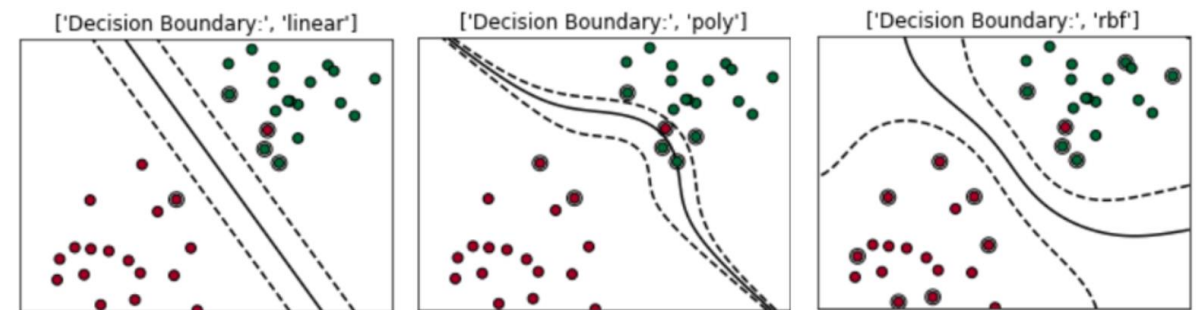
- Sigmoid / Logistic Kernel Function



$$K(X) = \theta^T X = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

## Support Vector Machine (SVM)

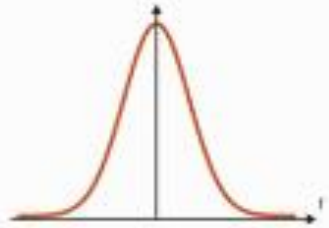
- **Class Partition** Kernels: Linear / Non-Linear



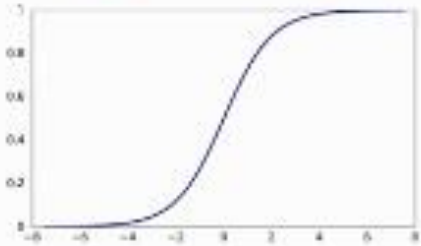
$$K(X) = x_i^T x_j$$

$$K(X) = \exp(-\gamma \|x_i - x_j\|^2)$$

# SVC: kernels



Gaussian RBF Kernel



Sigmoid Kernel



Polynomial Kernel

Gaussian kernel

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

Gaussian radial basis function (RBF)

$$K(x_i, x_j) = \exp\left(-\gamma\|x_i - x_j\|^2\right)$$

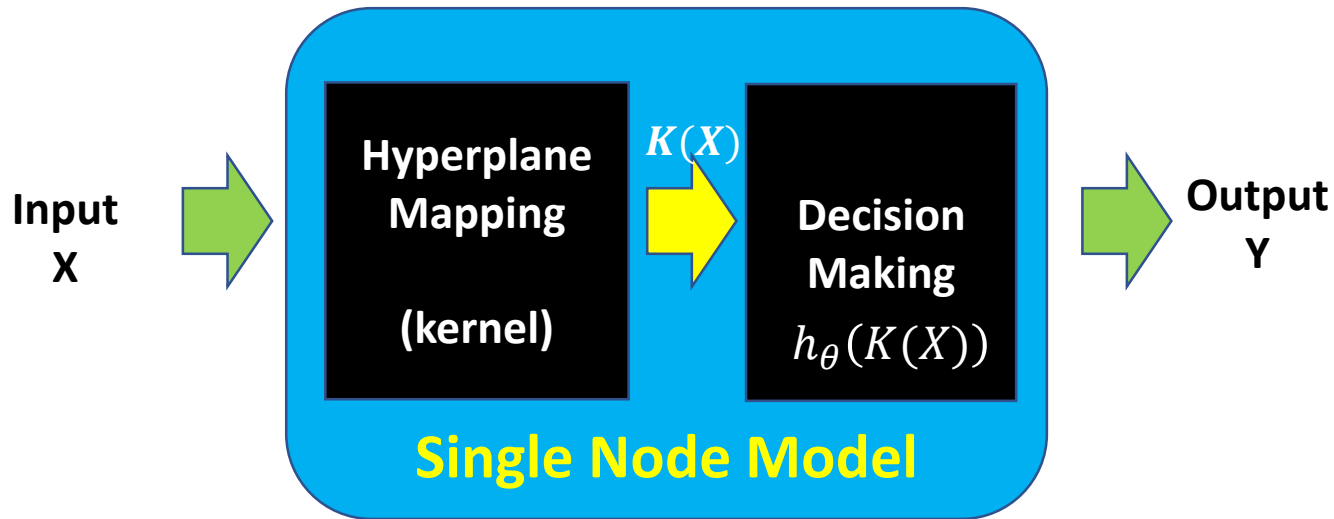
Polynomial kernel

$$K(x_i, x_j) = (x_i x_j)^d$$

Sigmoid kernel

$$K(x_i, x_j) = \tanh(\alpha x_i^T x_j + c)$$

# Single node Model Architecture



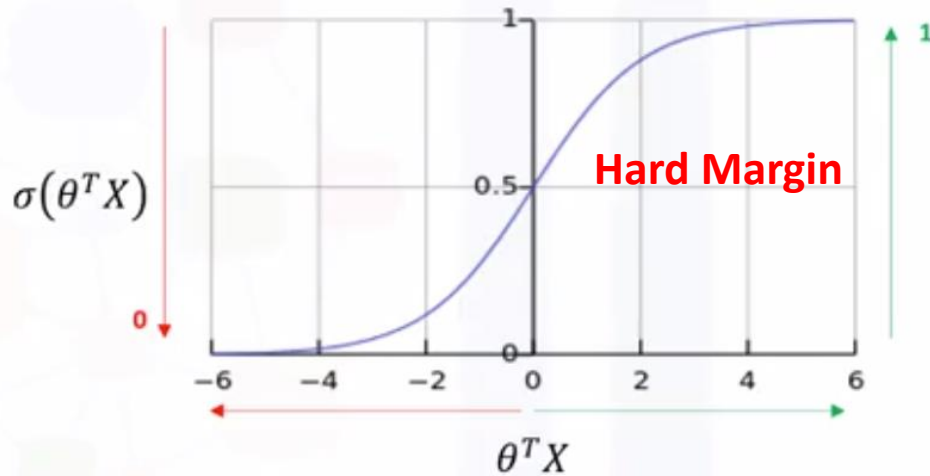
- **Model Architecture**

- Hyperplane Mapping through
  - Kernel (K)
- Decision Making through
  - Hypothesis function
    - Activation function
    - $h_{\theta}(K(X))$

# Model Kernel: Decision Making

## Logistic Regression

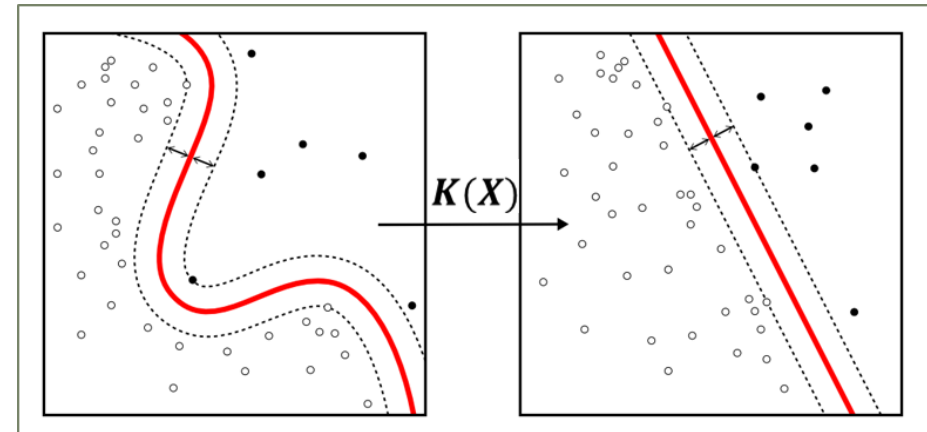
- Sigmoid / Logistic Kernel Function



$$K(X) = \theta^T X = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$
$$h_\theta(K(X)) = \sigma(\theta^T X) = \frac{1}{1 + \exp(\theta^T X)}$$
$$= \frac{1}{1 + \exp(\theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + b)}$$

## Support Vector Machine (SVM)

- **Class Partition** Kernels: Linear / Non-Linear



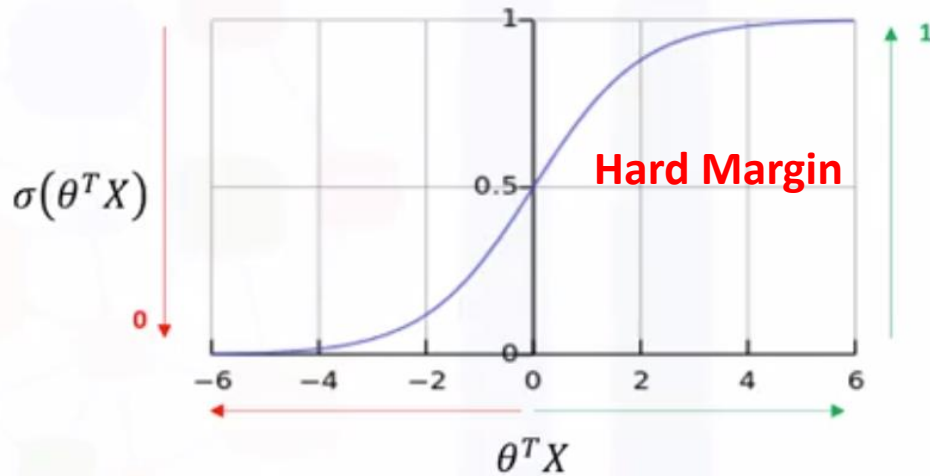
$$h_\theta(K(X)) = W^T K(X) - b = \theta^T K(X)$$

**Soft Margin**

# Model Kernel: Decision Making

## Logistic Regression

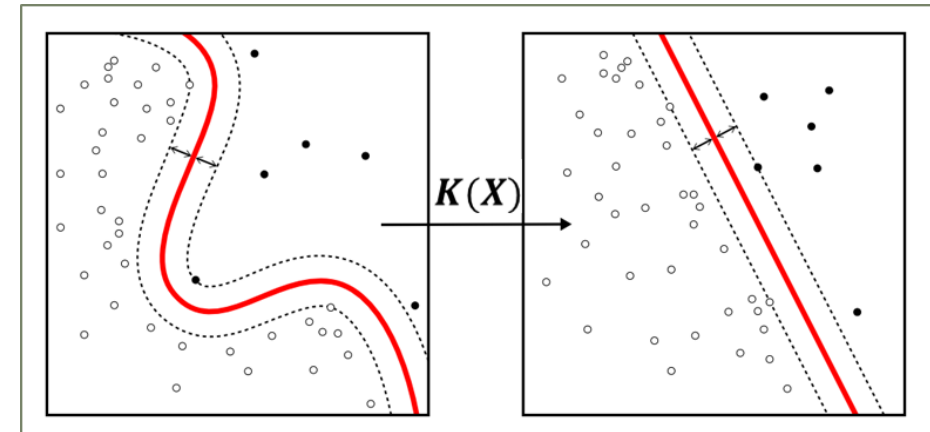
- Sigmoid / Logistic Kernel Function



$$P(y = 1|x) = h_{\theta}(K(X)) = \sigma(\theta^T X)$$
$$P(y = 0|x) = 1 - P(y = 1|x)$$
$$\hat{y} = \begin{cases} 1, & P(y = 1|x) \geq P(y = 0|x) \\ 0, & \text{Otherwise} \end{cases}$$

## Support Vector Machine (SVM)

- **Class Partition** Kernels: Linear / Non-Linear



$$h_{\theta}(K(X)) = \theta^T K(X) \quad \text{Soft Margin}$$

$$\hat{y} = \begin{cases} 1, & h_{\theta} \geq 0 \\ 0, & \text{Otherwise} \end{cases}$$

Best Optimizing  
hyperplane parameters  
( $\theta_0, \theta_1, \dots, \theta_N$  & kernel  
parameter)

Why do we need to  
tune these  
parameters?



# ML: Single Node Processor (Loss Function)

Hyperplane Parameter ( $\theta_i$ ) Optimization

$$L(x, y) = \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

- **Regression Loss** Function: **Error** Function
  - Normally use **L2 distance** between
    - $y$ : real output / desired output / ground truth
    - $h_{\theta}(K(X))$ : Hyperplane decision of kernel (K)

# ML: Single Node Processor (Loss Function)

## Hyperplane Parameter ( $\theta_i$ ) Optimization

- **Classification Loss** Function:

- distance measure
  - Dot product between  $y$  and  $h_\theta(K(X))$ 
    - $y \cdot h_\theta$

What's the difference between SVMs and Logistic Regression? (Revisited)

	<b>SVMs</b>	<b>Logistic Regression</b>
<b>Loss function</b>	Hinge loss	Log-loss
<b>Kernels</b>	Yes!	Yes!
<b>Solution sparse</b>	Often yes!	Almost always no!
<b>Semantics of learned model</b>	Linear model from "Margin"	Probability Distribution

# ML: Single Node Processor (Loss Function)

## Hyperplane Parameter ( $\theta_i$ ) Optimization

What's the difference between SVMs and Logistic Regression? (Revisited)

- **Classification Loss** Function:

- distance measure
  - Dot product between  $y$  and  $h_\theta(K(X))$ 
    - $y \cdot h_\theta = y \cdot z$

	<b>SVMs</b>	<b>Logistic Regression</b>
<b>Loss function</b>	Hinge loss	Log-loss

$$\ell(y, z) = \max(0, 1 - yz)$$

$$\ell_{\log}(y, z) = \ln(1 + e^{-yz})$$

# ML: Single Node Processor (Loss Function)

Hyperplane Parameter ( $\theta_i$ ) Optimization

$$\ell(y, z) = \max(0, 1 - yz)$$

$$\ell_{\log}(y, z) = \ln(1 + e^{-yz})$$

Optimizing on training data only

It is possible to cause **overfitted**

Not understand general inputs

# ML: Single Node Processor (Logistic)

Loss function with regularization

$$L(x, y) + \lambda (\textit{Regularization})$$

**Regularization Control ( $\lambda$ )**

$$C=1/\lambda$$

# ML: Single Node Processor

## Hyperplane Parameter Optimization

### Logistic Regression

#### Hard Margin

- Tuning Optimization Parameter:

- C: Regularization factor on
  - **weight magnitude**
- Kernel parameter
  - Linear -> No parameter

### Support Vector Machine (SVM)

#### Soft Margin

- Tuning Optimization Parameter:

- C: Regularization factor on
  - **Soft margin ( Slag variable)**
- Kernel parameter
  - Linear -> no parameter
  - RBF ->  $\gamma$ : gamma
  - *Polynomial -: d: polynomial degree*

# ML: Single Node Processor (Logistic)

Loss function with regularization

Regularization on Weight Magnitude

$$\ell_{\log}(y, z) = \ln(1 + e^{-yz})$$

$$L(x, y) \equiv \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2 + \lambda \sum_{i=1}^n \theta_i^2$$

L1 Regularization

$$\ell_{\log}(y, z) = \ln(1 + e^{-yz}) + \lambda \sum_{i=1}^n |\theta_i|$$

L2 Regularization

$$\ell_{\log}(y, z) = \ln(1 + e^{-yz})$$

C in Logistic

$$C=1/\lambda$$

Large C: less  $\lambda$  -> may be overfitted  
Less C: Large  $\lambda$  -> may be underfitted

# ML: Single Node Processor (Logistic)

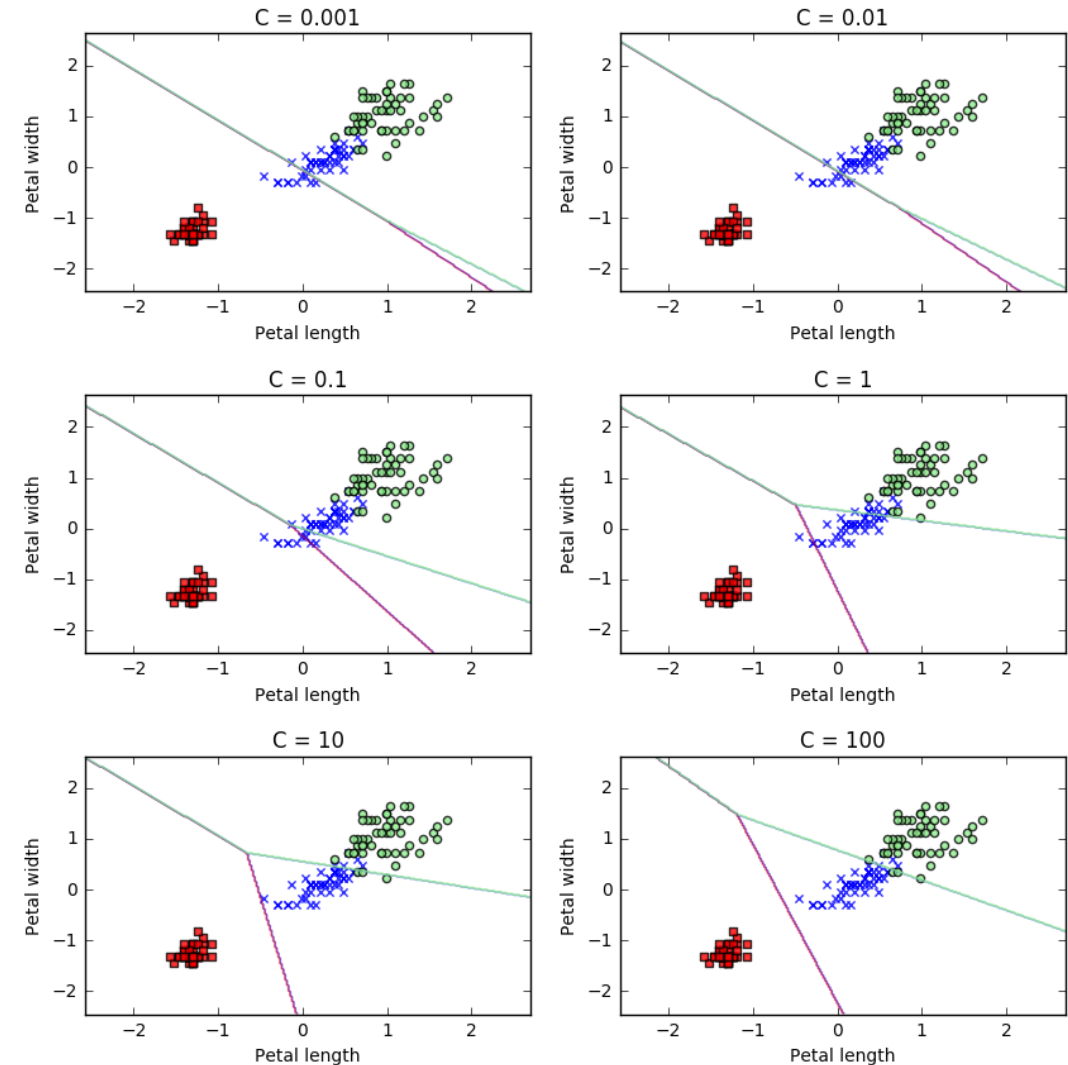
Loss function with regularization

$$\ell_{\log}(y, z) = \ln(1 + e^{-yz}) + \lambda \sum_{i=1}^n \theta_i^2$$

**C in Logistic**

**C=1/λ**

**Large C: less λ -> may be overfitted**  
**Less C: Large λ -> may be underfitted**





# ML: Single Node Processor (SVM)

Loss function with regularization

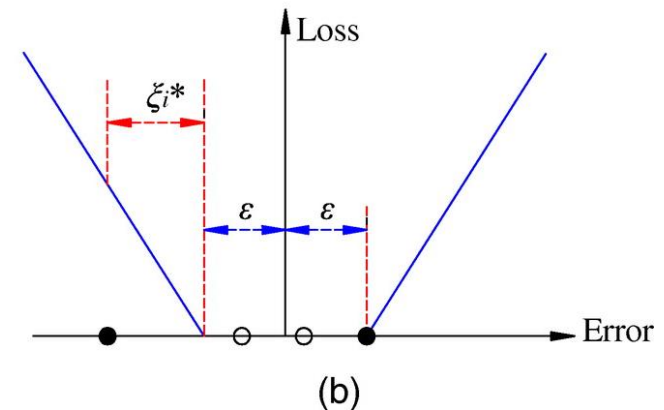
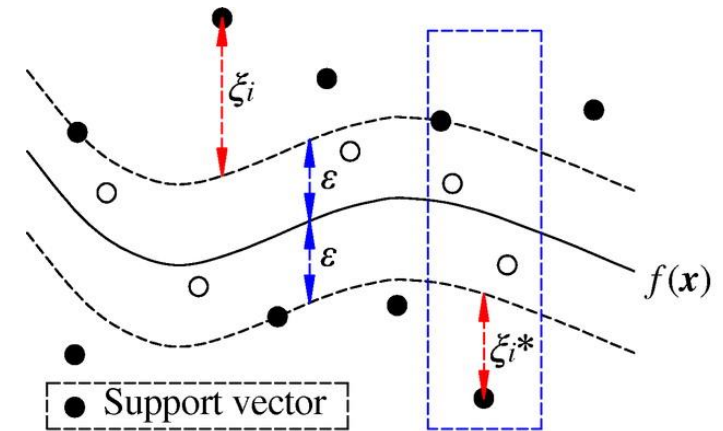
$$\ell(y, z) = \max(0, 1 - yz) + C \sum_i^N \xi_i$$

- L1 Regularization

C

Large C: may be overfitted  
Less C: may be underfitted

Regularization on Soft Margin



# ML: Single Node Processor (SVM)

Loss function with regularization

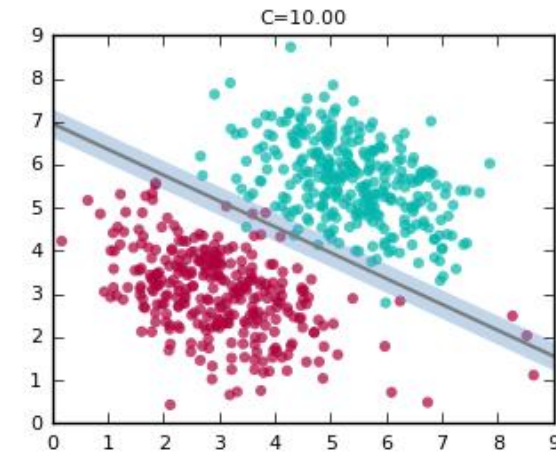
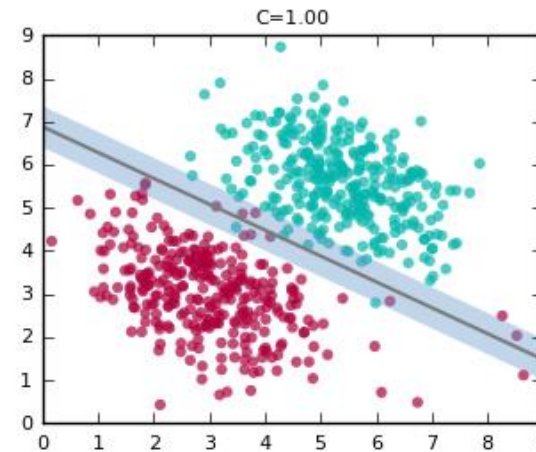
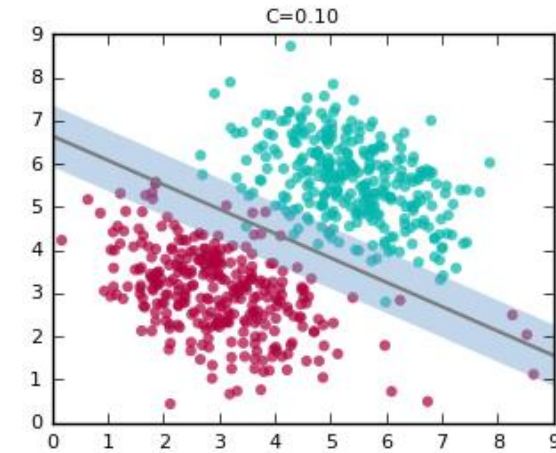
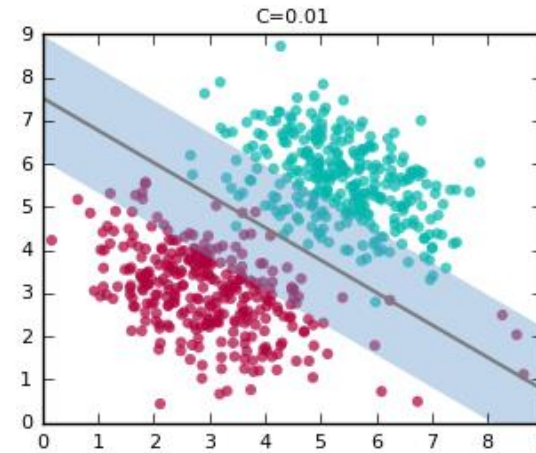
- L1 Regularization

$$\ell(y, z) = \max(0, 1 - yz) + C \sum_i^N \xi_i$$

**C in SVM**

Large C: may be overfitted  
Less C: may be underfitted

**Regularization on Soft Margin**



# ML: Single Node Processor (SVM)

## Loss function with regularization

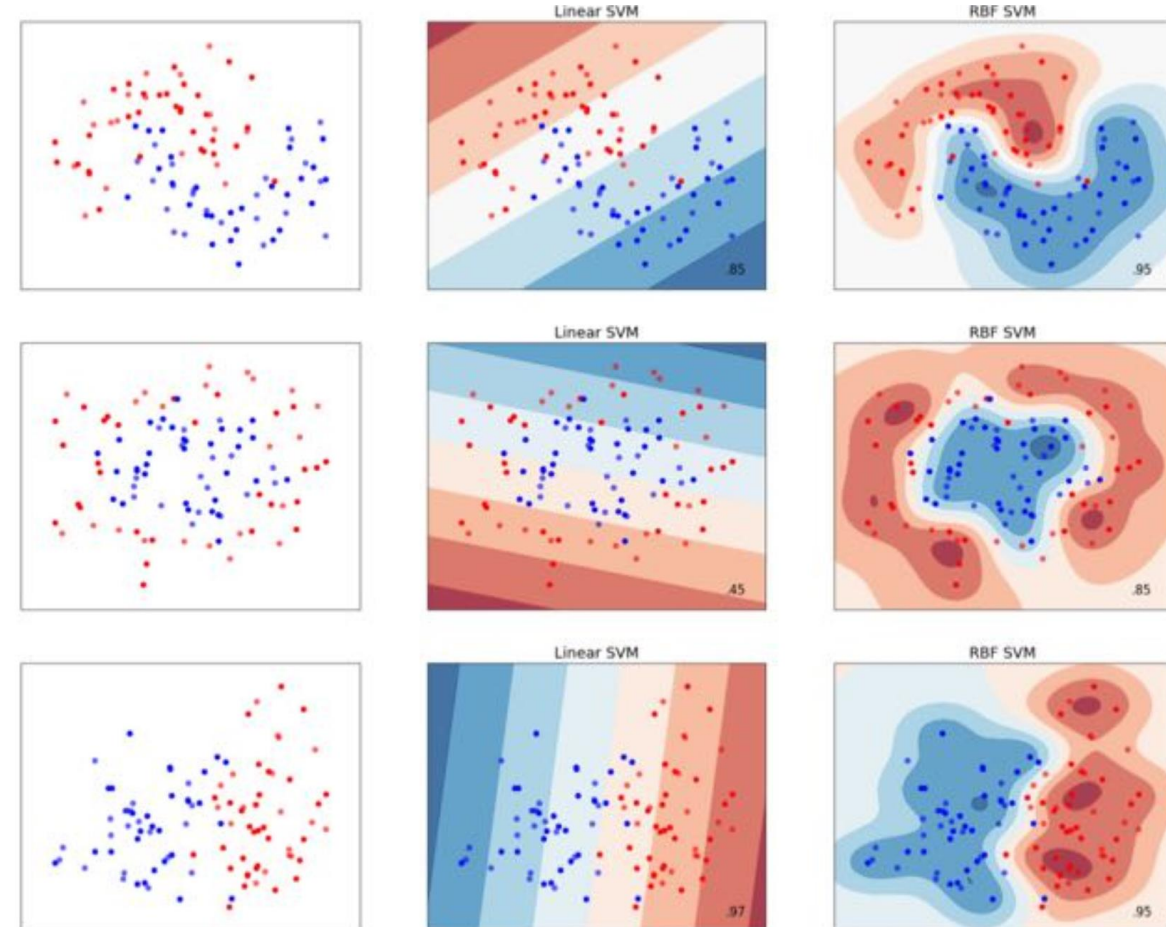
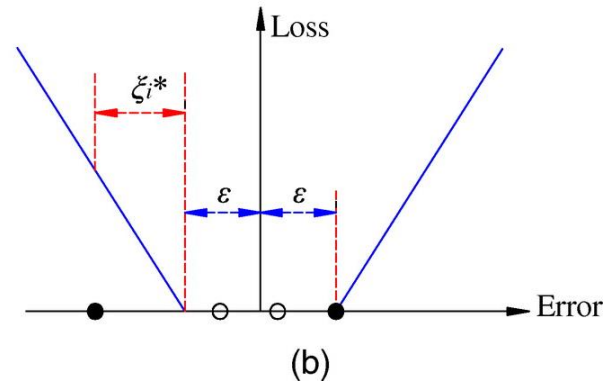
- L1 Regularization

$$\ell(y, z) = \max(0, 1 - yz) + \boxed{C \sum_i^N \xi_i}$$

**C in SVM**

Large C: may be overfitted  
Less C: may be underfitted

**Regularization on Soft Margin**



# ML: Single Node Processor (SVM)

## Loss function with regularization

- L1 Regularization

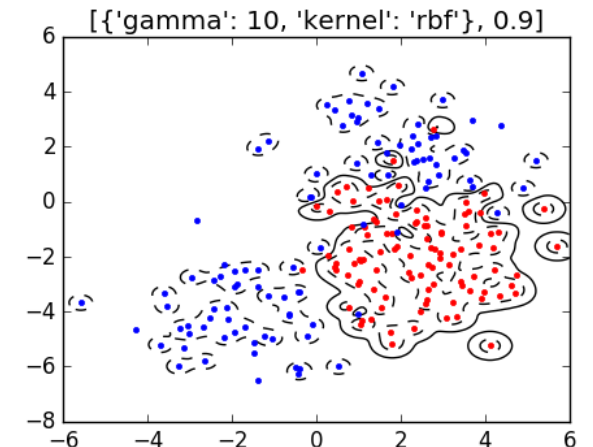
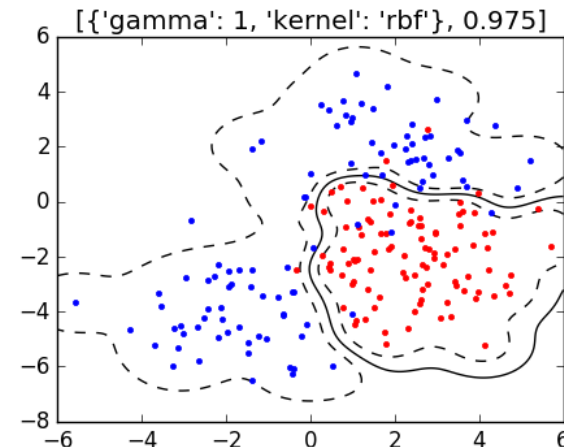
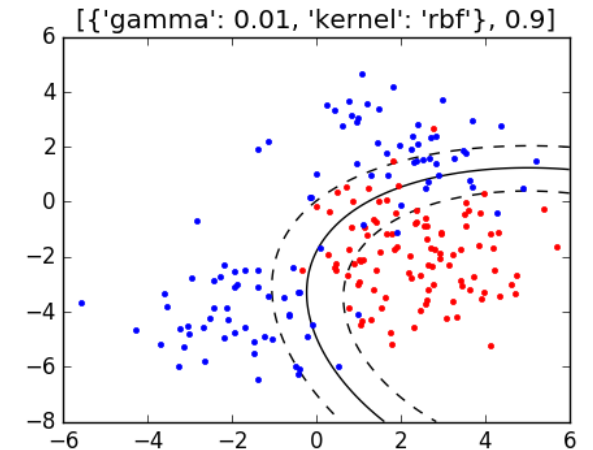
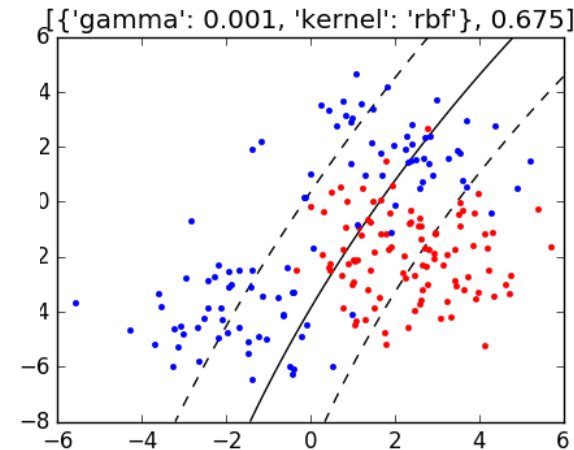
$$\ell(y, z) = \max(0, 1 - yz) + \boxed{C \sum_i^N \xi_i}$$

$$h_{\theta}(X) = \exp(-\gamma \|x_i - C_j\|^2)$$

**Gamma ( $\gamma$ ) in SVM**

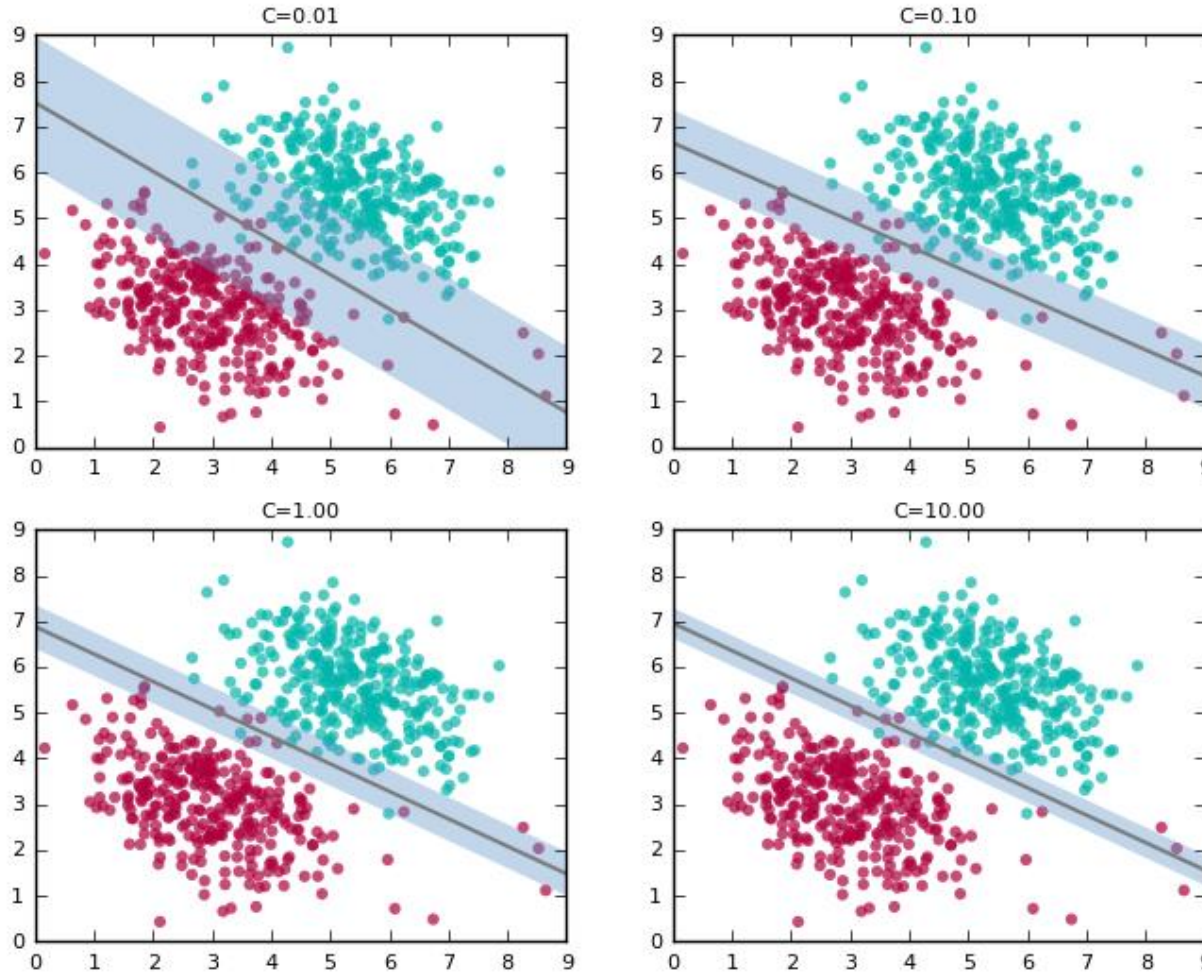
Large  $\gamma$  : may be overfitted  
Less  $\gamma$  : may be underfitted

**Regularization on Soft Margin**



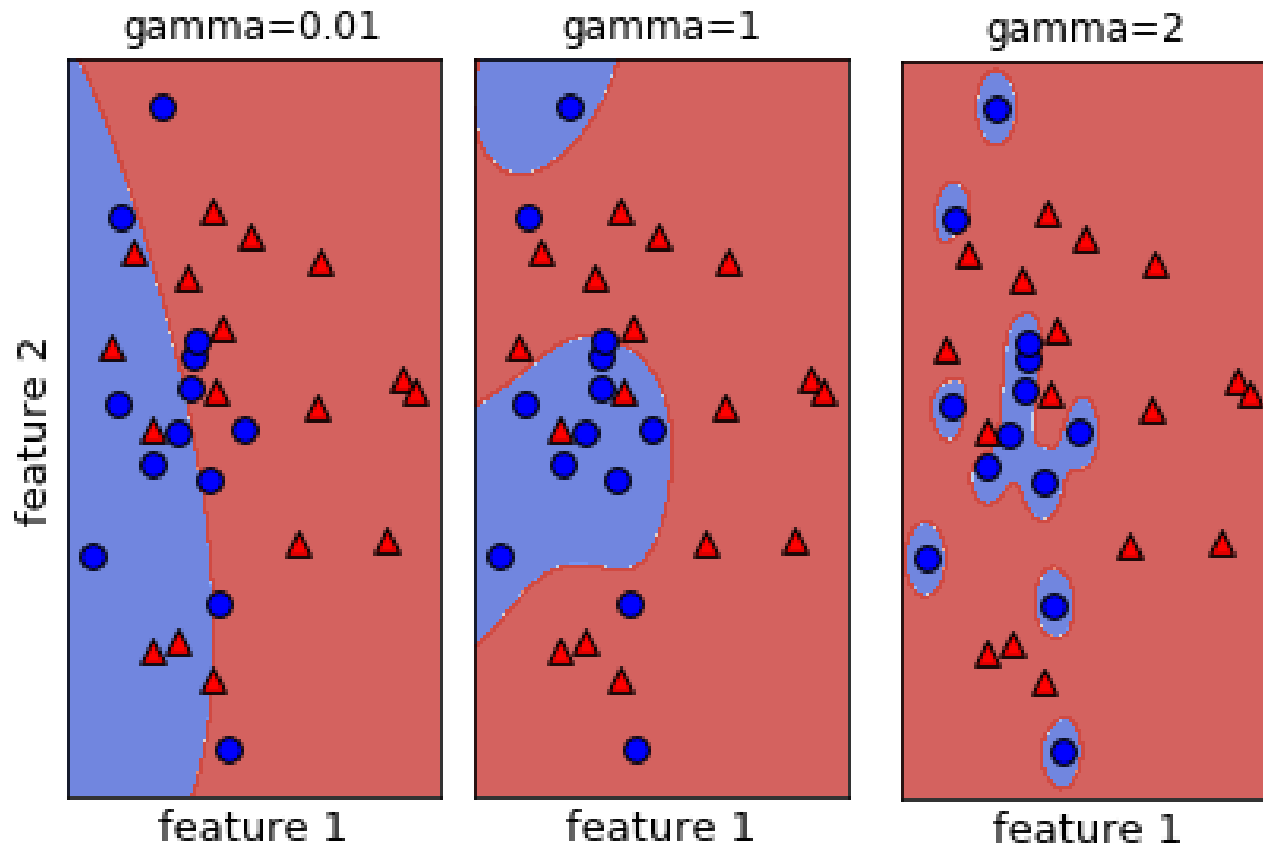


# SVC: Tuning parameters (Gamma)



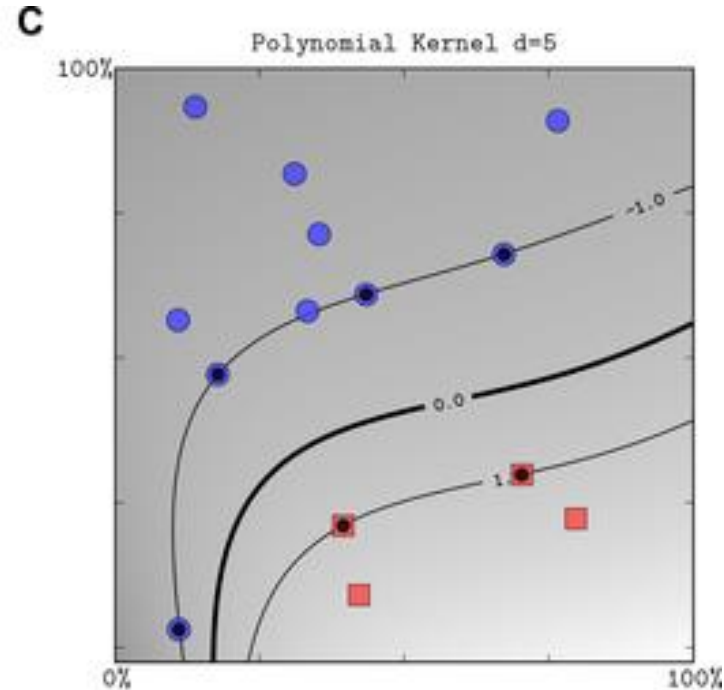
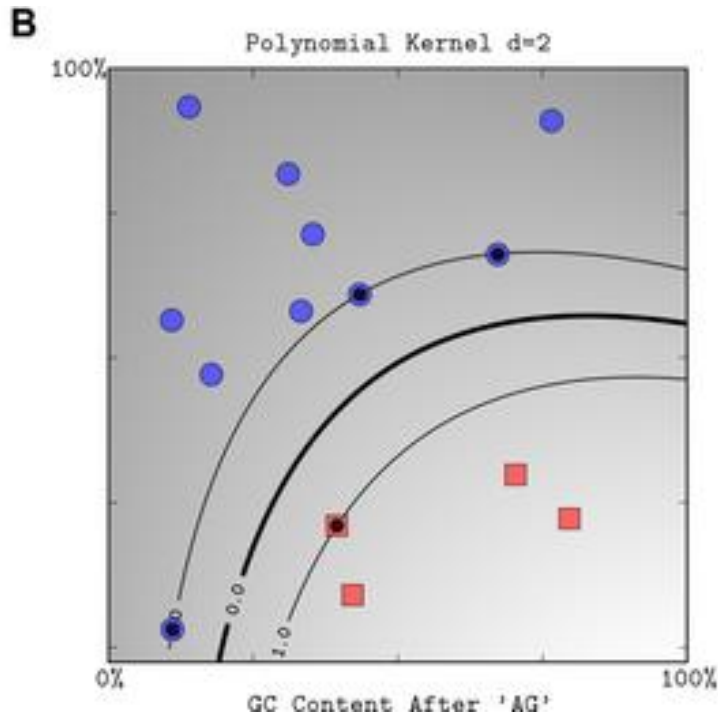
- General  $C$  parameters
  - Default:  $C=1$
- For **large** values of  $C$ ,
  - choose a **smaller-margin** hyperplane if that hyperplane does a better job of getting all the training points classified correctly.
- a very **small** value of  $C$ ,
  - look for a **larger-margin** separating hyperplane, even if that hyperplane misclassifies more points.

# SVC: Tuning parameters (Gamma)



- Gamma parameters
  - Kernel coefficient for 'rbf', 'poly' and 'sigmoid'
- Higher the value of gamma,
  - try to exact fit the as per training data set
  - i.e. generalization error and cause over-fitting problem.
- Default: 'auto' =  $1/n\_features$ 
  - 'scale' =  $1 / (n\_features * X.var())$

# SVC: Tuning parameters (Degree)



- Degree parameters
  - Kernel coefficient for 'poly'
- Higher the value of gamma,
  - try to exact fit the as per training data set
  - i.e. generalization error and cause over-fitting problem.
- Default: degree=3

# SVM Classification (SVC)

- `from sklearn.model_selection import train_test_split`
  - `X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20)`
- `from sklearn.svm import SVC`
  - `svclassifier = SVC(kernel='kernel_name', parameters )`
    - `svclassifier = SVC(kernel='rbf', C=100, gamma=0.01)`
  - `svclassifier.fit(X_train, y_train)`
  - `y_pred = svclassifier.predict(X_test)`
- `cv =` define cross validation you want to use with cross validation parameters
- `C_range = [0.01, 0.1, 10, 100]`
- `gamma_range = [0.001, 0.01, 0.1]`
- `param_grid = dict(gamma=gamma_range, C=C_range)`
- `grid = GridSearchCV(SVC(), param_grid=param_grid, cv=cv)`
- `grid.fit(X, y)`



# Multi-class Classification

# multi-class Classification

- Normally Logistic and SVC is only a binary classifier:
  - that is, it can only classify two classes at a time.
- Therefore, in order to classify multiple classes (more than two)
  - it has to train two or more binary classifiers
- Multi-class classification
  - “one-vs-one” approach
  - “one-vs-all” approach

# Multi-class SVM training 'one-vs-one'

- Assume  $n_{\text{class}}$  as  $n$  different classes
- Two pairs of classes are selected at a time and a binary classifier trained for them.
  - If  $n_{\text{class}}$  is the number of classes,
    - then  $n_{\text{class}} * (n_{\text{class}} - 1) / 2$  classifiers are constructed and
    - each one trains data from two classes.
- During the classification phases
  - all the binary classifiers are tested.
  - The class with the most votes wins.

# Multi-class SVM training 'one-vs-all'

- train one classifier per class
  - $n_{\text{class}}$  -> total  $n$  classifiers.
    - Ground truth = class  $i$ th
      - Class label  $i$ th -> positive
      - The rest -> negative
- Generic SVM might not work,
  - but still there are some workarounds.