Lecture 4

# Logistic Regression and Classification

**Physical Address for Prof. Ilias Tagkopoulos**
Computer Science:
Office: 3063 Kemper Hall
Phone: (530) 752-4821
Fax: (530) 752-4767

Genome and Biomedical Sciences Facility:
Office: 5313 GBSF
Phone: (530) 752-7707
Fax: (530) 754-9658

Instructor: Ilias Tagkopoulos
iliast@ucdavis.edu

- You were given a dataset with *m* samples $D = \{(x^{(i)}, y^{(i)}); i = 1 \dots m\}$.

  − Note that the superscript $x^{(i)}$ is the index of the sample.

  − Assume that each sample has *n* attributes (features)

$$D = \left\{ \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_n^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_n^{(2)} \\ \cdots & \cdots & \cdots & \cdots \\ 1 & x_1^{(m)} & \cdots & x_n^{(m)} \end{bmatrix} \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \cdots \\ y^{(m)} \end{bmatrix} \right\} \text{ or}$$

$$D = \left\{ \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \cdots \\ (x^{(m)})^T \end{bmatrix} \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \cdots \\ y^{(m)} \end{bmatrix} \right\} \text{ with } x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \cdots \\ x_n^{(i)} \end{bmatrix} \text{ being the } i^{th} \text{ sample}$$

$$D = \{X \quad Y\}$$

- By using the dataset D, you can **build a function that relates the input (x) to the output (y)**.

- You can then use this function to **predict** what the output (y) will be, based on a new, possibly never seen before, input.

$$f(x; w): X \rightarrow Y$$

With **weight vector** $w = \begin{bmatrix} w_0 \\ \cdots \\ w_n \end{bmatrix}$ and **input vector** $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \cdots \\ x_n^{(i)} \end{bmatrix}$

- The task is to *find the optimal parameter values* for this function, i.e. the $w$, so that the function $f(x; w)$ *best describes the relationship between X and Y*

- **Method 1: Ordinary Least Squares**

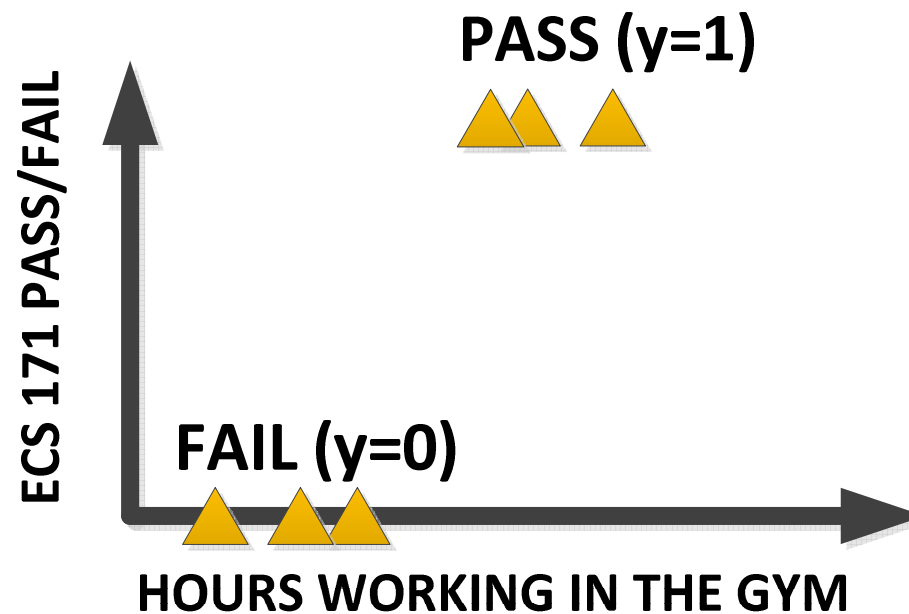$$\frac{\partial RSS}{\partial w} = 0 \Rightarrow$$
$$w = (X^T X)^{-1} X^T Y$$

- **Method 2: Gradient Descent**

$$w_j := w_j + a\left(y^{(i)} - \sum_{k=0}^{n} w_k x_k^{(i)}\right) x_j^{(i)}$$

**Next** $w_j$    **Previous** $w_j$
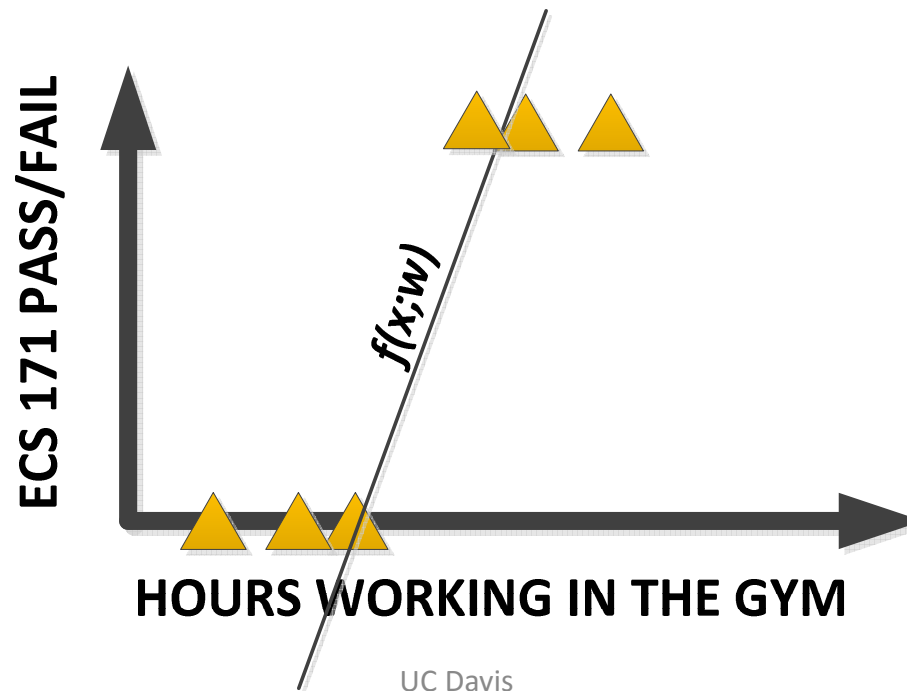
**constant**

**Update proportional to error**

## Logistic Regression: a classification method

- Suppose that all your aim is not to find the value of a continuous variable *per se*, but to categorize the samples into buckets or classes.

- Sure, you can still perform linear regression and then threshold it, but the solution will be sensitive to outliers and to the selected threshold.
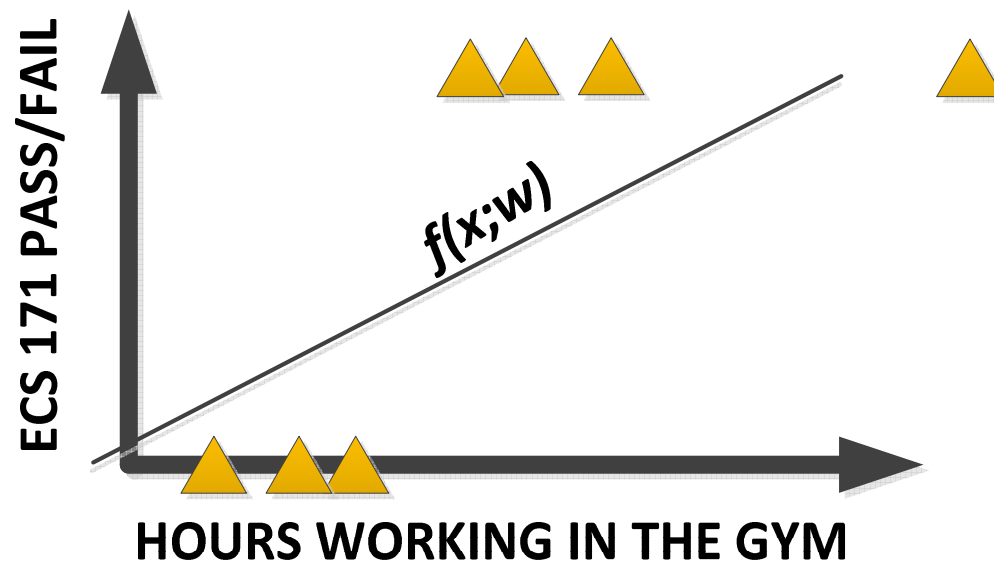
## Logistic Regression: a classification method

- Suppose that all your aim is not to find the value of a continuous variable *per se*, but to categorize the samples into buckets or classes.

- Sure, you can still perform linear regression and then threshold it, but the solution will be sensitive to outliers and to the selected threshold.

# Logistic Regression: a classification method
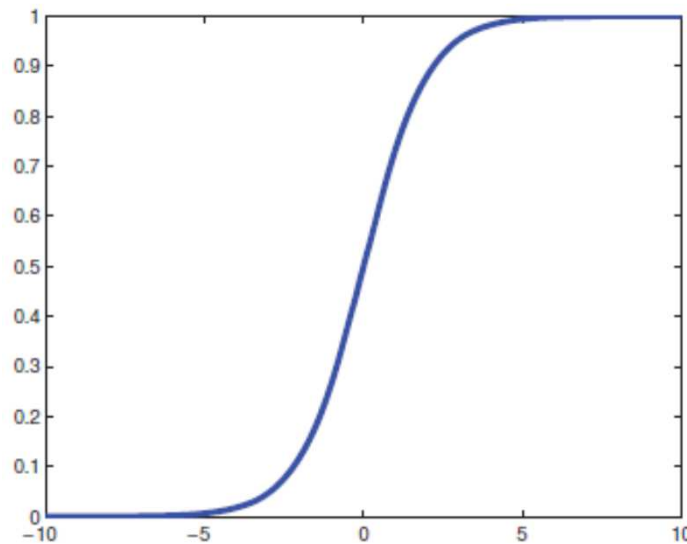
- Suppose that all your aim is not to find the value of a continuous variable *per se*, but to categorize the samples into buckets or classes.

- Sure, you can still perform linear regression and then threshold it, but the solution will be sensitive to outliers and to the selected threshold.

# Logistic Regression: a classification method

- A better way to perform classification when within the regression framework is logistic regression

  - Same formulation as with linear regression, but instead of a polynomial, use the sigmoid/logit/logistic function:

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

## Logistic Regression: a classification method
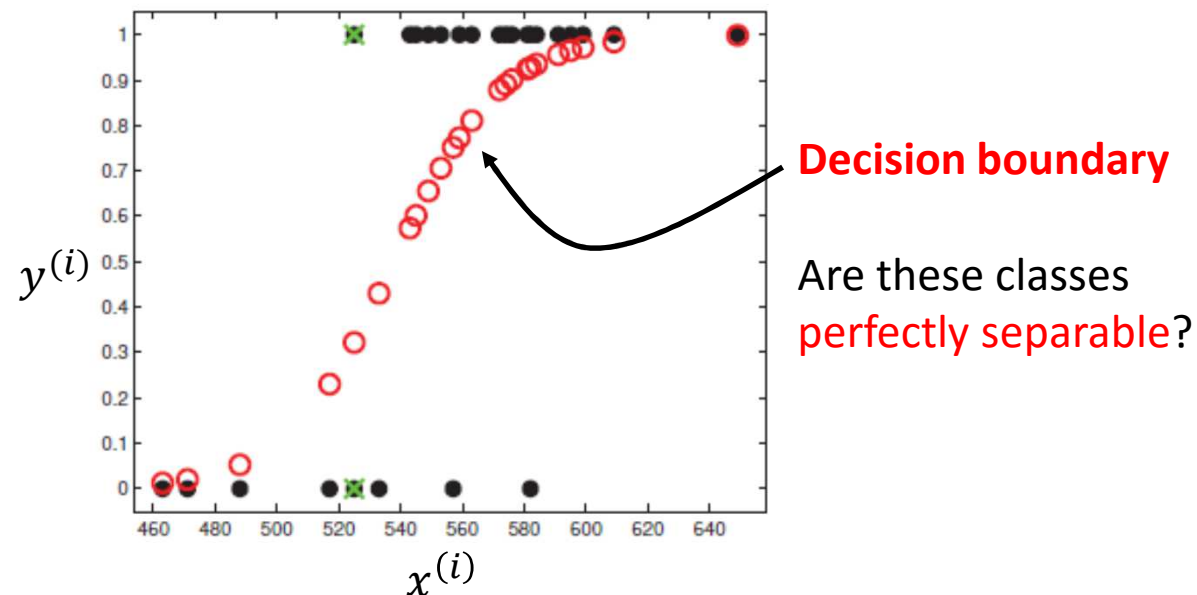
- As such, we can create the function

$$g(x; w) = sigm(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

And categorize into two classes (positive or negative) for any given sample $i$, by defining the label $y^{(i)}$ to be:

$$y^{(i)} = \begin{cases} 0 & if & g(w^T x^{(i)}) < threshold \\ 1 & if & g(w^T x^{(i)}) \geq threshold \end{cases}$$

With logistic regression, we can express it as probabilities too (instead of hard boundaries)



**Decision boundary**

Are these classes perfectly separable?

## Logistic Regression: MLE

- Great, but how do we find the optimal parameter/weight **w** set?

- There is no closed form solution (such as the OLS for linear regression) but we can first formulate it as a MLE problem and then use gradient descent to find the parameters.

- To do that, lets express the probability of each class:

$$P\big(y^{(i)} = 1 \,\big|\, x^{(i)}; \boldsymbol{w}\big) = g(x^{(i)}; w)$$
$$P\big(y^{(i)} = 0 \,\big|\, x^{(i)}; \boldsymbol{w}\big) = 1 - g(x^{(i)}; w)$$

Which can also be written as:

$$p\big(y^{(i)} \big| x^{(i)}; \boldsymbol{w}\big) = g(x^{(i)}; w)^{y^{(i)}} (1 - g(x^{(i)}; w))^{1 - y^{(i)}}$$

- Assuming the samples are i.i.d. then we can write the log Likelihood as :

$$l(\text{w}) \triangleq log \, p(D|w) = \sum_{i=1}^{M} \log p(y^{(i)}|x^{(i)}; \text{w}) =$$

$$= \sum_{i=1}^{M} \log\left(g(x^{(i)}; w)^{y^{(i)}}(1 - g(x^{(i)}; w))^{1-y^{(i)}}\right) =$$

$$= \sum_{i=1}^{M} y^{(i)} \log g\left(x^{(i)}; w\right) + (1 - y^{(i)})\log(1 - g(x^{(i)}; w))$$

■ To maximize the log likelihood, we first find <span style="color:red">the derivative of the log likelihood with respect to w:</span>

$$\frac{\partial l(\text{w})}{\partial w}$$

Which yields

$$\frac{\partial l(w)}{\partial w_j} = \left(y^{(i)} - g(x^{(i)}; w)\right) x_j^{(i)}$$

- With that, we can now apply gradient descent (ascent), which as we know updates the w based on the following rule:

$$w_j := w_j + a \frac{\partial l(w)}{\partial w_j}$$

Or

$$w_j := w_j + a\left(y^{(i)} - g(x^{(i)}; w)\right)x_j^{(i)}$$

Have you seen this before?

- Btw, gradient descent is not the only method to find the parameter w. E.g. Newton's method

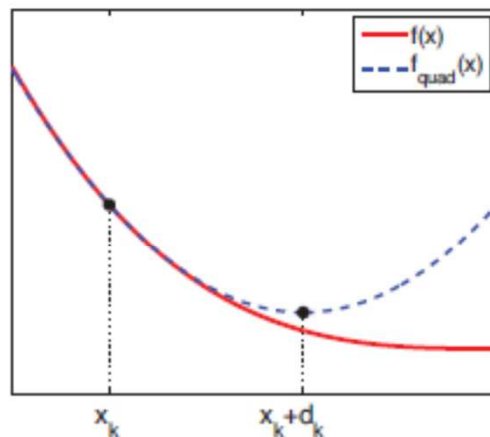$$w_j := w_j - \frac{\frac{\partial l(w)}{\partial w_j}}{\frac{\partial^2 l(w)}{(\partial w_j)^2}}$$

# Logistic Regression: Newton-Raphson

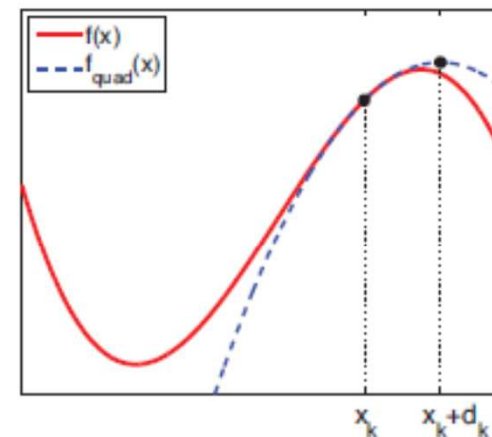**Algorithm 8.1:** Newton's method for minimizing a strictly convex function

1  Initialize $\theta_0$;
2  **for** $k = 1, 2, \ldots$ *until convergence* **do**
3      Evaluate $\mathbf{g}_k = \nabla f(\theta_k)$;
4      Evaluate $\mathbf{H}_k = \nabla^2 f(\theta_k)$;
5      Solve $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$ for $\mathbf{d}_k$;
6      Use line search to find stepsize $\eta_k$ along $\mathbf{d}_k$;
7      $\theta_{k+1} = \theta_k + \eta_k \mathbf{d}_k$;

$$w_j := w_j - \frac{\frac{\partial l(w)}{\partial w_j}}{\frac{\partial^2 l(w)}{(\partial w_j)^2}}$$

$$\theta_{k+1} = \theta_k - \eta_k \mathbf{H}_k^{-1} \mathbf{g}_k$$



(a)

(b)

- Actually if we use the same update rule but with hard boundaries, forcing the output to be {0,1}, we have the <span style="color:red">perceptron learning algorithm</span>

$$w_j := w_j + a\left(y^{(i)} - g(x^{(i)}; w)\right)x_j^{(i)}$$

with

$$g(z) = \begin{cases} 0 & if & z < threshold \\ 1 & if & z \geq threshold \end{cases}$$

Threshold can be any <span style="color:red">scalar</span> (e.g. 0).

- Generally <span style="color:red">Stochastic Gradient Descent</span> on <span style="color:red">logistic regression</span> is faster and has a better performance.

**End of Lecture 4**