

## CS 1675 Intro to Machine Learning Lecture 8

### Logistic regression

Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

### Classification

- **Data:**  $D = \{d_1, d_2, \dots, d_n\}$   
 $d_i = \langle \mathbf{x}_i, y_i \rangle$ 
    - $y_i$  represents a discrete class value
  - **Goal: learn**  $f : X \rightarrow Y$
  - **Binary classification**
    - A special case when  $Y \in \{0,1\}$
  - **First step:**
    - we need to devise a model of the function f
-

## Discriminant functions

- A common way to represent a **classifier is by using**
  - **Discriminant functions**
- **Works for both the binary and multi-way classification**
- **Idea:**
  - For every class  $i = 0, 1, \dots, k$  define a function  $g_i(\mathbf{x})$  mapping  $X \rightarrow \Re$
  - When the decision on input  $\mathbf{x}$  should be made choose the class with the highest value of  $g_i(\mathbf{x})$

$$y^* = \arg \max_i g_i(\mathbf{x})$$

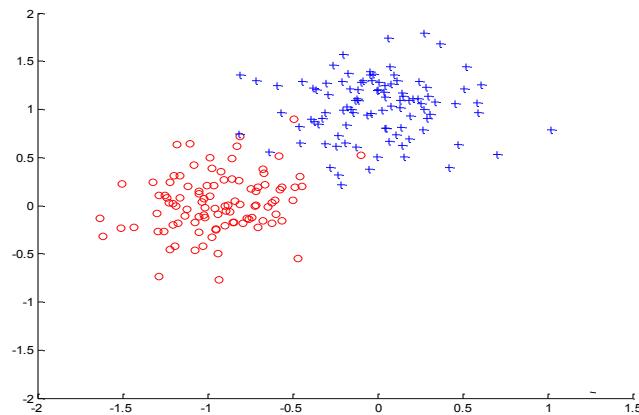
## Discriminant functions

- A common way to represent a **classifier is by using**
  - **Discriminant functions**
- **Works for both the binary and multi-way classification**
- **Idea:**
  - For every class  $i = 0, 1, \dots, k$  define a function  $g_i(\mathbf{x})$  mapping  $X \rightarrow \Re$
  - When the decision on input  $\mathbf{x}$  should be made choose the class with the highest value of  $g_i(\mathbf{x})$

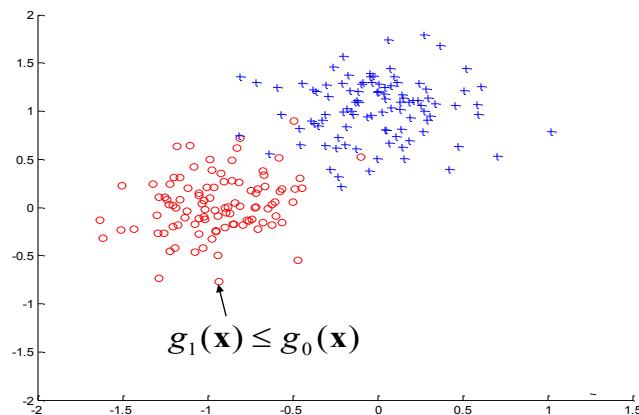
$$y^* = \arg \max_i g_i(\mathbf{x})$$

- So what happens with the input space? Assume a binary case.

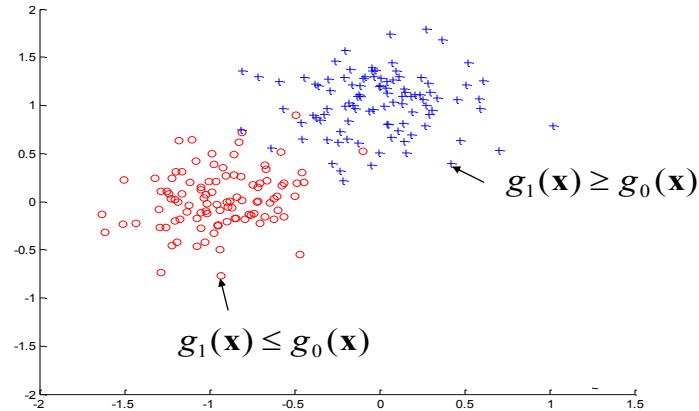
## Discriminant functions



## Discriminant functions

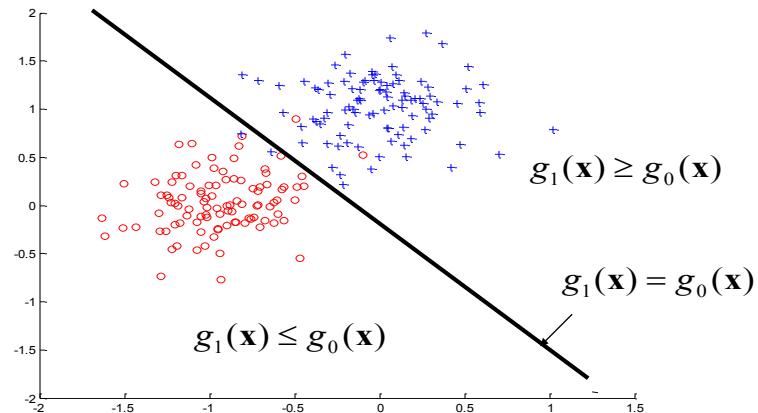


## Discriminant functions

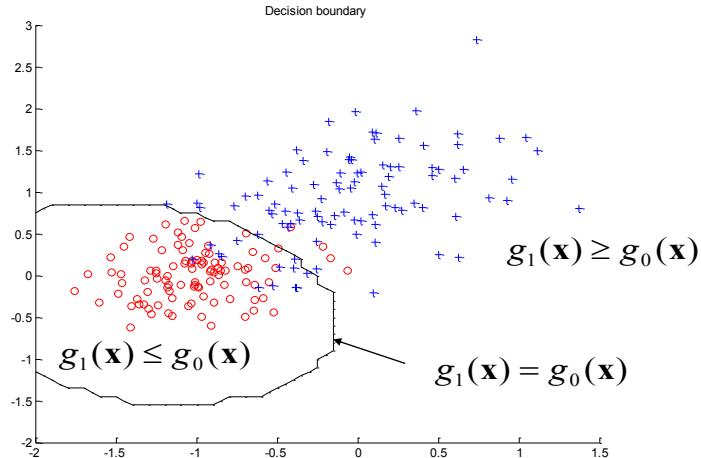


## Discriminant functions

- **Decision boundary:** discriminant functions are equal



## Quadratic decision boundary

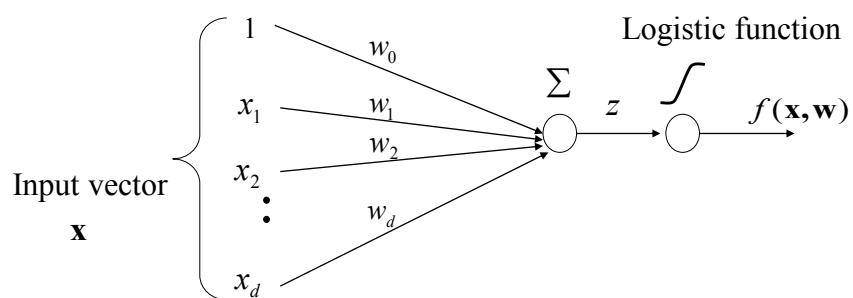


## Logistic regression model

- Defines a linear decision boundary
- Discriminant functions:  

$$g_1(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x}) \quad g_0(\mathbf{x}) = 1 - g(\mathbf{w}^T \mathbf{x})$$
- where  $g(z) = 1/(1 + e^{-z})$  - is a logistic function

$$f(\mathbf{x}, \mathbf{w}) = g_1(\mathbf{w}^T \mathbf{x}) = g(\mathbf{w}^T \mathbf{x})$$

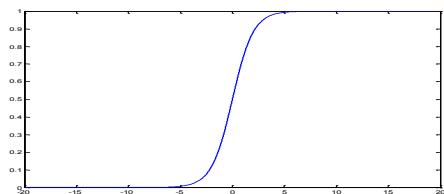


## Logistic function

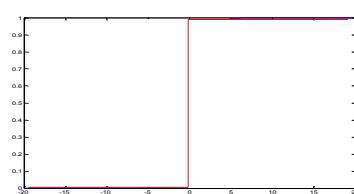
**Function:**

$$g(z) = \frac{1}{(1 + e^{-z})}$$

- Is also referred to as a **sigmoid function**
- takes a real number and outputs the number in the interval [0,1]
- Models a smooth switching function; replaces hard threshold function



Logistic (smooth) switching



Threshold (hard) switching

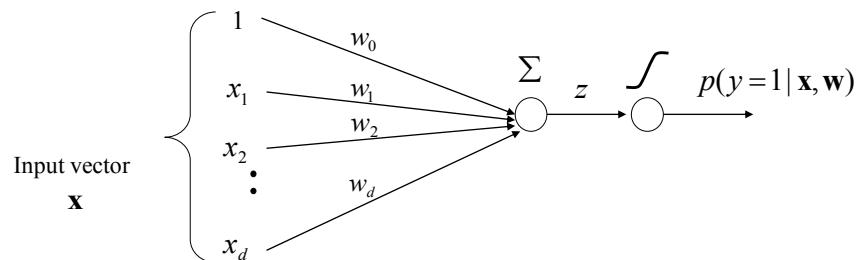
## Logistic regression model

- **Discriminant functions:**

$$g_1(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x}) \quad g_0(\mathbf{x}) = 1 - g(\mathbf{w}^T \mathbf{x})$$

- **Values of discriminant functions vary in interval [0,1]**
  - **Probabilistic interpretation**

$$f(\mathbf{x}, \mathbf{w}) = p(y=1 | \mathbf{w}, \mathbf{x}) = g_1(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x})$$



## Logistic regression

- We learn a **probabilistic function**

$$f : X \rightarrow [0,1]$$

– where  $f$  describes the probability of class 1 given  $\mathbf{x}$

$$f(\mathbf{x}, \mathbf{w}) = g_1(\mathbf{w}^T \mathbf{x}) = p(y=1 | \mathbf{x}, \mathbf{w})$$

**Note that:**

$$p(y=0 | \mathbf{x}, \mathbf{w}) = 1 - p(y=1 | \mathbf{x}, \mathbf{w})$$

- Making decisions with the logistic regression model:

?

## Logistic regression

- We learn a **probabilistic function**

$$f : X \rightarrow [0,1]$$

– where  $f$  describes the probability of class 1 given  $\mathbf{x}$

$$f(\mathbf{x}, \mathbf{w}) = g_1(\mathbf{w}^T \mathbf{x}) = p(y=1 | \mathbf{x}, \mathbf{w})$$

**Note that:**

$$p(y=0 | \mathbf{x}, \mathbf{w}) = 1 - p(y=1 | \mathbf{x}, \mathbf{w})$$

- Making decisions with the logistic regression model:

If  $p(y=1 | \mathbf{x}) \geq 1/2$  then choose **1**  
Else choose **0**

## Linear decision boundary

- Logistic regression model defines a **linear decision boundary**
- Why?
- **Answer:** Compare two **discriminant functions**.
- **Decision boundary:**  $g_1(\mathbf{x}) = g_0(\mathbf{x})$
- For the boundary it must hold:

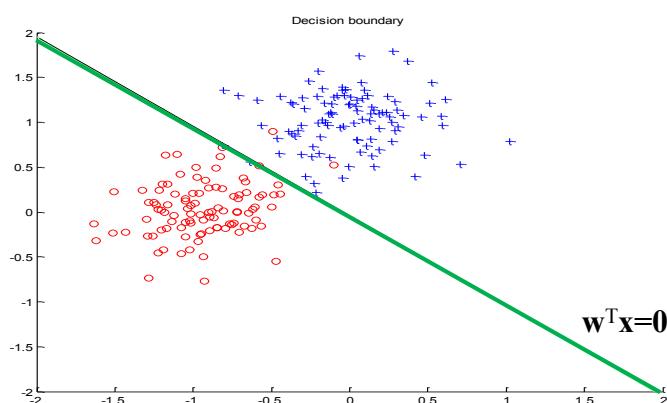
$$\log \frac{g_o(\mathbf{x})}{g_1(\mathbf{x})} = \log \frac{1 - g(\mathbf{w}^T \mathbf{x})}{g(\mathbf{w}^T \mathbf{x})} = 0$$

$$\log \frac{g_o(\mathbf{x})}{g_1(\mathbf{x})} = \log \frac{\frac{\exp - (\mathbf{w}^T \mathbf{x})}{1 + \exp - (\mathbf{w}^T \mathbf{x})}}{\frac{1}{1 + \exp - (\mathbf{w}^T \mathbf{x})}} = \log \exp - (\mathbf{w}^T \mathbf{x}) = \mathbf{w}^T \mathbf{x} = 0$$

## Logistic regression model. Decision boundary

- LR defines a linear decision boundary

Example: 2 classes (blue and red points)



## Logistic regression: parameter learning

### Likelihood of outputs

- Let  $D_i = \langle \mathbf{x}_i, y_i \rangle$   $\mu_i = p(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = g(z_i) = g(\mathbf{w}^T \mathbf{x})$

- Then

$$L(D, \mathbf{w}) = \prod_{i=1}^n P(y = y_i | \mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n \mu_i^{y_i} (1 - \mu_i)^{1-y_i}$$

- Find weights  $\mathbf{w}$  that maximize the likelihood of outputs

– Apply the log-likelihood trick. The optimal weights are the same for both the likelihood and the log-likelihood

$$\begin{aligned} l(D, \mathbf{w}) &= \log \prod_{i=1}^n \mu_i^{y_i} (1 - \mu_i)^{1-y_i} = \sum_{i=1}^n \log \mu_i^{y_i} (1 - \mu_i)^{1-y_i} = \\ &= \sum_{i=1}^n y_i \log \mu_i + (1 - y_i) \log (1 - \mu_i) \end{aligned}$$

## Logistic regression: parameter learning

- Notation:  $\mu_i = p(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = g(z_i) = g(\mathbf{w}^T \mathbf{x})$

- Log likelihood

$$l(D, \mathbf{w}) = \sum_{i=1}^n y_i \log \mu_i + (1 - y_i) \log (1 - \mu_i)$$

- Derivatives of the loglikelihood

$$\frac{\partial}{\partial w_j} l(D, \mathbf{w}) = \sum_{i=1}^n x_{i,j} (y_i - g(z_i))$$

Nonlinear in weights !!

$$\nabla_{\mathbf{w}} l(D, \mathbf{w}) = \sum_{i=1}^n \mathbf{x}_i (y_i - g(\mathbf{w}^T \mathbf{x}_i)) = \sum_{i=1}^n \mathbf{x}_i (y_i - f(\mathbf{w}, \mathbf{x}_i))$$

- Gradient descent:

$$\mathbf{w}^{(k)} \leftarrow \mathbf{w}^{(k-1)} - \alpha(k) \nabla_{\mathbf{w}} [-l(D, \mathbf{w})] \Big|_{\mathbf{w}^{(k-1)}}$$

$$\mathbf{w}^{(k)} \leftarrow \mathbf{w}^{(k-1)} + \alpha(k) \sum_{i=1}^n [y_i - f(\mathbf{w}^{(k-1)}, \mathbf{x}_i)] \mathbf{x}_i$$

## Derivation of the gradient

- **Log likelihood**  $l(D, \mathbf{w}) = \sum_{i=1}^n y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)$

- **Derivatives of the loglikelihood**

$$\frac{\partial}{\partial w_j} l(D, \mathbf{w}) = \sum_{i=1}^n \frac{\partial}{\partial z_i} [y_i \log g(z_i) + (1 - y_i) \log(1 - g(z_i))] \frac{\partial z_i}{\partial w_j}$$

**Derivative of a logistic function**

$$\frac{\partial z_i}{\partial w_j} = x_{i,j}$$

$$\frac{\partial g(z_i)}{\partial z_i} = g(z_i)(1 - g(z_i))$$

$$\begin{aligned} \frac{\partial}{\partial z_i} [y_i \log g(z_i) + (1 - y_i) \log(1 - g(z_i))] &= y_i \frac{1}{g(z_i)} \frac{\partial g(z_i)}{\partial z_i} + (1 - y_i) \frac{-1}{1 - g(z_i)} \frac{\partial g(z_i)}{\partial z_i} \\ &= y_i(1 - g(z_i)) + (1 - y_i)(-g(z_i)) \end{aligned}$$

$$\nabla_{\mathbf{w}} l(D, \mathbf{w}) = \sum_{i=1}^n -\mathbf{x}_i (y_i - g(\mathbf{w}^T \mathbf{x}_i)) = \sum_{i=1}^n -\mathbf{x}_i (y_i - f(\mathbf{w}, \mathbf{x}_i))$$

## Logistic regression. Online gradient descent

- **On-line component of the loglikelihood**

$$J_{\text{online}}(D_i, \mathbf{w}) = -[y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)]$$

- **On-line learning update for weight  $\mathbf{w}$**   $J_{\text{online}}(D_k, \mathbf{w})$

$$\mathbf{w}^{(k)} \leftarrow \mathbf{w}^{(k-1)} - \alpha(k) \nabla_{\mathbf{w}} [J_{\text{online}}(D_k, \mathbf{w})] |_{\mathbf{w}^{(k-1)}}$$

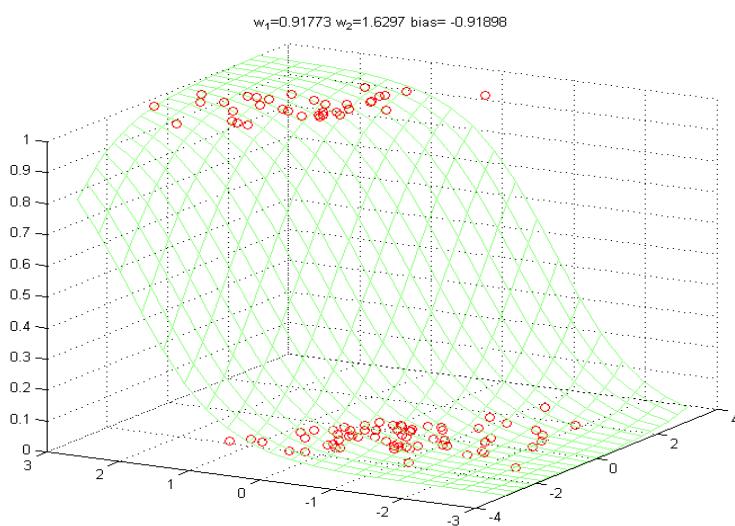
- **i-th update for the logistic regression** and  $D_k = \langle \mathbf{x}_k, y_k \rangle$

$$\mathbf{w}^{(k)} \leftarrow \mathbf{w}^{(k-1)} + \alpha(k) [y_k - f(\mathbf{w}^{(k-1)}, \mathbf{x}_k)] \mathbf{x}_k$$

## Online logistic regression algorithm

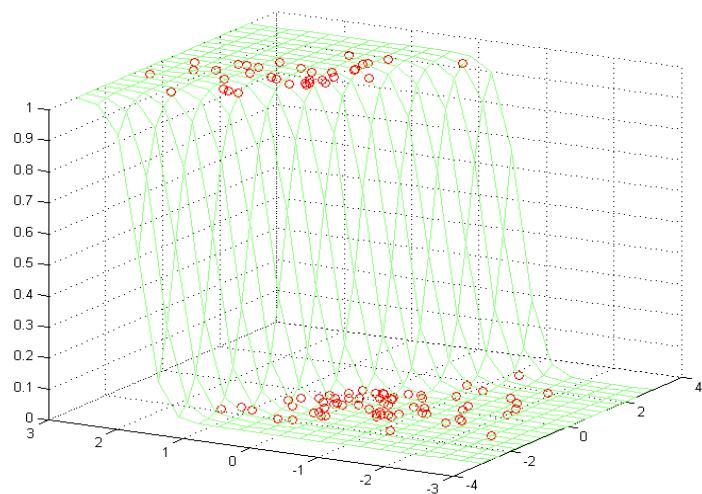
```
Online-logistic-regression (stopping_criterion)
    initialize weights    $\mathbf{w} = (w_0, w_1, w_2 \dots w_d)$ 
    while stopping_criterion = FALSE
        do      select next data point  $D_i = \langle \mathbf{x}_i, y_i \rangle$ 
            set       $\alpha(i)$ 
            update weights (in parallel)
                 $\mathbf{w} \leftarrow \mathbf{w} + \alpha(i)[y_i - f(\mathbf{w}, \mathbf{x}_i)]\mathbf{x}_i$ 
        end
    return weights  $\mathbf{w}$ 
```

## Online algorithm. Example.



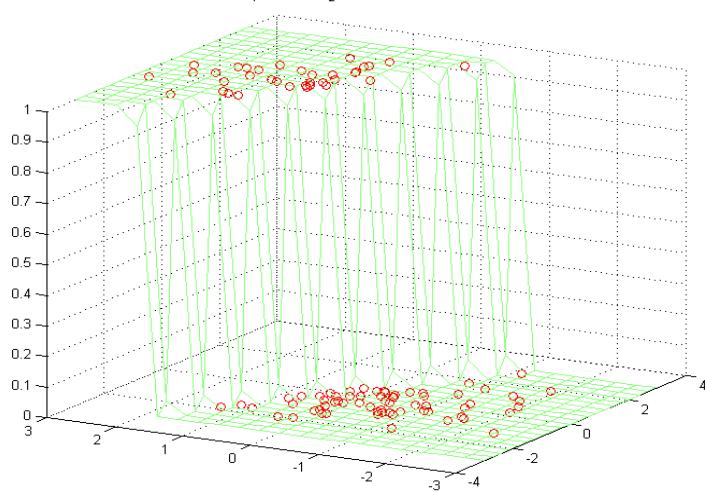
## Online algorithm. Example.

$w_1=3.5934$   $w_2=6.9126$  bias= -3.6709



## Online algorithm. Example.

$w_1=19.9144$   $w_2=39.7033$  bias= -20.8644



## Evaluation of classifiers

## Evaluation

For any data set we use to test the classification model on we can build a **confusion matrix**:

- Counts of examples with:
- class label  $\omega_j$  that are classified with a label  $\alpha_i$

		target	
		$\omega = 1$	$\omega = 0$
predict	$\alpha = 1$	140	17
	$\alpha = 0$	20	54

## Evaluation

For any data set we use to test the classification model on we can build a **confusion matrix**:

- Counts of examples with:
- class label  $\omega_j$  that are classified with a label  $\alpha_i$

		target	
		$\omega = 1$	$\omega = 0$
predict	$\alpha = 1$	140	17
	$\alpha = 0$	20	54

CS 2750 Machine Learning

## Evaluation

For any data set we use to test the model we can build a **confusion matrix**:

		target	
		$\omega = 1$	$\omega = 0$
predict	$\alpha = 1$	140	17
	$\alpha = 0$	20	54

agreement

Accuracy: ?

CS 2750 Machine Learning

## Evaluation

For any data set we use to test the model we can build a confusion matrix:

		target	
		$\omega = 1$	$\omega = 0$
		predict	
predict	$\alpha = 1$	140	17
	$\alpha = 0$	20	54

agreement

**Accuracy** = 194/231

CS 2750 Machine Learning

## Evaluation

For any data set we use to test the model we can build a confusion matrix:

		target	
		$\omega = 1$	$\omega = 0$
		predict	
predict	$\alpha = 1$	140	17
	$\alpha = 0$	20	54

**Accuracy** = 194/231

**Error** = ?

CS 2750 Machine Learning

## Evaluation

For any data set we use to test the model we can build a confusion matrix:

		target	
		$\omega = 1$	$\omega = 0$
		$\alpha = 1$	$\alpha = 0$
$\alpha = 1$	140	17	
$\alpha = 0$	20	54	

$$\text{Accuracy} = 194/231$$

$$\text{Error} = 37/231 = 1 - \text{Accuracy}$$

CS 2750 Machine Learning

## Evaluation for binary classification

Entries in the confusion matrix for binary classification have names:

		target	
		$\omega = 1$	$\omega = 0$
		$\alpha = 1$	$\alpha = 0$
$\alpha = 1$	$TP$	$FP$	
$\alpha = 0$	$FN$	$TN$	

$TP$ : True positive (hit)

$FP$ : False positive (false alarm)

$TN$ : True negative (correct rejection)

$FN$ : False negative (a miss)

## Additional statistics

- Sensitivity (recall)

$$SENS = \frac{TP}{TP + FN}$$

		target	
		$\omega = 1$	$\omega = 0$
predict	$\alpha = 1$	TP	FP
	$\alpha = 0$	FN	TN

## Additional statistics

- Sensitivity (recall)

$$SENS = \frac{TP}{TP + FN}$$

- Specificity

$$SPEC = \frac{TN}{TN + FP}$$

- Positive predictive value (precision)

$$PPT = \frac{TP}{TP + FP}$$

- Negative predictive value

$$NPV = \frac{TN}{TN + FN}$$

## Binary classification: additional statistics

- Confusion matrix

		target		
		1	0	
		predict		
1		140	10	$PPV = 140/150$
0		20	180	$NPV = 180/200$
		$SENS = 140/160$	$SPEC = 180/190$	

### Row and column quantities:

- Sensitivity (SENS)
- Specificity (SPEC)
- Positive predictive value (PPV)
- Negative predictive value (NPV)