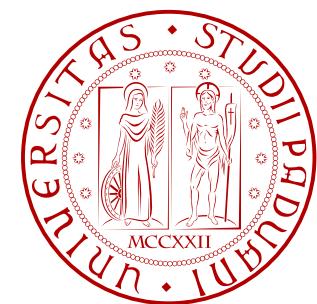


# A PRIMER ON REGRESSION AND CLASSIFICATION MODELS

---

Michele Rossi  
[rossi@dei.unipd.it](mailto:rossi@dei.unipd.it)

Dept. of Information Engineering  
University of Padova, IT



# Outline

- Regression
  - Linear space
  - Non-linear (feature) space
  - Logistic regression problem
  - Regularization
- Classification, class conditional probabilities
  - 2 classes
  - $K > 2$  classes

# REGRESSION

---

# Regression

- Fitting a given set of N data points

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$$

- Outputs are (here) scalar
- Input points are (column) vectors of D elements

$$\mathbf{x}_i = [x_{1,i} \ x_{2,i} \ \dots \ x_{D,i}]^T$$

- Linear regression models
  - Finding a best *fitting line* in 2D ( $D=1$ ), a best fitting *plane* in 3D ( $D=2$ ) or a *hyper-plane* in  $D+1$  dimensions when input vectors have D elements, with  $D>2$

# Linear regression objective

- We introduce a (column) vector of weights

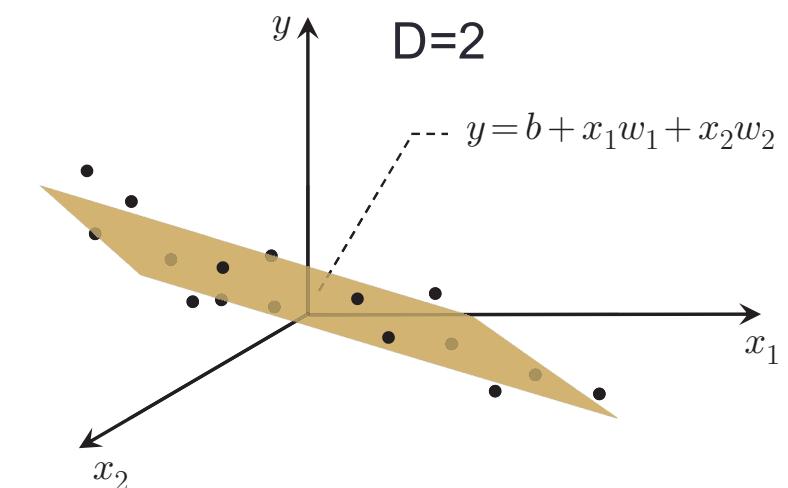
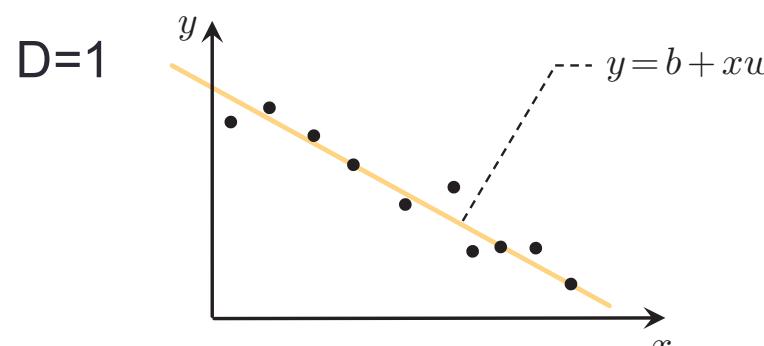
$$\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]^T$$

- Linear regression

- amount to finding the *best fitting hyper-plane*

$$b + \mathbf{x}_i^T \mathbf{w} \approx y_i, \quad i = 1, 2, \dots, N$$

- Input  $\mathbf{x}_i$  are referred to as **input features**
  - b is a constant, referred to as **the bias**



# Common cost function

- To find the best fitting hyper-plane
  - it is common to use the Least squares cost function

$$g(b, \mathbf{w}) = \sum_{i=1}^N (b + \mathbf{x}_i^T \mathbf{w} - y_i)^2$$

- Minimize sum of squared differences between
  - Input features and output points
  - Across the entire dataset
- Objective:

$$\underset{b, \mathbf{w}}{\text{minimize}} \sum_{i=1}^N (b + \mathbf{x}_i^T \mathbf{w} - y_i)^2$$

# Compact notation

$$\tilde{\mathbf{x}}_i = \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix} \quad \tilde{\mathbf{w}} = \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}$$

- This notation trick is used all the times with ANNs
- With this notation, we can rewrite the cost function as:

$$g(\tilde{\mathbf{w}}) = \sum_{i=1}^N (\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - y_i)^2$$

- This function is (quadratic) *always convex*
  - Regardless of the input data points

# The gradient

- Applying the **chain rule for derivatives**, we get:

$$\nabla g(\tilde{\mathbf{w}}) = 2 \sum_{i=1}^N \tilde{\mathbf{x}}_i (\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - y_i) = 2 \left( \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right) \tilde{\mathbf{w}} - 2 \sum_{i=1}^N \tilde{\mathbf{x}}_i y_i$$

- We can perform gradient descent to minimize the cost
- However, in this (rare) instance
  - We can set the gradient equal to zero and solve for  $\tilde{\mathbf{w}}$
  - This leads to the following **system of equations**:

$$\left( \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right) \tilde{\mathbf{w}} = \sum_{i=1}^N \tilde{\mathbf{x}}_i y_i$$

# Direct solution

- If matrix  $\sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T$  is invertible, we obtain:

$$\tilde{\mathbf{w}} = \left( \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right)^{-1} \sum_{i=1}^N \tilde{\mathbf{x}}_i y_i$$

- This is, in practice, **more computationally demanding** than solving the system of equations

# The efficacy of a learned model

- With the optimal parameters

$$\tilde{\mathbf{w}}^* = \begin{bmatrix} b^* \\ \mathbf{w}^* \end{bmatrix}$$

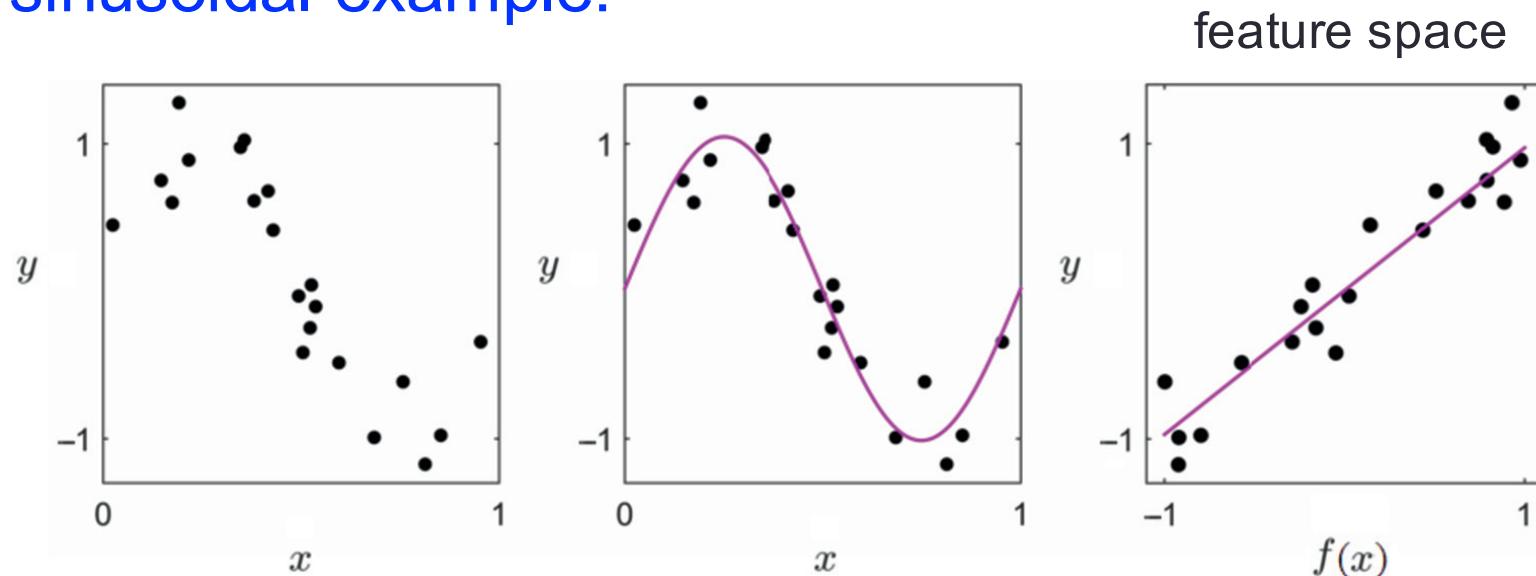
- We can compute the *efficacy* of a learned model
- Using the Mean Square Error (MSE)

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (b^* + \mathbf{x}_i^T \mathbf{w}^* - y_i)^2$$

- Sometimes it is a good idea to compute the MSE with *new testing data* (not used for training): to provide assurance that the model works well on future data points

# Non linear data (1/2)

- Often, data **do not follow a linear regression model**
- Solution: represent the data in a **suitable feature space**
- (nosy) sinusoidal example:



(left panel) A simulated regression dataset. (middle panel) A weighted form of a simple sinusoid  $y = b + f(x)$  (in magenta), where  $f(x) = \sin(2\pi x)$  and where  $b$  and  $w$  are tuned properly, describes the data quite well. (right panel) Fitting a sinusoid in the original feature space is equivalent to fitting a line in the transformed feature space where the input feature has undergone feature transformation  $x \rightarrow f(x) = \sin(2\pi x)$ .

## Non linear data (2/2)

- This amounts to using a non-linear transformation

$$\mathbf{f}_i \triangleq f(\mathbf{x}_i) \quad \tilde{\mathbf{f}}_i = \begin{bmatrix} 1 \\ \mathbf{f}_i \end{bmatrix}$$

- Use this for the linear fitting in place of the original data
- The **new cost function** is:

$$g(\tilde{\mathbf{w}}) = \sum_{i=1}^N (\tilde{\mathbf{f}}_i^T \tilde{\mathbf{w}} - y_i)^2$$

- Is still *quadratic* and *convex*

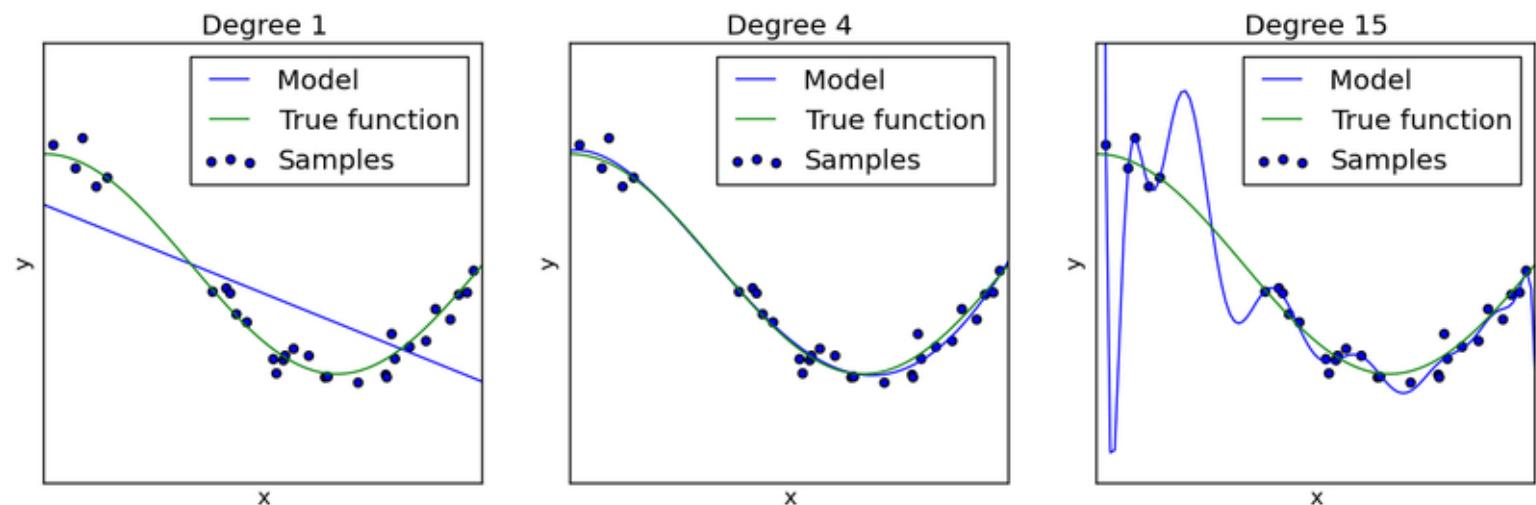
# Take away message

A properly designed feature (or set of features) for linear regression provides a good *nonlinear* fit in the original feature space and, simultaneously, a good *linear* fit in the transformed feature space.

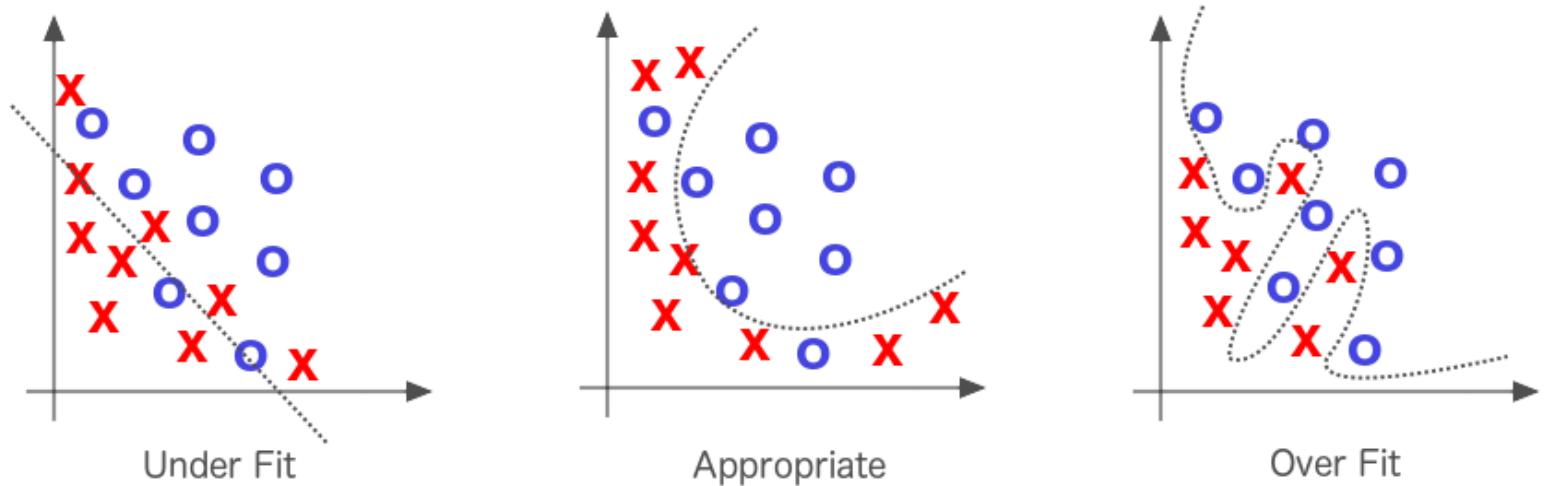
- Meaning of “good” deserves some attention
  - We do not want to “overfit” the data
  - **Overfitting**
    - production of a model that corresponds *closely* or *exactly* to a particular set of data, and may therefore fail to fit additional (new) data or *predict future observations reliably*
    - **essence of overfitting** - having extracted some of the residual variation (representing noise), as if that variation represented *underlying model structure*. This produces models with *too many parameters*

# Overfitting: visual examples

Regression:



Classification:



# Non-linear regression

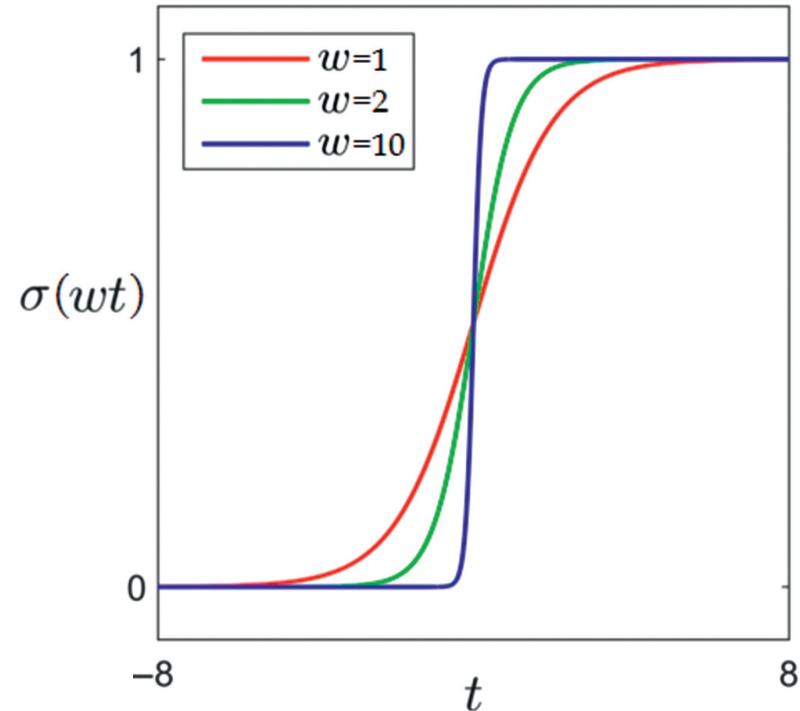
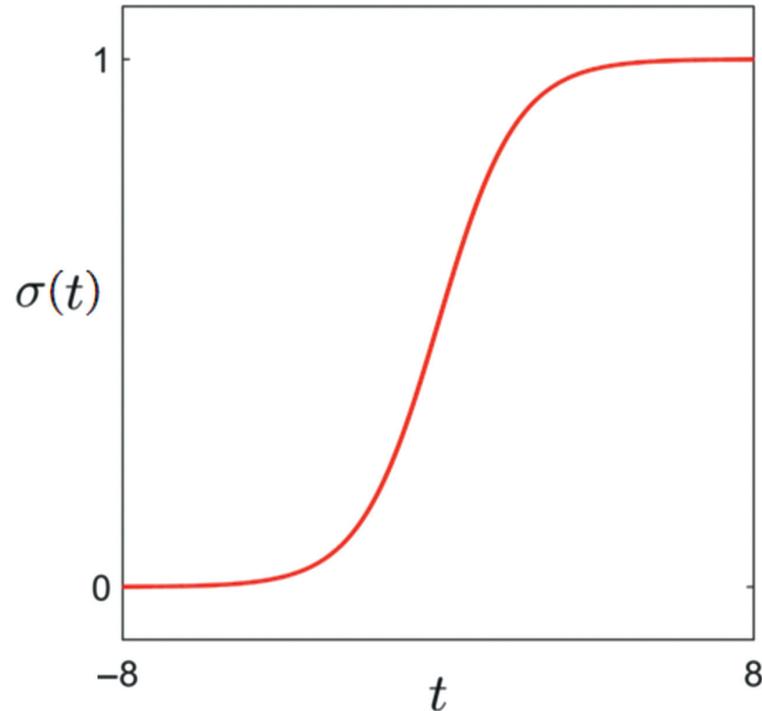
- Previous model
  - Nonlinear relation was easy to assess
  - Cost function remained convex and
  - Input ( $x$ ) / output ( $y$ ) relationship was linear the parameters ( $b$ ,  $\mathbf{w}$ ):
    - *It still is a linear regression problem*
- Next, we explore
  - Nonlinear models for regression where
    - Cost function is no longer convex
    - Input ( $x$ ) / output ( $y$ ) relationship nonlinear in its parameters
    - $\ell_2$  regularization
- To be concrete: we focus on *logistic regression*, but the discussion applies to a broader set of problems

$$b + f(\mathbf{x}_i)^T \mathbf{w} \approx y_i$$

  
linear

# Logistic sigmoid function

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$



- Devised by Pierre François Verhulst in the early 19<sup>th</sup> century
  - To model how a population (bacteria, animal species, etc.) grows in time
  - **Intuition:** regardless of the environment, the system only has a finite amount of resources (space, food, etc.) → there must be a cap on the total population in any biological system → saturation point

# Logistic regression problem (1/2)

- If a dataset is distributed according to the **logistic function**
- Then, this data must satisfy:

$$\sigma(b + \mathbf{x}_i^T \mathbf{w}) \approx y_i , \quad i = 1, 2, \dots, N$$

- So we have a nonlinear system of P equations
- Non-linearity leads to a **non-convex Least Squares cost function**

$$g(b, \mathbf{w}) = \sum_{i=1}^N (\sigma(b + \mathbf{x}_i^T \mathbf{w}) - y_i)^2$$

$$g(\tilde{\mathbf{w}}) = \sum_{i=1}^N (\sigma(\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}}) - y_i)^2 \quad \text{compact notation}$$

# Logistic regression problem (2/2)

- Note that (**key property**):

$$\sigma'(t) \triangleq \frac{d\sigma(t)}{dt} = -\frac{(1 + e^{-t})'}{(1 + e^{-t})^2} = \frac{e^{-t}}{(1 + e^{-t})^2} = \sigma(t)(1 - \sigma(t))$$

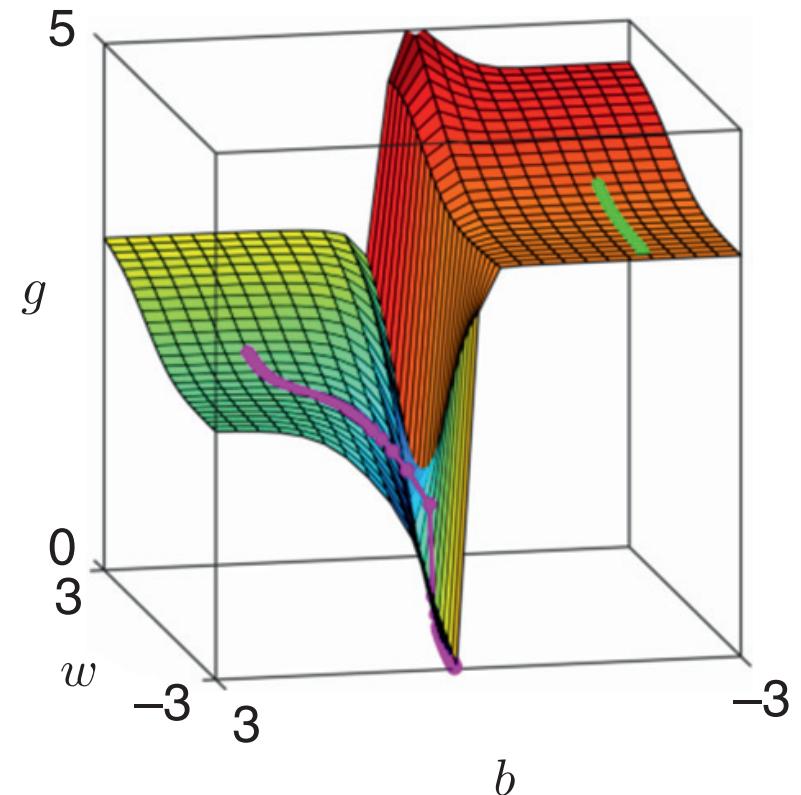
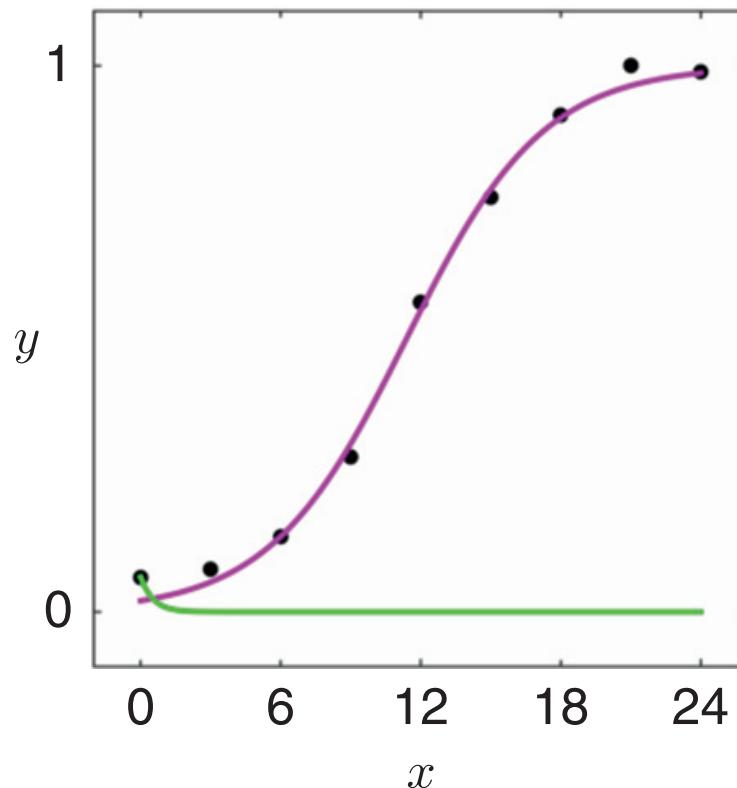
- Using this and the chain rule for derivatives,

$$\nabla g(\tilde{\mathbf{w}}) = \sum_{i=1}^N 2 \left( \sigma(\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}}) - y_i \right) \sigma(\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}})(1 - \sigma(\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}})) \tilde{\mathbf{x}}_i$$

- Due to the involved non-linearities

- Finding directly a minimum of the cost function *is a fruitless venture*
- A numerical technique (e.g., Newton's method) may be used

# Finding the minimum ...



- Two fits found minimizing the least square cost
- Using **gradient descent** with a different initialization
- The cost function is *no longer convex*

# Regularization

- The problematic flat areas *of non-convex functions*
  - Can be ameliorated using a regularizer
  - Is a simple convex function that is added to the original cost function
  - Has the effect of *slightly convexifying the cost function*
  - The most common regularizer is the squared  $\ell_2$  norm

$$\|\mathbf{w}\|^2 = \sum_{i=1}^D w_i^2$$

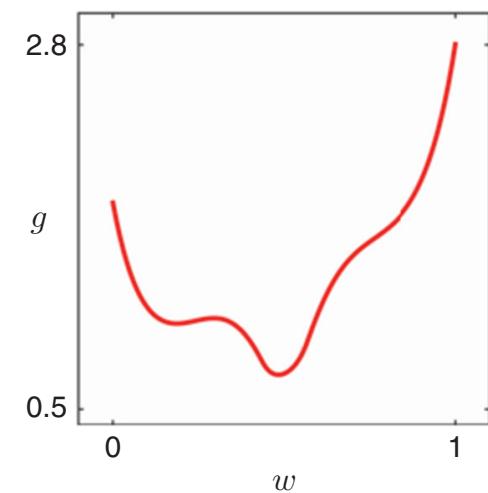
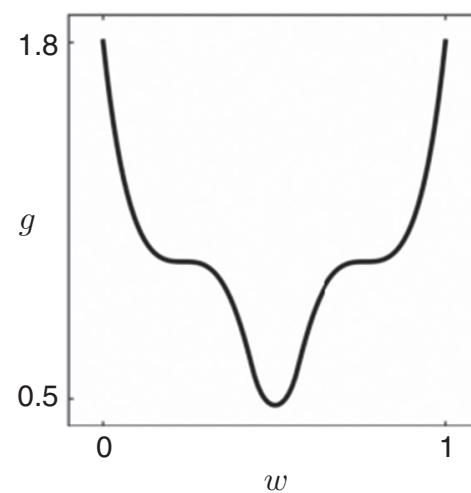
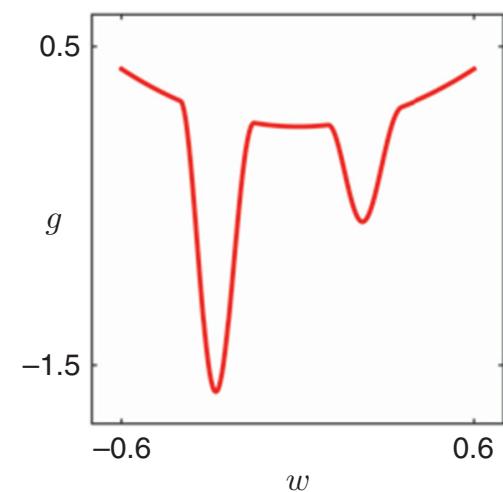
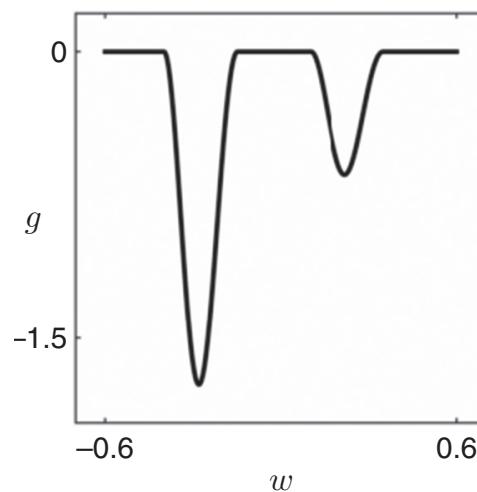
- The *regularized cost function* is:

$$g(b, \mathbf{w}) + \lambda \|\mathbf{w}\|^2$$

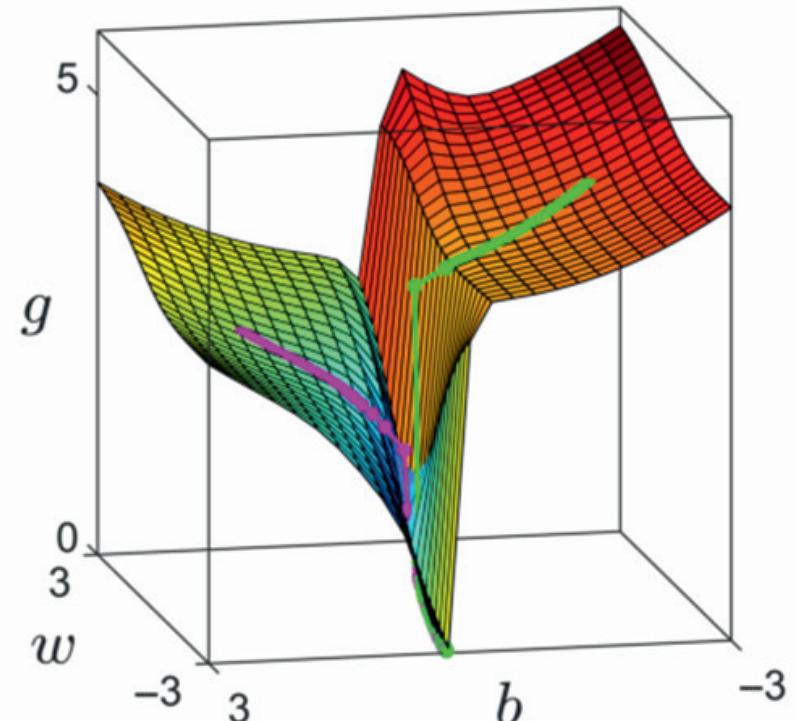
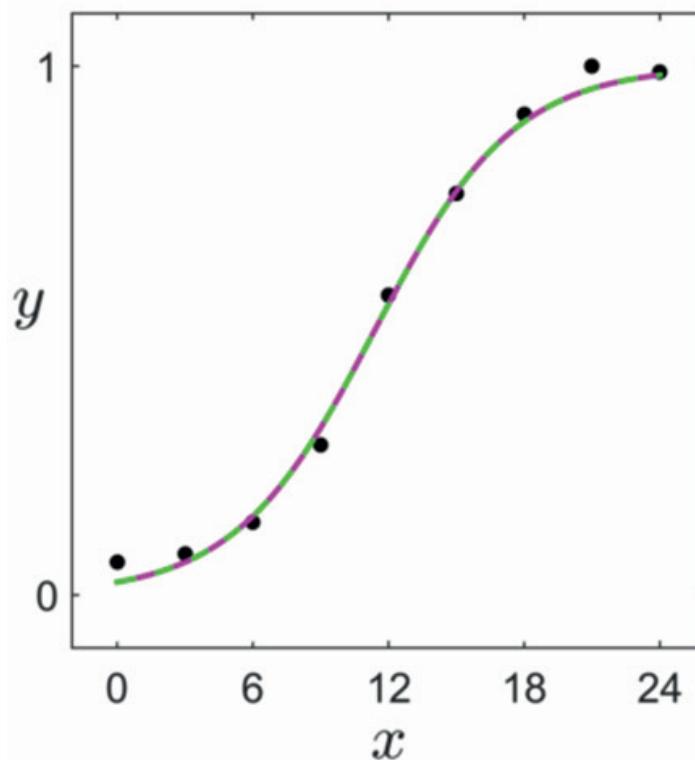
- $\lambda$  is set by the user and weighs the importance of the two terms

# Regularization - examples

- $\lambda=1$  (pretty large)
- The regularizer helps
  - remove flat areas
  - remove saddle points
  - remove local minima



# Regularized logistic regression



- $\lambda=0.1$
- Now, both initializations lead to the same global minimum
- Surface is still non-convex, but the large flat area that led to a poor solution for the green path has been curved upwards

# CLASSIFICATION

---

# Two Classes (1/2)

- Classes are: C1 and C2
- Probabilistic view of clustering
  - we model the class conditional densities  $p(\mathbf{x}|C_i)$ ,  $i = 1, 2$
  - and the class priors  $p(C_i)$
- Posterior probability of class C1 can be written as:

$$\begin{aligned} p(C_1|\mathbf{x}) &= \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} = \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned}$$

- where we have defined:

$$a \triangleq \ln \left( \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)} \right)$$

# Two Classes (2/2)

- Important property of logistic function:  $\sigma(-a) = 1 - \sigma(a)$
- Inverse of logistic function:

$$a = \ln \left( \frac{\sigma}{1 - \sigma} \right) = \ln \left( \frac{p(\text{C}_1 | \mathbf{x})}{p(\text{C}_2 | \mathbf{x})} \right)$$

- It it called the **logit function**
- Represents the log of the ratio of probs. for the two classes
- **To see this:**

$$a = \ln \left( \frac{p(\mathbf{x} | \text{C}_1)p(\text{C}_1)}{p(\mathbf{x} | \text{C}_2)p(\text{C}_2)} \right) = \ln \left( \frac{\sigma(a)}{\sigma(-a)} \right) = \ln \left( \frac{\sigma(a)}{1 - \sigma(a)} \right)$$

# Multiple Classes

- For the case of  $K > 2$  classes, we have:

$$\begin{aligned} p(C_k | \mathbf{x}) &= \frac{p(\mathbf{x}|C_k)p(C_k)}{\sum_j p(\mathbf{x}|C_j)p(C_j)} = \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

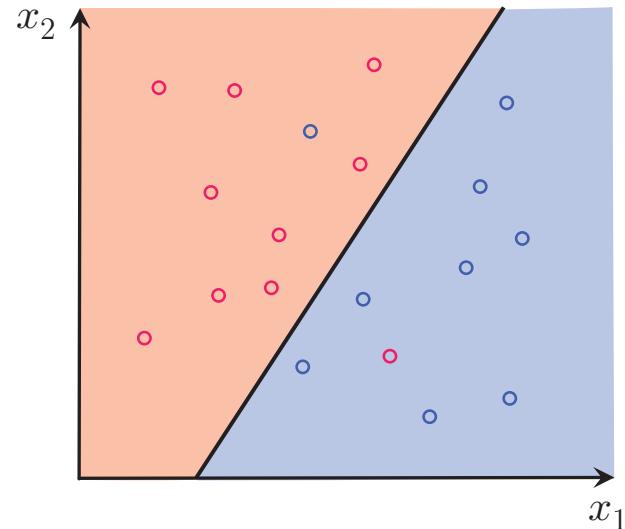
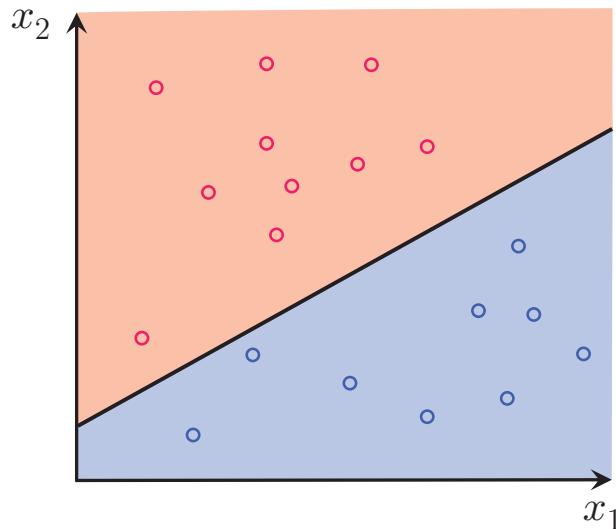
- which is known as the **normalized exponential**
- is a **multiclass generalization of the logistic function**, with:

$$a_k = \ln p(\mathbf{x}|C_k)p(C_k)$$

- It is also known as the **softmax function**
  - generalizes the **max function** as if:  $a_k \gg a_j$  for all  $j \neq k$

$$\rightarrow p(C_k | \mathbf{x}) = 1, p(C_j | \mathbf{x}) = 0$$

# From regression to classification



With linear classification we aim to learn a hyperplane  $b + \mathbf{x}^T \mathbf{w} = 0$  (shown here in black) to separate feature representations of the two classes, colored red (class “+1”) and blue (class “−1”), by dividing the feature space into a red half-space where  $b + \mathbf{x}^T \mathbf{w} > 0$ , and a blue half-space where  $b + \mathbf{x}^T \mathbf{w} < 0$ . (left panel) A linearly separable dataset where it is possible to learn a hyperplane to perfectly separate the two classes. (right panel) A dataset with two overlapping classes. Although the distribution of data does not allow for perfect linear separation, we can still find a hyperplane that minimizes the number of misclassified points that end up in the wrong half-space.

# Classification vs regression

- Linear regression: aim is to represent (approximate) a dataset
- Classification
  - Aim is to **separate two (or more) classes** of the input/output data **through a learned hyperplane**. In other words, we want to learn a hyperplane of the form:

$$b + \mathbf{x}_p^T \mathbf{w} = 0$$

- That separates the two classes of points as much as possible with
  - one class (C1) of points lying “above” the hyperplane with

$$b + \mathbf{x}_p^T \mathbf{w} > 0$$

- another class (C2) of points lying “below” the hyperplane with:

$$b + \mathbf{x}_p^T \mathbf{w} < 0$$

# Formally (1/2)

- We have a dataset  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ 
  - where:  
$$\mathbf{x}_i = [x_{1,i} \ x_{2,i} \ \dots \ x_{D,i}]^T$$
  - With each entry representing a feature (as with regression)
- For classification
  - The labels  $y_i$  are no longer continuous, but take values on a discrete set
  - They indicate the class membership of the corresponding point  $\mathbf{x}_i$
  - For two classes C1 and C2

$$y_i = +1 \text{ if } b + \mathbf{x}_i^T \mathbf{w} > 0 \text{ } (\mathbf{x}_i \text{ belongs to class C1})$$

$$y_i = -1 \text{ if } b + \mathbf{x}_i^T \mathbf{w} < 0 \text{ } (\mathbf{x}_i \text{ belongs to class C2})$$

## Formally (2/2)

- These relations

$$y_i = +1 \text{ if } b + \mathbf{x}_i^T \mathbf{w} > 0 \text{ (\mathbf{x}_i \text{ belongs to class C1})}$$

$$y_i = -1 \text{ if } b + \mathbf{x}_i^T \mathbf{w} < 0 \text{ (\mathbf{x}_i \text{ belongs to class C2})}$$

- Can be expressed compactly: by multiplying both expressions by their relative label value  $-y_i$ , leading to the equivalent expression for a correctly classified point  $\mathbf{x}_i$ :

$$-y_i(b + \mathbf{x}_i^T \mathbf{w}) < 0$$

- By taking the maximum between this expression and zero, a point (feature vector)  $\mathbf{x}_i$  is classified correctly if

$$\max(0, -y_i(b + \mathbf{x}_i^T \mathbf{w})) = 0$$

- In fact, if a point is classified correctly, this max must give zero, and it is  $>0$  only for a misclassified point

# Cost function

- Since for a **correctly classified point  $\mathbf{x}_i$** , we have:

$$\max(0, -y_i(b + \mathbf{x}_i^T \mathbf{w})) = 0$$

- For a **misclassified point  $\mathbf{x}_i$** , we have:

$$\max(0, -y_i(b + \mathbf{x}_i^T \mathbf{w})) > 0$$

- We can define a cost function as:

$$g(b, \mathbf{w}) = \sum_{n=1}^N \max(0, -y_n(b + \mathbf{x}_n^T \mathbf{w}))$$

- which amounts to “weighing” (somewhat counting) the number of *misclassified points*, as only these points generate a value strictly greater than zero

# Optimization problem

- The cost function is

$$g(b, \mathbf{w}) = \sum_{n=1}^N \max(0, -y_n(b + \mathbf{x}_n^T \mathbf{w}))$$

- We can use it to formulate an optimization problem as:

$$\underset{(b, \mathbf{w})}{\text{minimize}} \sum_{n=1}^N \max(0, -y_n(b + \mathbf{x}_n^T \mathbf{w}))$$

- Problem: this function is continuous, but *is not everywhere differentiable*, preventing the use of gradient and Newton's methods. Moreover has a trivial minimum at  $(b=0, \mathbf{w}=0)$

# Solution

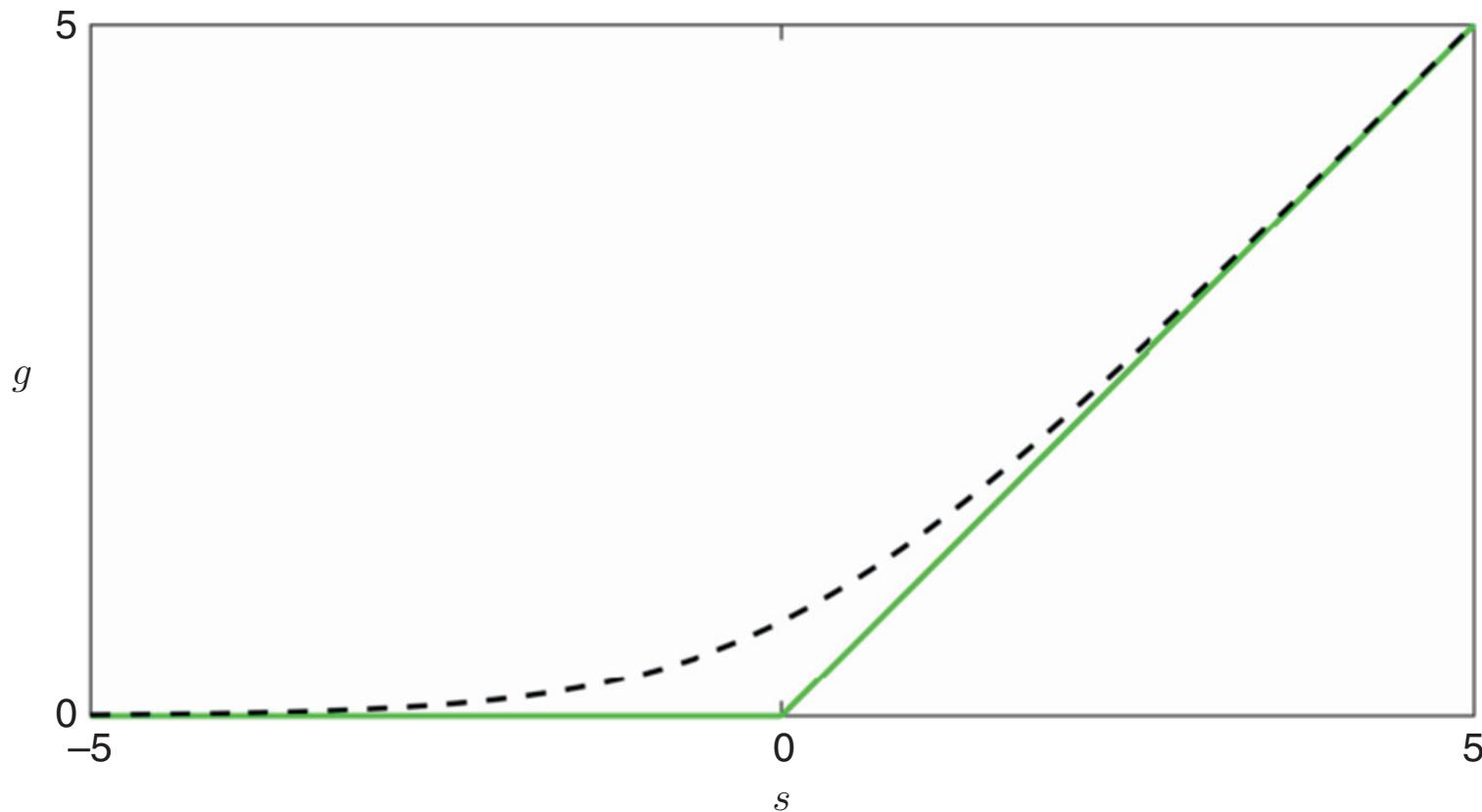
- Use a *smooth version* of the cost function

$$g(b, \mathbf{w}) = \sum_{n=1}^N \max(0, -y_n(b + \mathbf{x}_n^T \mathbf{w}))$$

- In particular, we replace the non-differentiable max function
  - with the differentiable softmax function

$$\text{soft}(s_1, s_2) \stackrel{\text{def}}{=} \log(e^{s_1} + e^{s_2})$$

# Softmax example



Plots of the non-differentiable perceptron or hinge cost  $g(s) = \max(0, s)$  (shown in green) as well as its smooth softmax approximation  $g(s) = \text{soft}(0, s) = \log(1 + e^s)$  (shown in dashed black).

$$\text{soft}(s_1, s_2) \stackrel{\text{def}}{=} \log(e^{s_1} + e^{s_2}) \approx \max(s_1, s_2)$$

# Cost function revisited

- Using the *softmax* function we approximate the max as

$$\begin{aligned}\max(0, -y_n(b + \mathbf{x}_n^T \mathbf{w})) &\approx \text{soft}(0, -y_n(b + \mathbf{x}_n^T \mathbf{w})) = \\ &= \log(1 + e^{-y_n(b + \mathbf{x}_n^T \mathbf{w})})\end{aligned}$$

- The cost function can thus be rewritten as:

$$\begin{aligned}g(b, \mathbf{w}) &= \sum_{n=1}^N \max(0, -y_n(b + \mathbf{x}_n^T \mathbf{w})) \\ &\approx \sum_{n=1}^N \log(1 + e^{-y_n(b + \mathbf{x}_n^T \mathbf{w})})\end{aligned}$$

# Softmax function for K inputs

- The *softmax* function generalizes to K inputs as:

$$\max(s_1, s_2, \dots, s_K) \approx \text{soft}(s_1, s_2, \dots, s_K) = \log \left( \sum_{k=1}^K e^{s_k} \right)$$

# Sources

- Chapter 3 and Chapter 4 (Section 4.1) of [1]

[1] J. Watt, R. Borhani, A. K. Katsaggelos, “Machine learning refined: foundations, algorithms and applications,” Cambridge University Press, 2016.

- Section 4.2 of [2]

[2] C. Bishop, “Pattern Recognition and Machine Learning,” Springer Verlag, 2006.

# A PRIMER ON REGRESSION AND CLASSIFICATION MODELS

---

Michele Rossi  
[rossi@dei.unipd.it](mailto:rossi@dei.unipd.it)

Dept. of Information Engineering  
University of Padova, IT

