# ESS2222

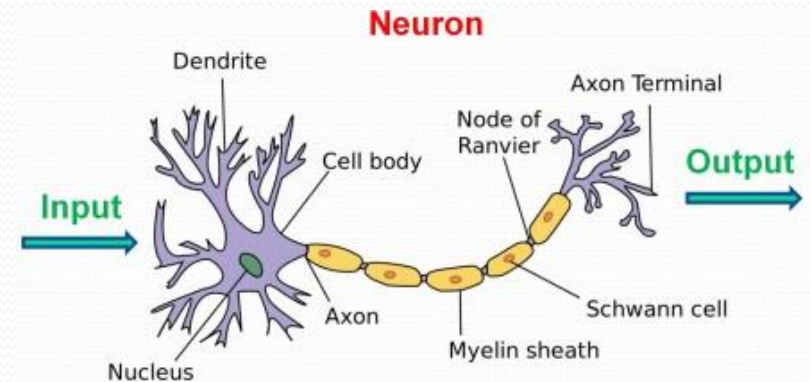# Lecture 5 – Support Vector Machine

*Hosein Shahnas*

*University of Toronto, Department of Earth Sciences,*

# Outline

- ❑ **Support Vector Machine (SVM)**
- ❑ **Soft Margin SVM**
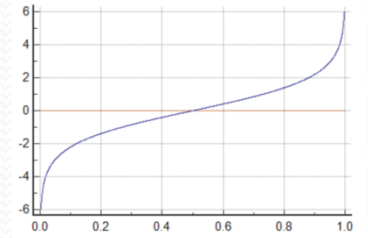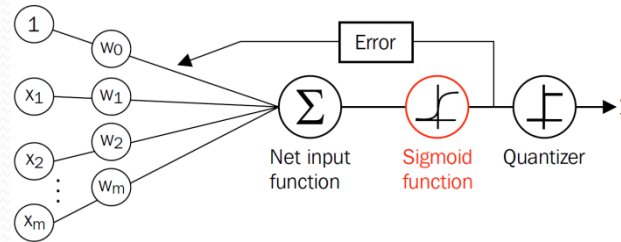- ❑ **Multiclass Problems**
- ❑ **Image Recognition**

## Logistic Regression

$$logit(p) = log\left[\frac{p}{(1-p)}\right], \qquad \text{p: } (0-1) \rightarrow \text{ logit(p): } (-\infty - \infty)$$
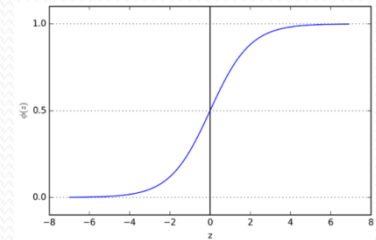
**The inverse function**

Logistic function: $\phi(z) = \frac{1}{1+e^{-z}}$, $\qquad$ z: $(-\infty - \infty) \rightarrow \phi(z)$: $(0-1)$ $\qquad$ logit(p)

$$\Delta wj = -\eta\frac{\partial J}{\partial w_j} = \eta\sum_i(y^i - \phi(z^i))x^i_j$$



$\phi(z)$

## Predicting Continuous Target Variables

**Classification:** **Credit approval (good/bad)**
**($x_i$, $y_i$)** $\qquad\qquad$ **$y_i \in [0, 1]$**

**Regression:** **Credit line (dollar amount)**
**($x_i$, $y_i$)** $\qquad\qquad$ **$y_i \in \mathbb{R}$**

3

$$E_{in}(w) = \equiv \frac{1}{N} \sum_{n=1}^{N}(w^T x_n - y_n)^2$$

$$\nabla_w E_{in}(w) = 0 \quad \rightarrow \quad X^T X w = X^T y$$

$$w = X^\dagger y \qquad \text{where} \quad X^\dagger = (X^T X)^{-1} X^T \qquad \text{pseudo-inverse}$$

## Linear Regression For Classification

1- Solve $w = X^\dagger y \quad (y \in R)$

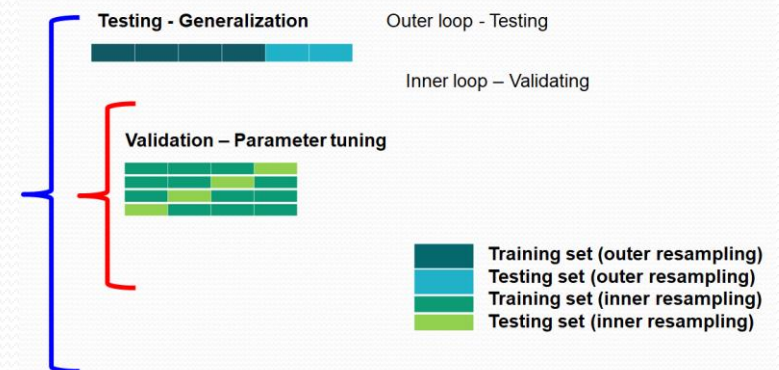2 - **Use the initial values for w obtained by the <span style="color:red">linear regression</span> method cas good starting point for <span style="color:red">classification</span>**

## K-Nearest Neighbours (KNN) Algorithm

$$d(X, Y) = \left(\sum_{i=1}^{n}|x_i - y_i|^p\right)^{1/p} \qquad \text{n features}$$



## Cross-Validation

# Support Vector Machine - More Details

In SVM the objective is to maximize the margin between different classes.
It can be shown that a model with a large margin Can show **better performance** on out-of-sample.

In other words in SVM we want to find **the weight vector w** such that not only classifies the samples **correctly**, but also **maximizes** the margin.

For any point on the class boundary we have $w_0 + w^T x = 0$.
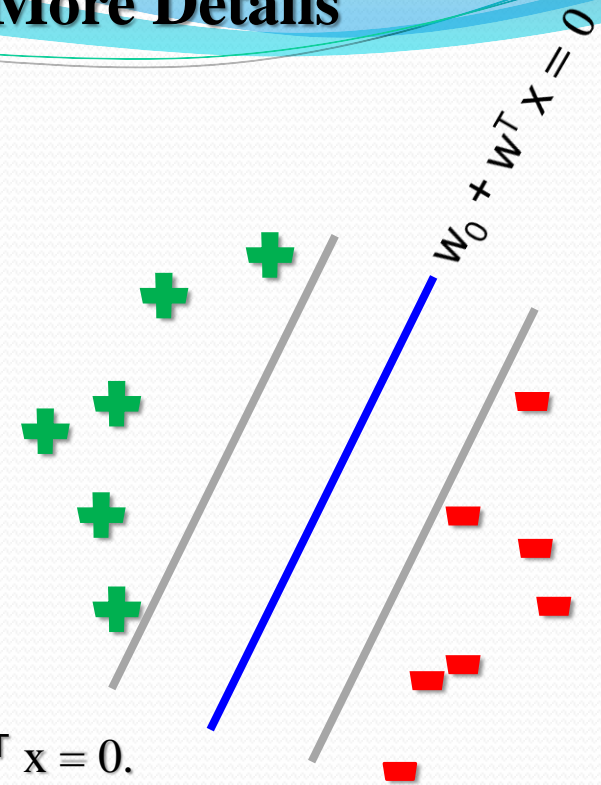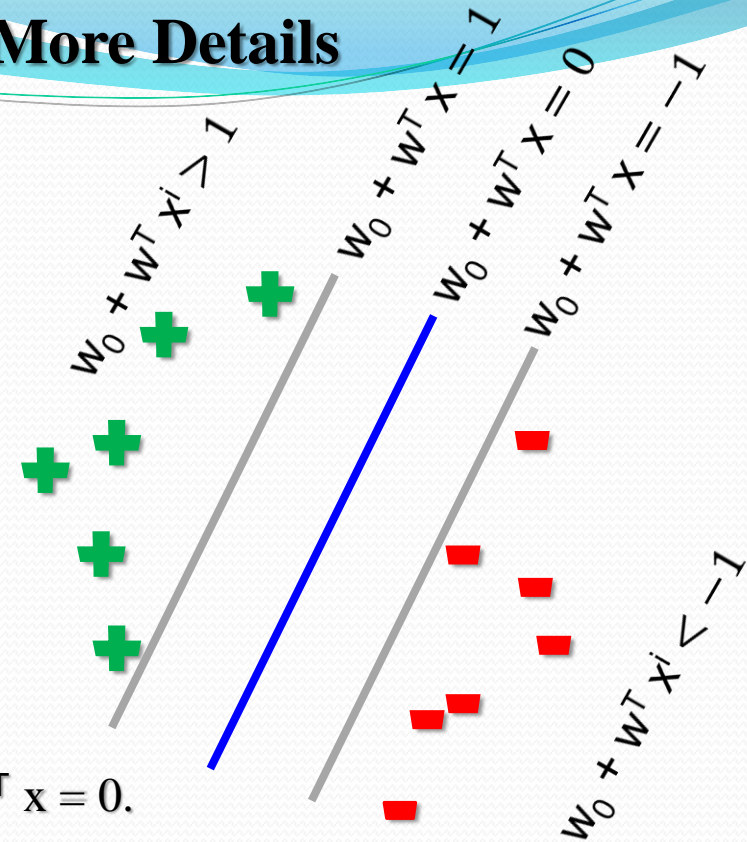
$w_0 + w^T x = 0$

# Support Vector Machine - More Details

In SVM the objective is to maximize the margin between different classes.
It can be shown that a model with a large margin Can show **better performance** on out-of-sample.

In other words in SVM we want to find **the weight vector w** such that not only classifies the samples **correctly**, but also **maximizes** the margin.

For any point on the class boundary we have $w_0 + w^T x = 0$.

$w_0 + w^T x^i \geq 1$

$w_0 + w^T x = 1$

$w_0 + w^T x = 0$

$w_0 + w^T x = -1$

$w_0 + w^T x^i \leq -1$

Suppose $x^n$ is a data point at the margin: $\left| w_0 + w^T x_n \right| = 1$
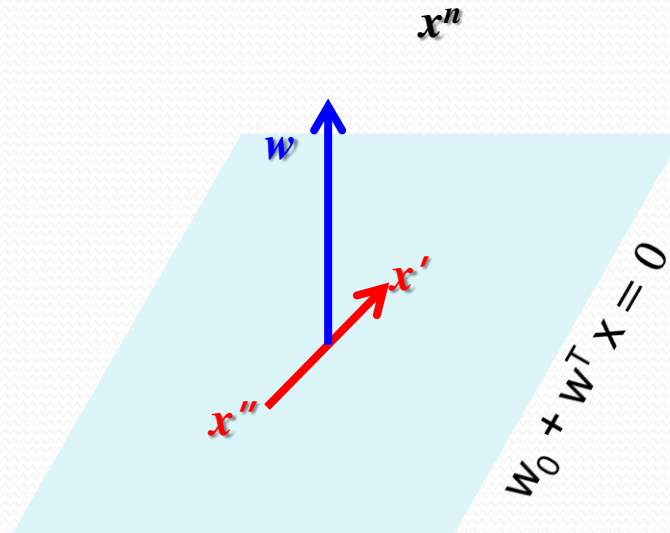
What is the distance between $x^n$ and the plane $w_0 + w^T x = 0$?

**The vector w is <span style="color:red">perpendicular</span> to the plane in the X-space.**

$w_0 + wT\ x' = 0$
$w_0 + wT\ x'' = 0$
$w^T\ (x' - x'') = 0 \quad \rightarrow \quad w^T \perp (x' - x'')$

$x^n$

$w$

$x'$

$x''$

$w_0 + w^T x = 0$

**Projection of (xⁿ - x) on w**

$$\widehat{w} = \frac{w}{\|w\|}$$

**Distance:** $\quad d = |\widehat{w}^T(x_n - x)| = \frac{1}{\|w\|}|w^T x^n - w^T x|$

$$d = \frac{1}{\|w\|}|w^T x^n + w_0 - w^T x - w_0| = \frac{1}{\|w\|}|w^T x^n + w_0|$$

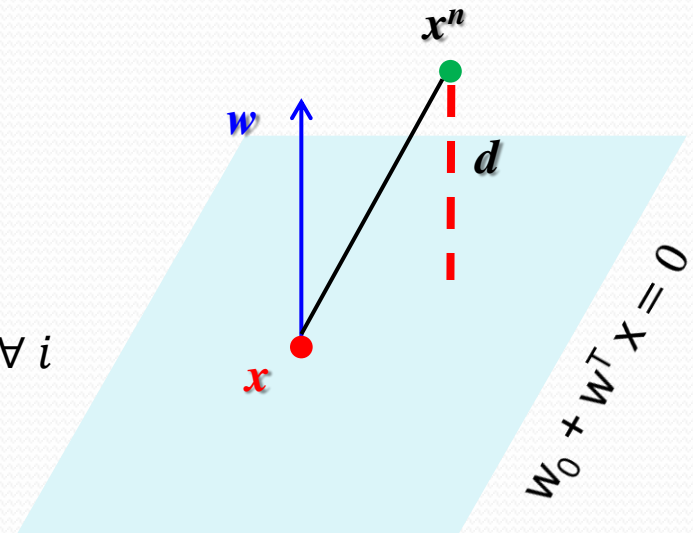**For xⁿ at the margin:**

$$d = = \frac{1}{\|w\|}|w^T x^n + w_0| = \frac{1}{\|w\|}$$

**This can be achieved by minimizing** $\frac{1}{2}\|w\|^2$

**Subject to the condition:** $\quad y^i(w_0 + w^T x^i) \geq 1 \quad \forall\, i$

**This is the confidence condition.**

$x^n$

$w$

$d$

$x$

$w_0 + w^T x = 0$

Suppose that there is a margin violation. Note that the sample may still be correctly classified with zero error.

With this violation the condition
$y^i ( w_0 + w^T x^i ) \geq 1 \quad \forall i$
will change to:
$y^i ( w_0 + w^T x^i ) \geq 1 - \xi^i \quad \forall i$
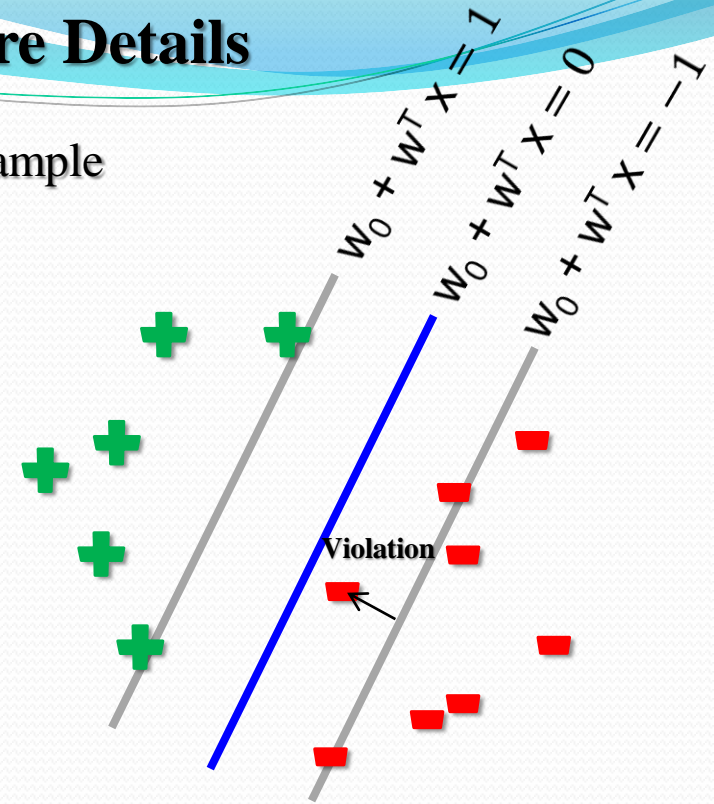where $\xi^i \geq 0$   slack vaiable

Total violation $= \sum_i \xi^i$

**New optimization:**

$$\frac{1}{2} \|w\|^2 \rightarrow \frac{1}{2} \|w\|^2 + C \sum_i \xi^i , \xi^i \geq 0$$

C (the slack coefficient) determines the relative importance of the first term wrt the second term.
C is obtained by cross-validation.

$w_0 + w^T x = 1$

$w_0 + w^T x = 0$

$w_0 + w^T x = -1$

**Violation**

# Soft Margin SVM - More Details

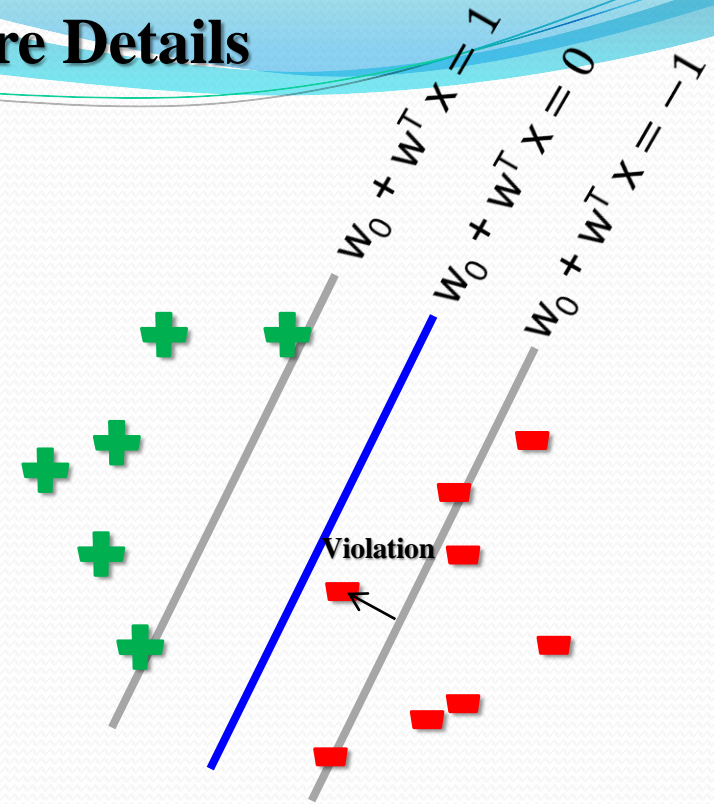**C $\to$ $\infty$ : Not violate the margins (hard margin)**
**C $\to$ 0 : Margin violations are allowed**

**Minimize** $\quad \frac{1}{2} \|w\|^2 + C \sum_i \xi^i$ , $\xi^i \geq 0$

**Subject to:**

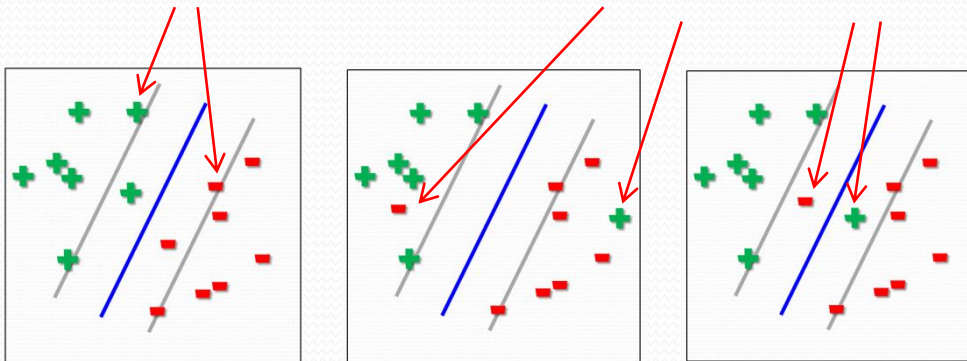$y^i ( w_0 + w^T x^i ) \geq 1 - \xi^i \quad \forall i \quad$ **where** $\xi^i \geq 0$

$\sum_i \xi^i = \sum_i \max\{0, (1 - yi(w.xi))\}$

**Violation**

**Types of violations:**

**Margin support vectors**   **Non-margin support vectors**
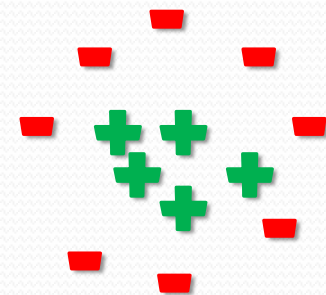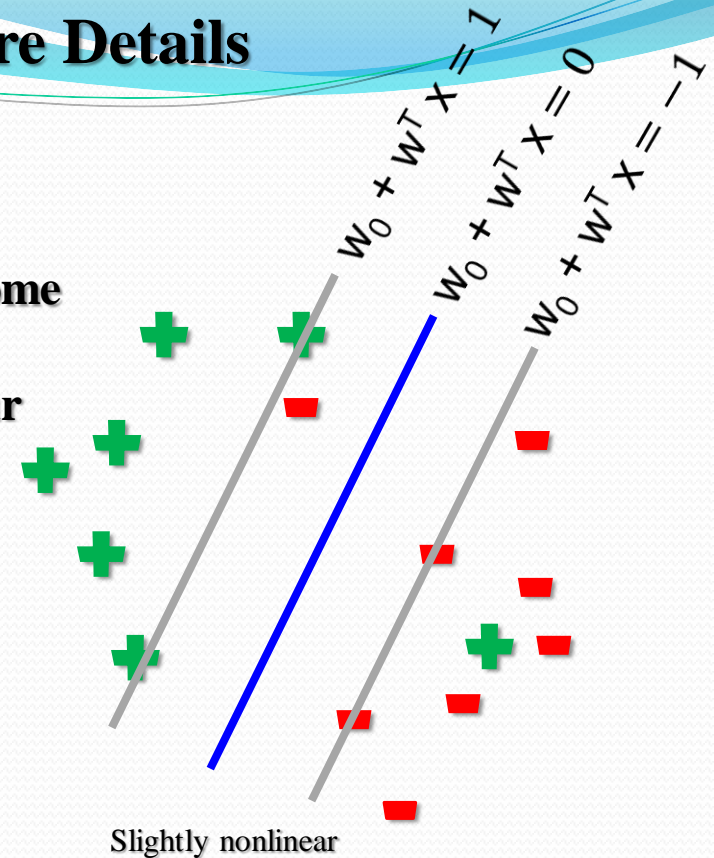
# Soft Margin SVM - More Details

**Remarks**

The method is called soft margin because it **allows** some **misclassifications.**
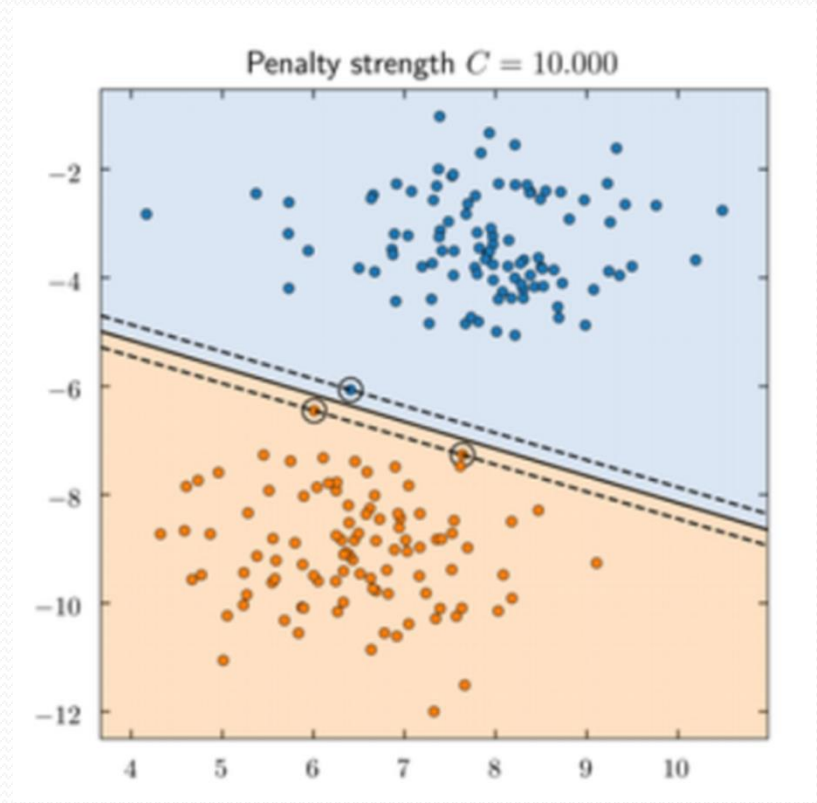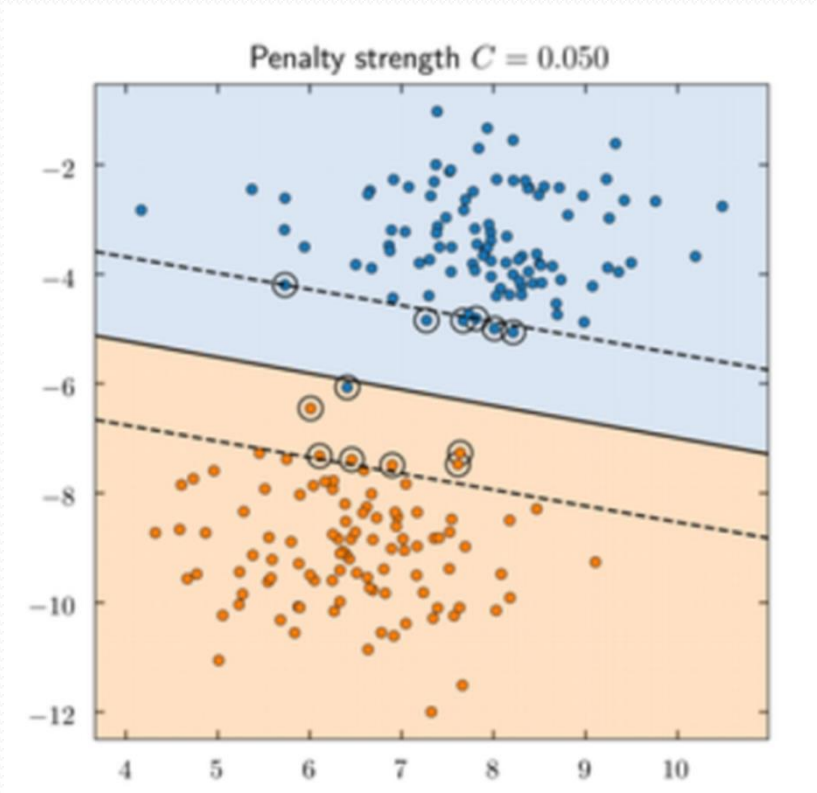Suppose the data is **slightly** nonlinear. This may occur when there is **noise** in data.

**Soft margin SVM** deals with **slightly** nonlinear problems.

**Kernel method** deals with **seriously** nonlinear problems.

But in reality we deal with practical problems where most datasets have the aspects of both, so we usually combine Kernel and soft margin SVM in almost all problems.

$w_0 + w^T x = 1$

$w_0 + w^T x = 0$

$w_0 + w^T x = -1$
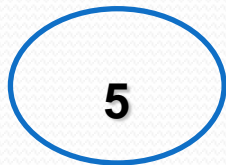
Slightly nonlinear
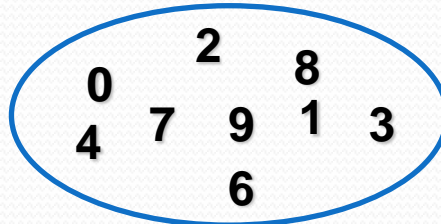
Seriously nonlinear

# Soft Margin SVM - More Details

# Multiclass Problems

## One-vs.-All (OvA)

The perceptron algorithm can be extended to multi-class classification-for example, through the One-vs.-All (OvA) technique.
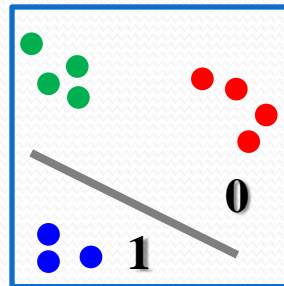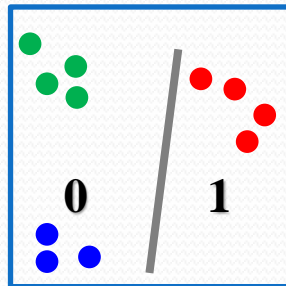


Class 1                 Class 0

**For 3-class problem**

# Multiclass Problems

## Softmax

This is a generalization of the logistic function to compute meaningful class-probabilities in multi-class settings (multinomial logistic regression).

$$\Phi(z_j) = \frac{1}{1+e^{-z_j}} \qquad \Longrightarrow \qquad \text{softmax}(z_j) = \frac{e^{z_j}}{\sum_{k=1}^{K} e^{z_k}}$$

$$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = \text{softmax} \begin{vmatrix} W_{1,1}x_1 + W_{1,2}x_1 + W_{1,3}x_1 + b_1 \\ W_{2,1}x_2 + W_{2,2}x_2 + W_{2,3}x_2 + b_2 \\ W_{3,1}x_3 + W_{3,2}x_3 + W_{3,3}x_3 + b_3 \end{vmatrix}$$

# Image Recognition

**Each pixel of image has an intensity in the range (0-255)**



**28x28 pixels**

## Handwritten Digits Classification
## Learning numbers 0-9

**n samples (images)**



**2D: 28x28 pixels**

**2D to 1D array**

**1D: 784**

**784**

**n images, one per line**

1
2
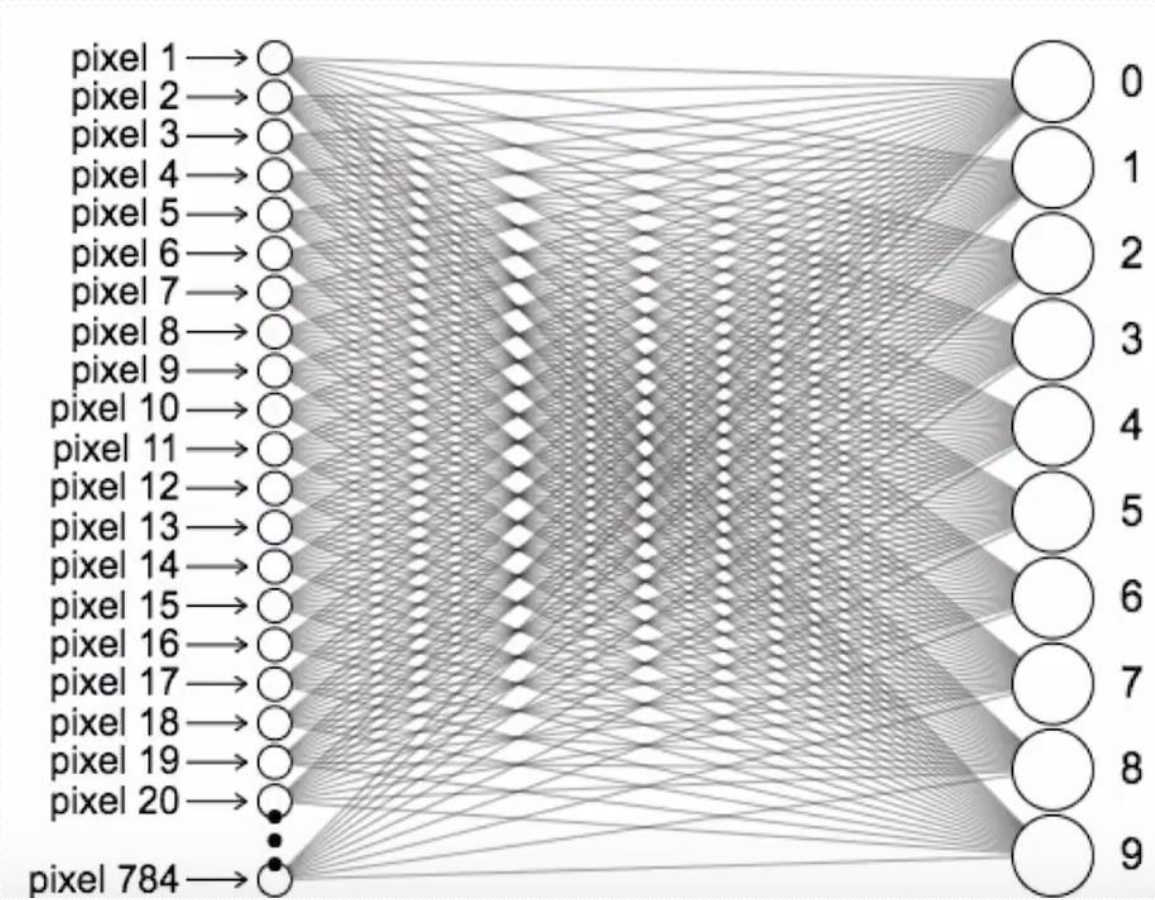3
...
...
...

...
n

**Learning**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

$$Q \begin{bmatrix} P_0 \\ P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \\ P_7 \\ P_8 \\ P_9 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} =$$

$$\sum_i P_i = 1$$

# SVM-Nonlinear Problems

**sklearn.svm.SVC**

**class sklearn.svm.SVC(C=1.0, kernel='rbf', degree=3, gamma='auto_deprecated', coef0=0.0, shrinking=True, probability=False, tol=0.001, cache_size=200, class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr', random_state=None)**

**C:** Penalty parameter of the error term **(default=1.0).**
**Kernel:** 'linear', 'poly', 'rbf', 'sigmoid', 'precomputed' or a callable **(default='rbf').**
**degree :** Degree of the polynomial kernel function ('poly'). Ignored by all other kernels **(default=3).**
**gamma:** Kernel coefficient **(default ='auto' which uses 1 / n_features).**
**Coef0:** Independent term in kernel function. It is only significant in 'poly' and 'sigmoid'.
**shrinking :** To save the training time, the shrinking technique tries to identify and remove some bounded elements **(default=True).**
**Probability:** Whether to enable probability estimates **(default=False).**
**tol :** Tolerance for stopping criterion **(default=1e-3).**
**cache_size:** The size of the kernel cache (in MB).
**Verbose:** Enable verbose output **(default: False)**
**max_iter:** Hard limit on iterations within solver, or -1 for no limit **(default=-1).**
**decision_function_shape:** 'ovo', 'ovr', **(default='ovr').**
**random_state:** The seed of the pseudo random number generator used when shuffling **(default=None; random generator is np.random).**