



**ESS2222**

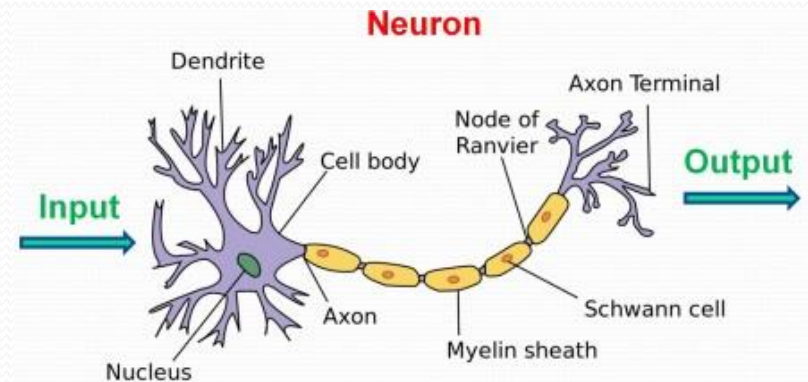
**Lecture 4 – Linear model**

*Hosein Shahnas*

*University of Toronto, Department of Earth Sciences,*

# Outline

- ❑ Logistic Regression
- ❑ Predicting Continuous Target Variables
- ❑ Support Vector Machine (Some Details)
- ❑ Nested Cross-Validation
- ❑ KNN - Algorithm



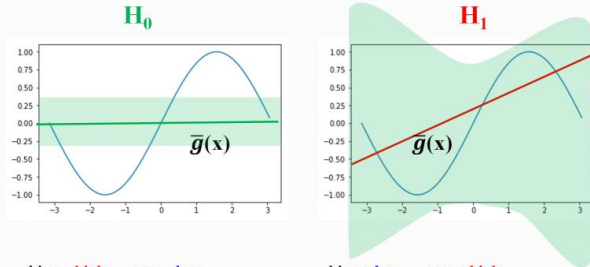
# Review of Lecture 3

## Bias-Variance Trade-off

$$E_D[E_X[(g^D(x) - f(x))^2]] = E_X[E_D[(g^D(x) - \bar{g}(x))^2]] + E_X[(\bar{g}(x) - f(x))^2]$$

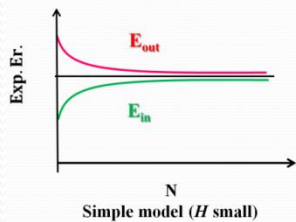
$$\bar{g}(x) = E_D[g^D(x)] \approx \frac{1}{k} \sum_k g^{D_k}(x) \quad \text{var} \quad \text{bias}$$

$$F(x) = \sin(2\pi x)$$

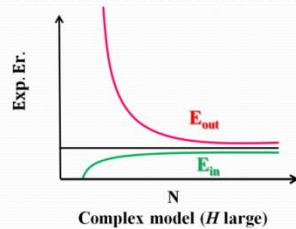


bias = high var = low

bias = low var = high



Simple model ( $H$  small)

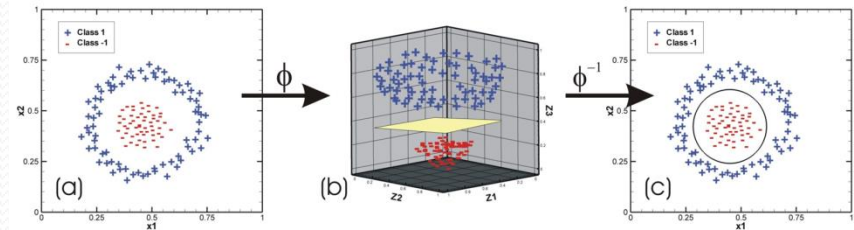


Complex model ( $H$  large)

## Cross-Validation

16%	16%	16%	16%	16%	20%
Training data	Training data	Training data	Training data	Validation	Testing data
Training data	Training data	Training data	Validation	Training data	Testing data
Training data	Training data	Validation	Training data	Training data	Testing data
Training data	Validation	Training data	Training data	Training data	Testing data
Validation	Training data	Training data	Training data	Training data	Testing data

## Nonlinearity



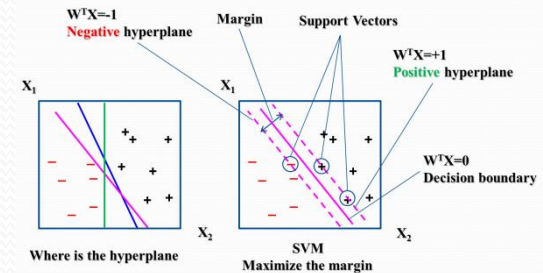
## L2-regression (Ridge)

$$J'_r(w) = \frac{1}{2} \sum_i (y^i - \phi(z^i))^2 + \frac{\lambda}{2} \sum_i w^2_i$$

## L1-regression (Lasso)

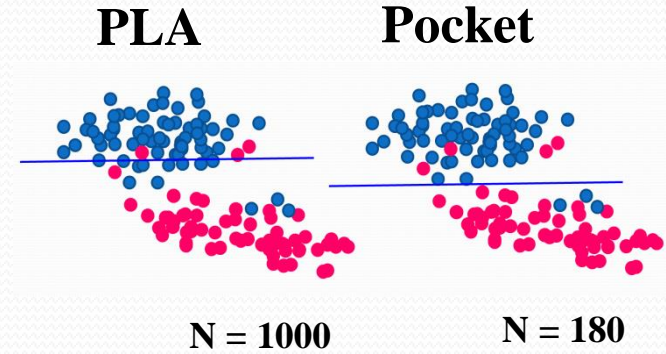
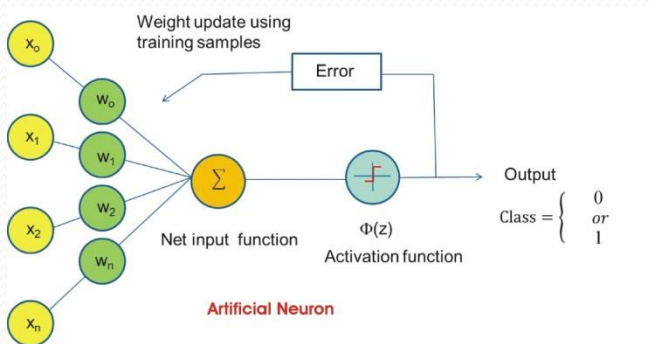
$$J'_l(w) = \frac{1}{2} \sum_i (y^i - \phi(z^i))^2 + \frac{\lambda}{2} \sum_i |w_i|$$

$$0 \leq \lambda < \infty$$

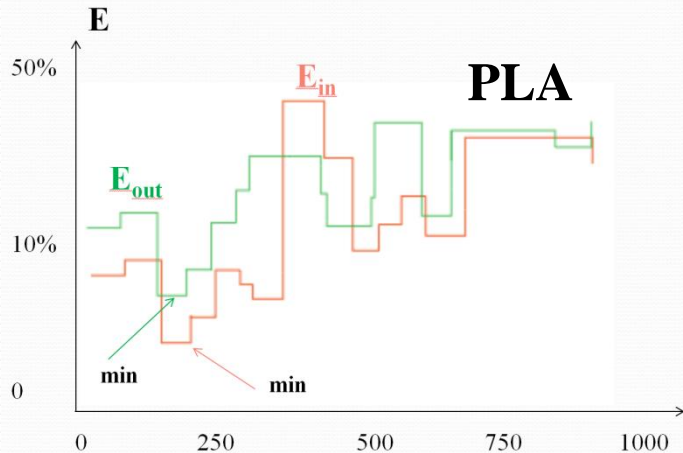


# Perceptron Learning Algorithm Versus Pocket Algorithm

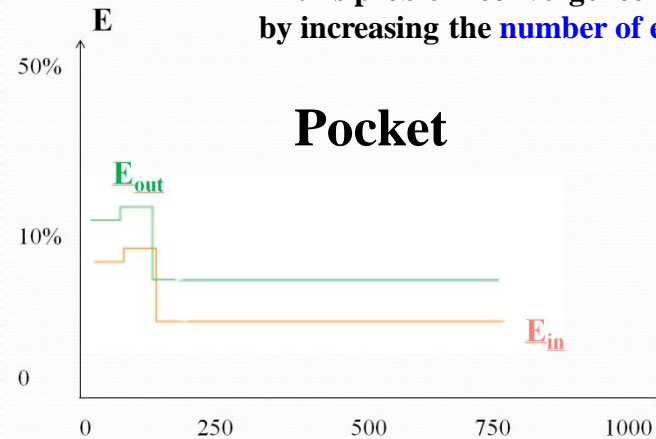
Suppose there exists small nonlinearity.



$$\mathbf{Z} = \mathbf{X} \cdot \mathbf{W} = \sum_0^n w_i x_i = w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$



In this problem convergence **cannot** be achieved by increasing the **number of epochs**



Store the best output and change it only for better outputs.

# Logistic Regression

## Modeling Class Probabilities via Logistic Regression

In order to overcome the convergence problem in Perceptron algorithm, we consider another **simple** yet **more powerful** algorithm for **linear and binary classification** problems. The method can be extended to multiclass classification via the **One-vs.-Rest (OvR)**.

Note that, this is an algorithm for **classification**, not **regression** (despite the name).

### Definition:

**Odds for (or odds of, or odds in favour):** Odds of event reflect the likelihood that the event will take place, while odds against reflect the likelihood that it will not.

### In gambling:

The odds are the ratio of payoff to stake, and do not necessarily reflect exactly the probabilities'

$$O_f = \frac{W}{L} = \frac{1}{O_a}, \quad O_a = \frac{W}{L} = \frac{1}{O_f}, \quad O_a \cdot O_a = 1 \quad W: \text{winning}, L: \text{loosing}$$

$$P = \frac{W}{(W+L)} = 1 - q, \quad q = \frac{L}{(W+L)} = 1 - p, \quad P + q = 1$$

# Logistic Regression

$$O_f = \frac{W}{L} = \frac{1}{O_a}, \quad O_a = \frac{L}{W} = \frac{1}{O_f}, \quad O_f \cdot O_a = 1 \quad \text{Odds}$$

$$p = \frac{W}{(W+L)} = 1 - q, \quad q = \frac{L}{(W+L)} = 1 - p, \quad p + q = 1 \quad \text{Probabilities}$$

$$O_f = \frac{p}{q} = \frac{p}{(1-p)} = \frac{(1-q)}{q}, \quad O_a = \frac{q}{p} = \frac{(1-p)}{p} = \frac{q}{(1-q)},$$

## Logistic Regression as a Probabilistic Model

The odds ratio for an event:  $O_f = \frac{p}{(1-p)}$

The **logit** function is defined as:  $\text{logit}(p) = \log \left[ \frac{p}{(1-p)} \right]$ , *usually natural log.*

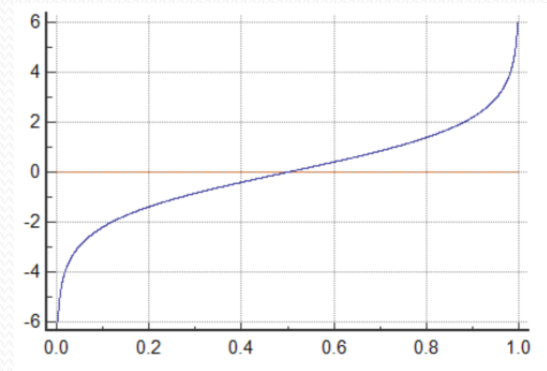
$p: (0 - 1) \rightarrow \text{logit}(p): (-\infty - \infty)$

# Logistic Regression

The **inverse** function of logit function is **logistic** function:

The inverse form of **logit** function is **logistic** function:  $\phi(z) = \frac{1}{1 + e^{-z}}$   
 $z: (-\infty - \infty) \rightarrow \phi(z): (0 - 1)$

$$\left\{ \begin{array}{l} y = \log\left(\frac{x}{1-x}\right), \quad y \rightarrow x \quad x \rightarrow y \quad x = \log\left(\frac{y}{1-y}\right) \\ \rightarrow e^x = \frac{y}{1-y} \rightarrow y = \frac{1}{1 + e^{-x}} \end{array} \right.$$



logit(p)

# Logistic Regression

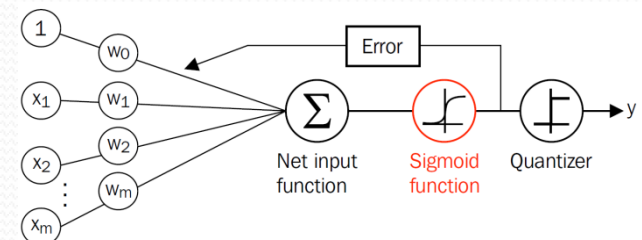
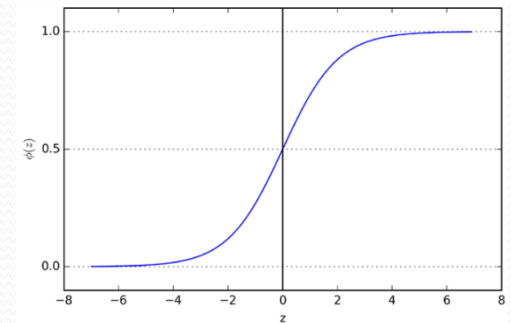
**Logistic function**, sometimes simply abbreviated as **sigmoid** function due to its characteristic **S-shape**.

$$z = x \cdot w = \sum_0^n w_i x_i = w_0 x_0 + w_1 x_1 + \dots + w_m x_m \quad \text{Net input}$$

**How do we interpret the sigmoid function?**

The output of the sigmoid function is interpreted as the probability of particular sample belonging to  $y$ , given  $x$  (i.e., class 1, given features  $x$  parameterized by the weights  $w$ ):

$$\phi(z) = P(y=1|x,w), \quad \phi(z) = \frac{1}{1 + e^{-z}}$$
$$\hat{y} = \begin{cases} 1 & \text{if } \Phi(z) \geq 0.5 \\ 0 & \text{otherwise} \end{cases} \quad \text{or} \quad \hat{y} = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{otherwise} \end{cases}$$





# Logistic Regression – The Cost Function

In Adeline algorithm we defined sum-squared-error cost function:

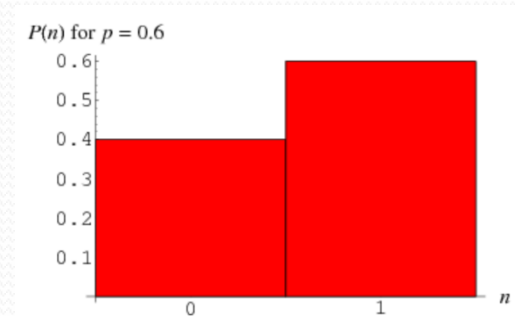
$$J(w) = \frac{1}{2} \sum_i (y^i - \phi(z^i))^2$$

In order to learn the **weights w** we **minimized** this function.

## Bernoulli Distribution

$$P(n) = \begin{cases} 1 - p & \text{for } n = 0 \\ p & \text{for } n = 1 \end{cases}$$

$$P(n) = p^n (1-p)^{(1-n)}$$



# Logistic Regression – The Cost Function

To derive the cost function for logistic regression, we define the likelihood function  $L$ , assuming that the individual samples in our dataset are independent of one another:

$$L(\mathbf{w}) = P(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \prod_{i=1}^n P(y^i, \mathbf{x}^i; \mathbf{w}) = \prod_{i=1}^n (\phi(\mathbf{z}^i))^{y^i} (1 - \phi(\mathbf{z}^i))^{(1-y^i)}$$

**It is easy to work with (natural) log of this function for two reasons:**

- Applying the log function **reduces** the potential for numerical **underflow**, which can occur if the likelihoods are very small,
- We can convert the product of factors into a summation of factors, which makes it **easier to obtain the derivative** of this function.

$$l(\mathbf{w}) = \log L(\mathbf{w}) = \sum_i^n y^i \log(\phi(\mathbf{z}^i)) + (1 - y^i) \log(1 - \phi(\mathbf{z}^i))$$

# Logistic Regression – The Cost Function

We can use an optimization algorithm such as **gradient ascent** to **maximize** this log-likelihood function. **Alternatively**, we can use **gradient decent** to **minimize** the following cost function:

$$J(\mathbf{w}) = -l(\mathbf{w}) = \sum_i^n -y^i \log(\phi(\mathbf{z}^i)) - (1 - y^i) \log(1 - \phi(\mathbf{z}^i))$$

To get a better grasp on this cost function, let's calculate for one single-sample instance:

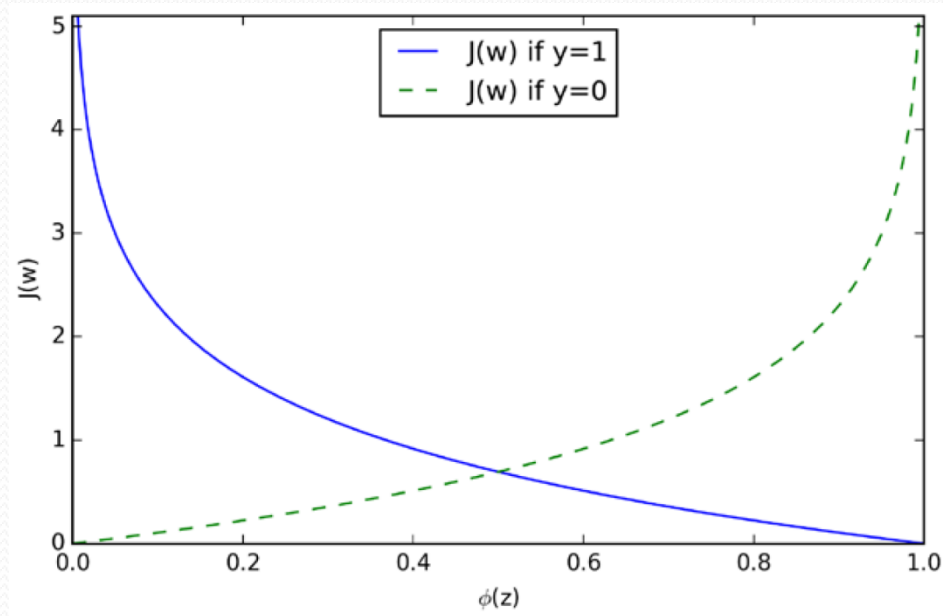
$$J(\mathbf{w}) = -y \log(\phi(\mathbf{z})) - (1 - y) \log(1 - \phi(\mathbf{z}))$$

It can be seen that the first term becomes zero if  $y = 0$  and the second term becomes zero if  $y = 1$ , respectively:

$$J(\phi(\mathbf{z}), y; \mathbf{w}) = \begin{cases} -\log(\phi(\mathbf{z})) & \text{if } y = 1 \\ -\log(1 - \phi(\mathbf{z})) & \text{if } y = 0 \end{cases}$$

# Logistic Regression – The Cost Function

$$J(\phi(z), y; w) = \begin{cases} -\log(\phi(z)) & \text{if } y = 1 \\ -\log(1 - \phi(z)) & \text{if } y = 0 \end{cases}$$



Note that  $J \rightarrow 0$  (blue line) if the sample is correctly predicted as class 1. Similarly,  $J \rightarrow 0$  (green dashed line) if the sample is correctly predicted as class 0. However, for wrong prediction  $J \rightarrow \infty$ . To minimize  $J$ , the wrong predictions should be penalized with an increasingly larger cost.

# Logistic Regression – The Cost Function

We can show that the weight update in **logistic regression** via gradient descent is the same we obtained in **Adaline**:

$$J(w) = -y \log(\phi(z)) - (1 - y) \log(1 - \phi(z))$$

$$\frac{\partial J}{\partial w_j} = \frac{\partial J}{\partial \phi(z)} \frac{\partial \phi(z)}{\partial w_j} = - \left( y \frac{1}{\phi(z)} - (1 - y) \frac{1}{(1 - \phi(z))} \right) \frac{\partial \phi(z)}{\partial w_j}$$

$$\frac{\partial \phi(z)}{\partial w_j} = \frac{\partial \phi(z)}{\partial z} \frac{\partial z}{\partial w_j}, \quad \text{for } \phi(z) = \frac{1}{1 + e^{-z}}$$

$$\frac{\partial \phi(z)}{\partial w_j} = \frac{e^{-z}}{(1 + e^{-z})^2} \frac{\partial z}{\partial w_j} = \frac{1}{1 + e^{-z}} \left( 1 - \frac{1}{1 + e^{-z}} \right) \frac{\partial z}{\partial w_j} = \phi(z)(1 - \phi(z)) \frac{\partial z}{\partial w_j}$$

$$\frac{\partial z}{\partial w_j} = x_j$$

$$\frac{\partial J}{\partial w_j} = - \left( y \frac{1}{\phi(z)} - (1 - y) \frac{1}{(1 - \phi(z))} \right) \phi(z)(1 - \phi(z)) x_j$$

# Logistic Regression – The Cost Function

$$\frac{\partial J}{\partial w_j} = - \left( y \frac{1}{\phi(z)} - (1 - y) \frac{1}{(1 - \phi(z))} \right) \phi(z) (1 - \phi(z)) x_j$$

$$= - y (1 - \phi(z)) x_j + ((1 - y) \phi(z) x_j = - (y - \phi(z)) x_j$$

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

## Using sklearn:

### Example -

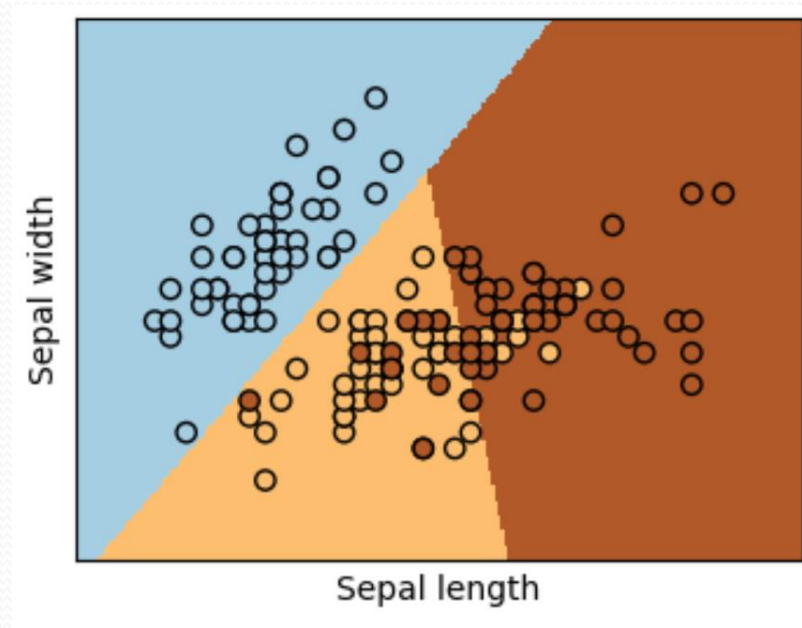
```
from sklearn.linear_model import LogisticRegression
Log_reg = LogisticRegression(C, random_state)
Log_reg.fit(X_train, y_train)
Log_reg.predict_proba(X_test[i])
```

>> 0.000, 0.063, 0.937 The probability that the test sample belongs to class **0**, **1**, and **2** are respectively **0.000**, **0.063**, and **0.937**.

Log\_reg.predict(X\_test[i]) will return class **2** for this example,

# Logistic Regression

## Iris Example



[https://scikit-learn.org/stable/auto\\_examples/linear\\_model/plot\\_iris\\_logistic.html](https://scikit-learn.org/stable/auto_examples/linear_model/plot_iris_logistic.html)

[https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html#sklearn.linear\\_model.LogisticRegression](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html#sklearn.linear_model.LogisticRegression)

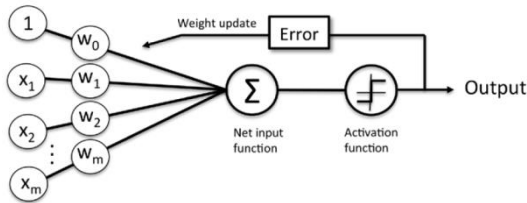
# Predicting Continuous Target Variables

$$z(\mathbf{x}) = \sum_{i=0}^m w_i x_i = \mathbf{w}^T \mathbf{x}$$

$$\phi(z) \equiv h(\mathbf{x})$$

$$\phi(z) = \text{sign}(z)$$

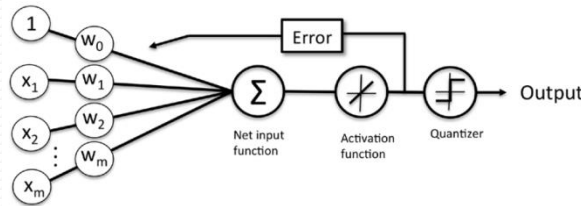
Perceptron



Discrete labels

$$\phi(z) = z + \text{Quantizer}$$

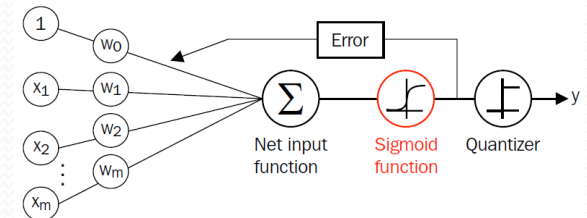
Linear regression  
(Adeline)



Discrete labels

$$\phi(z) = \frac{1}{1 + e^{-z}}, + \text{Quantizer}$$

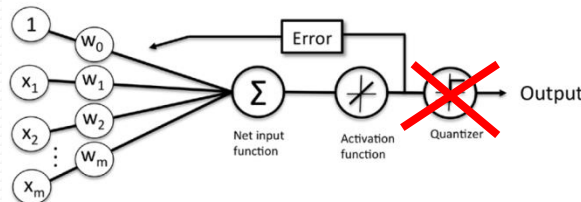
Logistic regression



Discrete labels

$$\phi(z) = z + \text{Quantizer}$$

Linear regression



Continuous outputs



# Linear Regression

## Examples:

**Classification:** Credit approval (good/bad)

$$(\mathbf{x}_i, y_i) \quad y_i \in [0, 1]$$

**Regression:** Credit line (dollar amount)

$$(\mathbf{x}_i, y_i) \quad y_i \in \mathbb{R}$$

## The error:

□ **In classification problem** we just count the number of wrong prediction (the frequency) compared to the total number of estimations.

$$E_{\text{in}}(\mathbf{h}) = \frac{1}{N} \sum_1^N \mathbb{I}[\mathbf{h}(x_n) \neq f(x_n)]$$

□ **In regression problem:**

$$E_{\text{in}}(\mathbf{h}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{h}(x_n) - y_n)^2$$

$$E_{\text{in}}(\mathbf{w}) = \equiv \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \equiv \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \quad \text{in-sample error}$$

# Linear Regression

## Minimizing $E(h)_{in}$

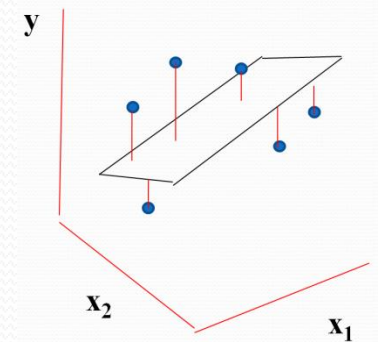
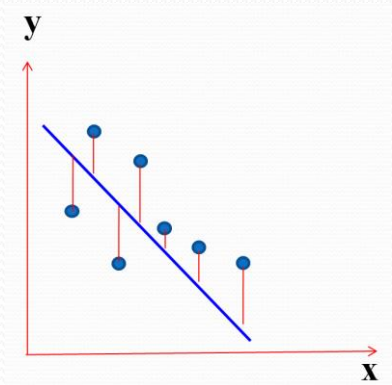
$$E_{in}(w) \equiv \frac{1}{N} \sum_{n=1}^N (w^T x_n - y_n)^2 \equiv \frac{1}{N} \|Xw - y\|^2$$

$$\nabla_w E_{in}(w) = \frac{2}{N} X^T (Xw - y) = 0 \quad \rightarrow \quad X^T Xw = X^T y$$

$$(X^T X)^{-1} X^T Xw = (X^T X)^{-1} X^T y \quad \rightarrow \quad w = (X^T X)^{-1} X^T y$$

$$w = (X^T X)^{-1} X^T y$$

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



# Linear Regression

## Minimizing $E(h)_{in}$

$$X^T = (X^T X)^{-1} X^T$$

$$x_i = [x_{i0}, x_{i1}, \dots, x_{im}]$$
$$w = [w_0, w_1, \dots, w_m]$$

$$X = \begin{bmatrix} x_{10}, x_{11}, \dots & x_{1m} \\ x_{20}, x_{21}, \dots & x_{2m} \\ \dots & \dots \\ x_{i0}, x_{i1}, \dots & x_{im} \\ \dots & \dots \\ x_{N0}, x_{N1}, \dots & x_{Nm} \end{bmatrix}$$
$$X^T = \begin{bmatrix} x_{10}, x_{20}, \dots & x_{N0} \\ x_{11}, x_{21}, \dots & x_{N1} \\ \dots & \dots \\ x_{1k}, x_{2k}, \dots & x_{Nk} \\ \dots & \dots \\ x_{1m}, x_{2m}, \dots & x_{Nm} \end{bmatrix}$$

$$wX^T = [y_1, y_2, \dots, y_N]$$

$$Xw^T = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{bmatrix}$$

# Linear Regression

## Minimizing $E(\mathbf{h})_{\text{in}}$

$$\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

$$v_1 \begin{bmatrix} \phantom{h} \\ \phantom{h} \\ \mathbf{h} \end{bmatrix} \mathbf{h} \begin{bmatrix} \phantom{h} \\ \phantom{h} \\ \mathbf{h}_2 \end{bmatrix} = \begin{bmatrix} \phantom{h} \\ \phantom{h} \\ \mathbf{h}_2 \end{bmatrix} v_1$$

$$\mathbf{X}^\dagger = \left( \begin{array}{c} \left[ \begin{array}{c} \phantom{h} \\ \phantom{h} \\ \mathbf{h} \end{array} \right] \left[ \begin{array}{c} \phantom{h} \\ \phantom{h} \\ \mathbf{h}_2 \end{array} \right] \\ \underbrace{\hspace{10em}}_{(m+1) \times N} \quad \underbrace{\hspace{10em}}_{N \times (m+1)} \end{array} \right)^{-1} \left[ \begin{array}{c} \phantom{h} \\ \phantom{h} \\ \mathbf{h}_2 \end{array} \right] \underbrace{\hspace{10em}}_{(m+1) \times N}$$

$N$ : Number of samples  
 $m$ : Number of features

# Linear Regression

## Minimizing $E(h)_{in}$

$$\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

$$\mathbf{X}^\dagger = \left( \begin{array}{c} \left[ \begin{array}{c} \phantom{\mathbf{X}^\dagger} \\ \phantom{\mathbf{X}^\dagger} \\ \phantom{\mathbf{X}^\dagger} \end{array} \right] \\ \underbrace{\hspace{10em}} \\ (m+1) \times (m+1) \end{array} \right)^{-1} \left[ \begin{array}{c} \phantom{\mathbf{X}^\dagger} \\ \phantom{\mathbf{X}^\dagger} \\ \phantom{\mathbf{X}^\dagger} \end{array} \right] \\ \underbrace{\hspace{10em}} \\ (m+1) \times N \end{array} \right)$$

$(m+1) \times N$

$N$ : Number of samples  
 $m$ : Number of features

# Linear Regression

## Minimizing $E(h)_{in}$

### Method:

- Construct matrix  $X$  and vector  $y$  using data  $(X_1, y_1), (X_2, y_2), \dots, (X_m, y_m)$

$$X^T = \begin{bmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix} \begin{matrix} \\ \\ \\ \\ \end{matrix} \begin{matrix} \\ \\ \\ \\ \end{matrix} \begin{matrix} \\ \\ \\ \\ \end{matrix} \begin{matrix} \\ \\ \\ \\ \end{matrix} \quad \begin{matrix} \\ \\ \\ \\ \\ \end{matrix} \begin{matrix} \\ \\ \\ \\ \\ \end{matrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ w_1 \\ \\ \\ w_m \end{bmatrix}$$

**N**

- Compute  $(X^T X)^{-1} X^T$

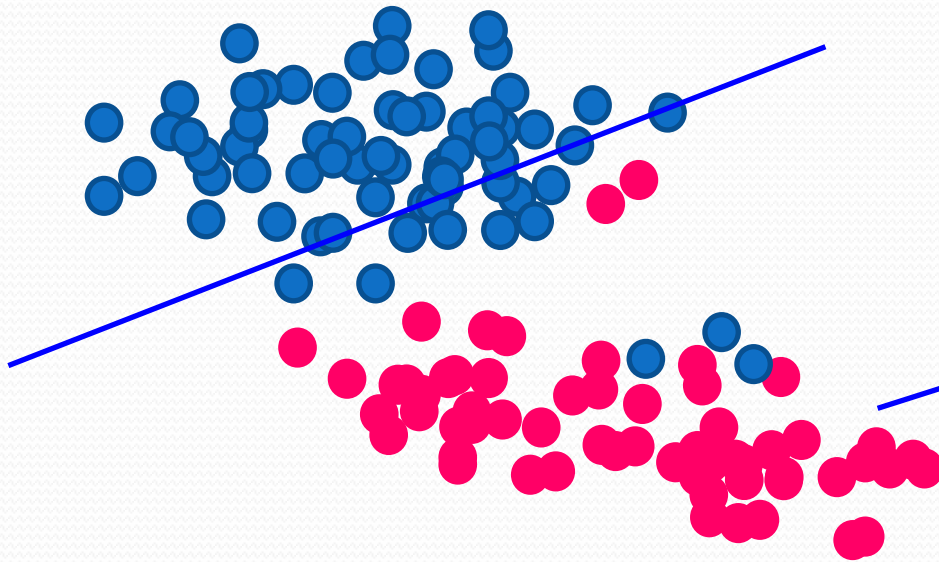
- Return  $w = X^T y$

# Linear Regression For Classification

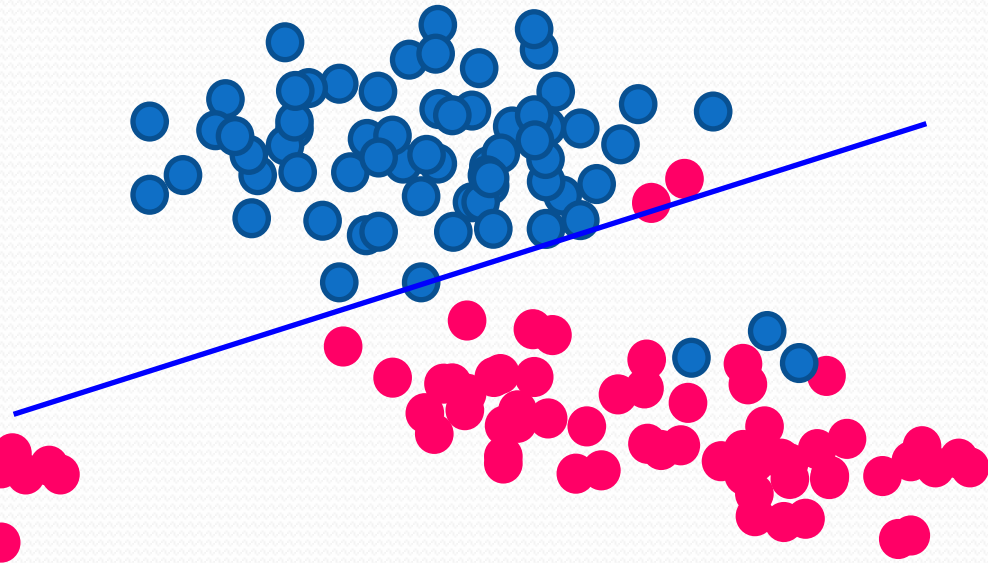
- ❑ In linear regression learning the model learn a real valued function  $y = f(x) \in R$
- ❑ In classification learning the binary-valued functions are also real valued  $y = \pm 1 \in R$
- ❑ We can use linear regression to obtain  $w$  ( $w^T x = y$ )
- ❑  $\text{Sign}(w^T x_j)$  likely agrees with  $y_j = \pm 1$
- ❑ The initial values for  $w$  obtained by the **linear regression** method can be good starting point for **classification** which minimizes the computation time (compared to the case we start from random values for  $w$ ).

# Linear Regression For Classification

Start with linear regression



Continue with classification





# Support Vector Machine

## Maximum margin classification with support vector machines

The objective in SVM-algorithm is to maximize the margin.

Margin: The distance between the separating hyperplane (decision boundary) and the training samples that are **closest** to this hyperplane, which are the so-called **support vectors**.

Models with decision boundaries with **large margins** tend to have a **lower generalization error** whereas models with small margins are **more prone to overfitting**.

Example:

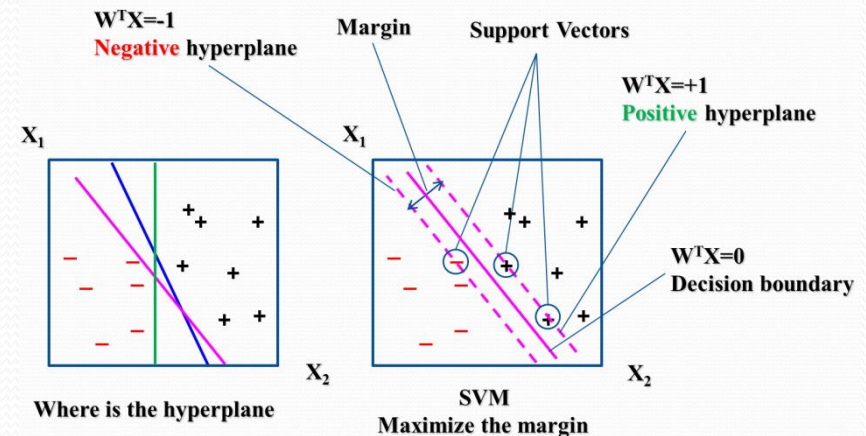
$$w_0 + w^T x_{\text{pos}} = 1$$

$$w_0 + w^T x_{\text{neg}} = -1$$

$$w^T (x_{\text{pos}} - x_{\text{neg}}) = 2$$

$$\frac{w^T (x_{\text{pos}} - x_{\text{neg}})}{\|w\|} = \frac{2}{\|w\|}$$

$$\|w\| = \sqrt{\sum_{j=1}^m w_j^2}$$



# Support Vector Machine

$\frac{w^T (x_{pos} - x_{neg})}{\|w\|}$  **margin**: distance between positive and negative hyperplanes

The **objective** of SVM algorithm is **maximization of** this margin by maximizing of  $\frac{2}{\|w\|}$  under the **constraint** that the samples are classified **correctly**:

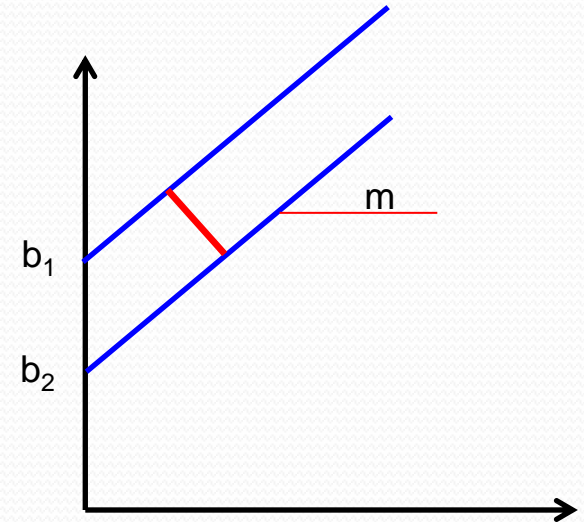
$$w_0 + w^T x^i \geq 1 \quad \text{if } y^i = 1 \quad (1)$$

$$w_0 + w^T x^i \leq -1 \quad \text{if } y^i = -1 \quad (2)$$

$$y^i (w_0 + w^T x^i) \geq 1 \quad \forall i \quad (3) \text{ compact form}$$

**Maximizing**  $\frac{2}{\|w\|}$   $\rightarrow$  **Minimizing**  $\frac{1}{2} \|w\|$

So  $\frac{1}{2} \|w\|$  is minimized with the condition of  $y^i (w_0 + w^T x^i) \geq 1 \quad \forall i$  using **quadratic Programming**.



$$d = \frac{|b_1 - b_2|}{\sqrt{1+m^2}}$$

# Soft-margin Classification

## Dealing with the Nonlinearly Separable Case Using Slack Variables

Introduced by Vladimir Vapnik in 1995

By introducing the **positive slack** variable, the linear constraints are **relaxed** for nonlinearly separable data to allow convergence of the optimization in the presence of misclassifications under the **appropriate cost penalization**.

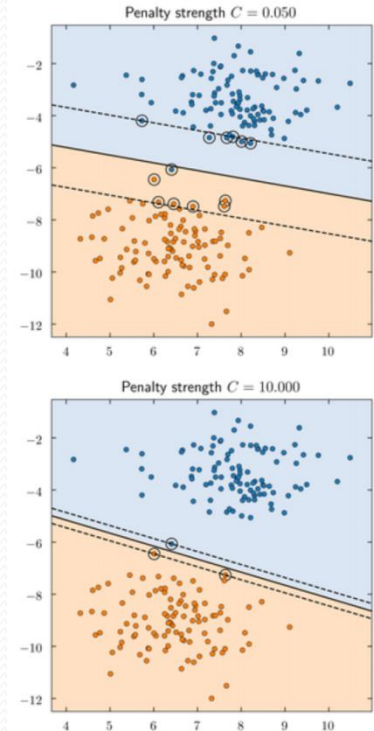
$$\begin{aligned} w^T x^i &\geq 1 && \text{if } y^i = 1 - \xi^i \\ w^T x^i &\leq -1 && \text{if } y^i = 1 + \xi^i \end{aligned}$$

The new objective to be minimized is:

$$\frac{1}{2} \|w\| + C \sum_i \xi^i$$

C(penalty strength): Controls the penalty for misclassification.

Increasing the value of C increases the bias and lowers the variance of the model. If the **penalty is small** the number of training points that define the **separating hyperplane is large**.



# Nested Cross-Validation

## Cross-Validation (CV):

- 1- **Train** a model on a subset of data, and **validate** the trained model on the remaining data (validation data).
- 2- **Repeat** this for all splits and average the validation error. This gives an estimate of the generalization performance of the model. **Chose the parameter leading to the best score.**

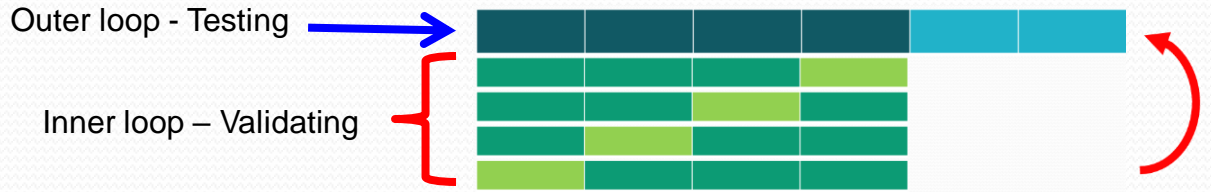
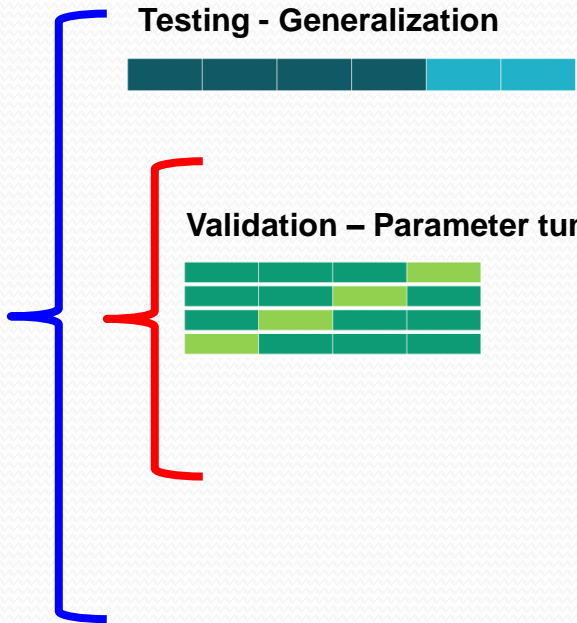
Model selection in CV uses the same data to tune model parameters and evaluate model performance. Information may thus “**leak**” into the model and overfit the data.

Choosing the parameters that maximize the scores in non-nested CV **biases** the model to the dataset, yielding an overly-optimistic score.

20%	20%	20%	20%	20%
Training data	Training data	Training data	Training data	Validation
Training data	Training data	Training data	Validation	Training data
Training data	Training data	Validation	Training data	Training data
Training data	Validation	Training data	Training data	Training data
Validation	Training data	Training data	Training data	Training data

5-fold validation

# Nested Cross-Validation



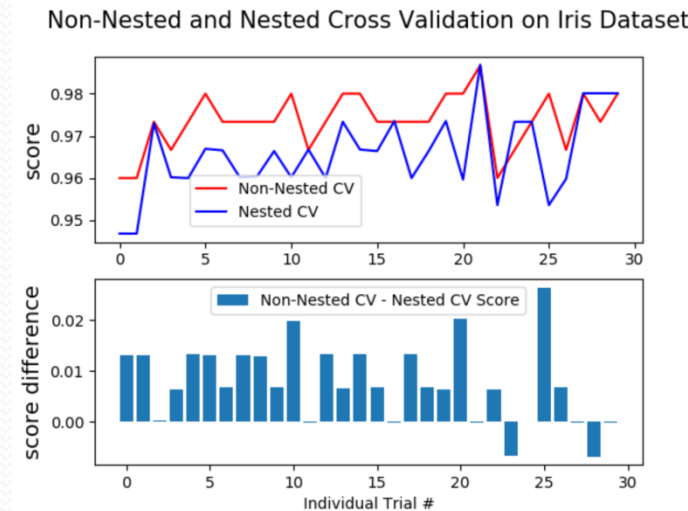
- Training set (outer resampling)
- Testing set (outer resampling)
- Training set (inner resampling)
- Testing set (inner resampling)

# Nested Cross-Validation

## Nested Cross-Validation:

- 1- The **hyper-parameter tuning** is carried out in the inner loop.
- 2- In the outer loop the **generalization** error of the underlying layer is estimated. The inner loop is responsible for model selection / hyper-parameter tuning., while the outer loop is for error estimation (test set).

In **sklearn** the hyper-parameters are turned by **GridSearchCV** in the inner loop. In the outer loop generalization error is estimated by averaging test set scores (using **cross\_val\_score**) over several dataset splits



Non-nested and nested cross-validation strategies on a classifier of the iris data set.

[https://scikit-learn.org/stable/auto\\_examples/model\\_selection/plot\\_nested\\_cross\\_validation\\_iris.html](https://scikit-learn.org/stable/auto_examples/model_selection/plot_nested_cross_validation_iris.html)

# K-Nearest Neighbours (KNN) Algorithm

## K-Nearest Neighbours – A Lazy Learning Algorithm

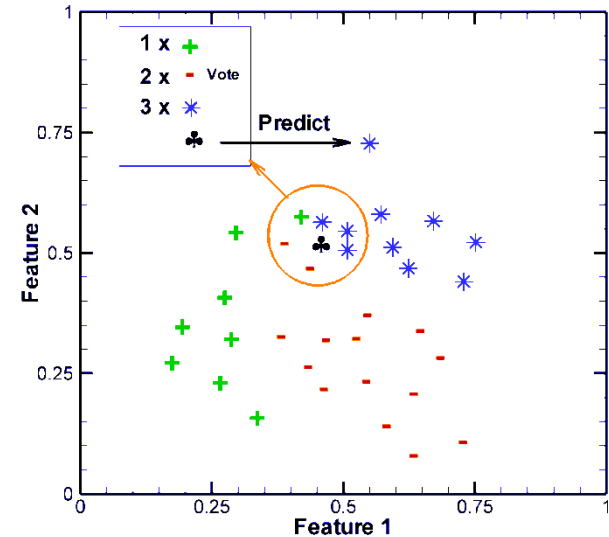
It is called lazy because it **doesn't learn** a discriminative function from the training data but **memorizes** the training dataset instead.

### Steps:

1. Choose the number of  $k$  and a distance metric.
2. Find the  $k$  nearest neighbours of the sample to be classified.
3. Assign the class label by majority vote.

The right choice of  $k$  is crucial to find a good balance between **over-** and **underfitting**.

$$D(X, Y) = \left( \sum_{i=1}^n |x_i - y_i|^p \right)^{1/p} \quad \text{Minkowski distance}$$

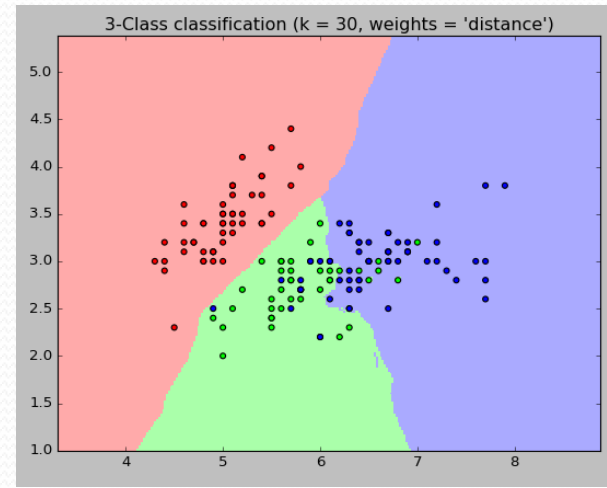
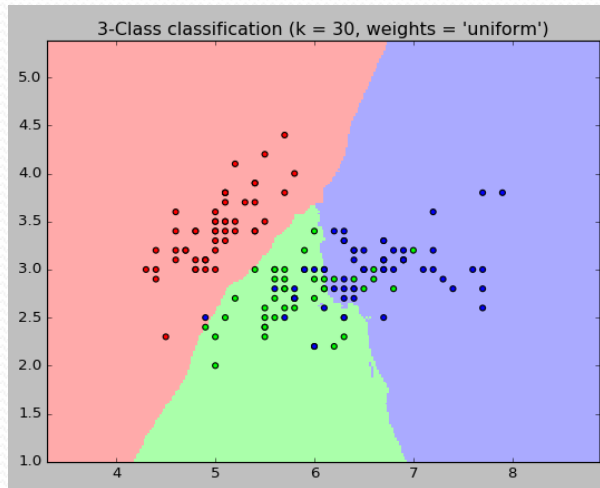
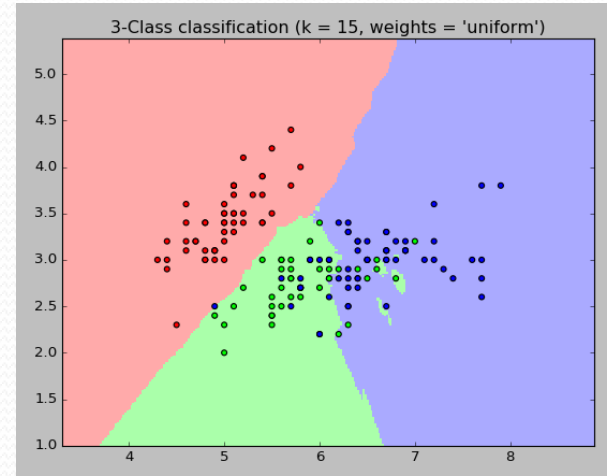
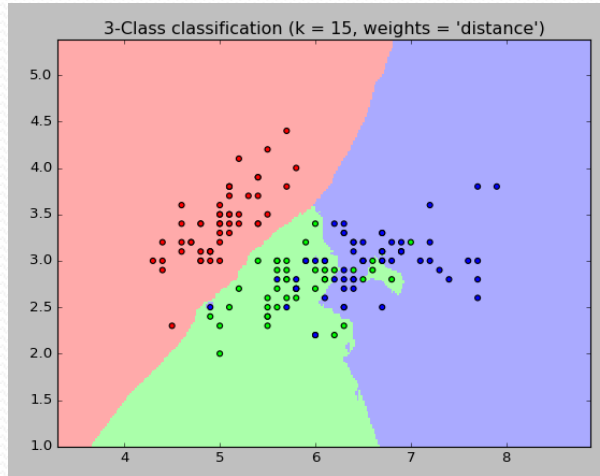


Metrics intended for real-valued vector spaces:

identifier	class name	args	distance function
“euclidean”	EuclideanDistance		$\text{sqrt}(\text{sum}((x - y)^2))$
“manhattan”	ManhattanDistance		$\text{sum}( x - y )$
“minkowski”	MinkowskiDistance	$p$	$\text{sum}( x - y ^p)^{1/p}$

# K-Nearest Neighbours (KNN) Algorithm

## Iris Example



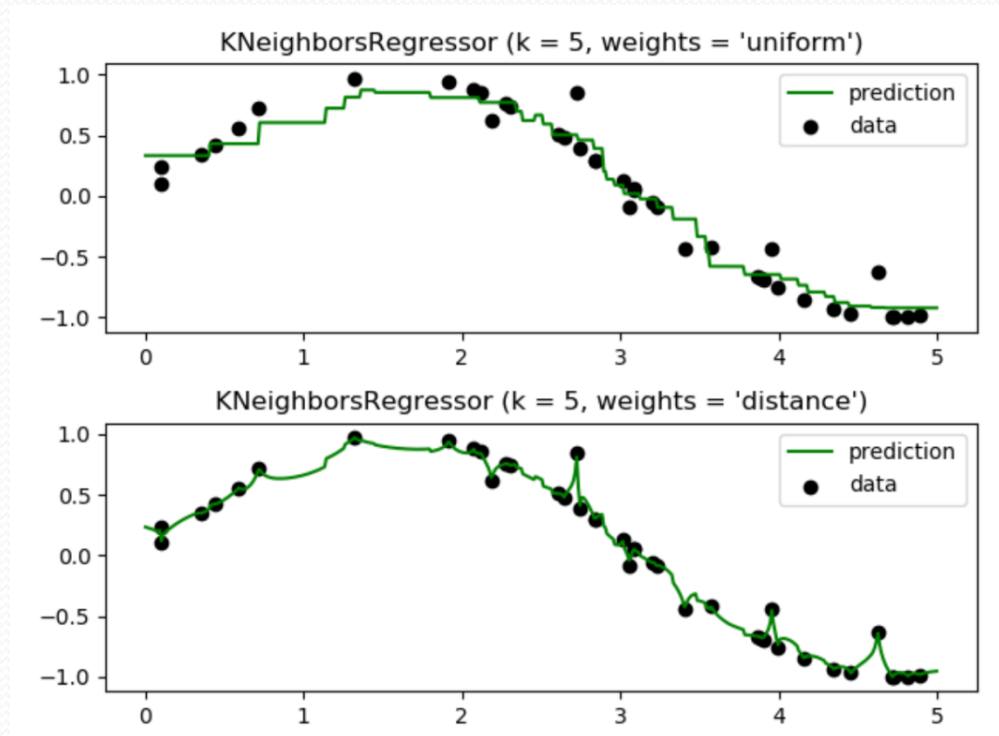
[https://scikit-learn.org/stable/auto\\_examples/neighbors/plot\\_classification.html#sphx-glr-auto-examples-neighbors-plot-classification-py](https://scikit-learn.org/stable/auto_examples/neighbors/plot_classification.html#sphx-glr-auto-examples-neighbors-plot-classification-py)

<https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html#sklearn.neighbors.KNeighborsClassifier>



# K-Nearest Neighbours (KNN) Algorithm

## Iris Example



### Sklearn:

```
clf = neighbors.KNeighborsClassifier(n_neighbors, weights=distance, p=2, metric='minkowski')  
clf.fit(X, y)
```