

# Optimal Linear Classifiers: Support Vector Machines (SVM)

Nicolás Pérez de la Blanca

DECSAI-UGR

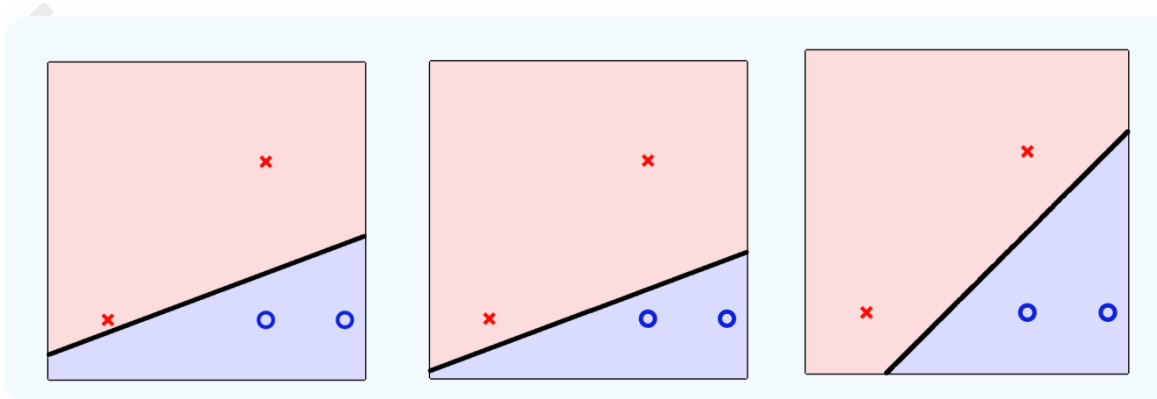
# Background

- Linear Models:
  - Powerful when NLT are used
  - VC dimension increase very fast with NLT
  - Difficult of tune with high order transformations
- Neural networks ( FeedForward)
  - Very high expressive power
  - Prone to overfitting
  - High computation time (training)
- Can we get the expressive power without paying the price ?
- Yes !! The Support Vector Machine does it
- How ? Implementing the SRM criteria as:

Keep the empirical error constant (small) y minimize the VC dimension

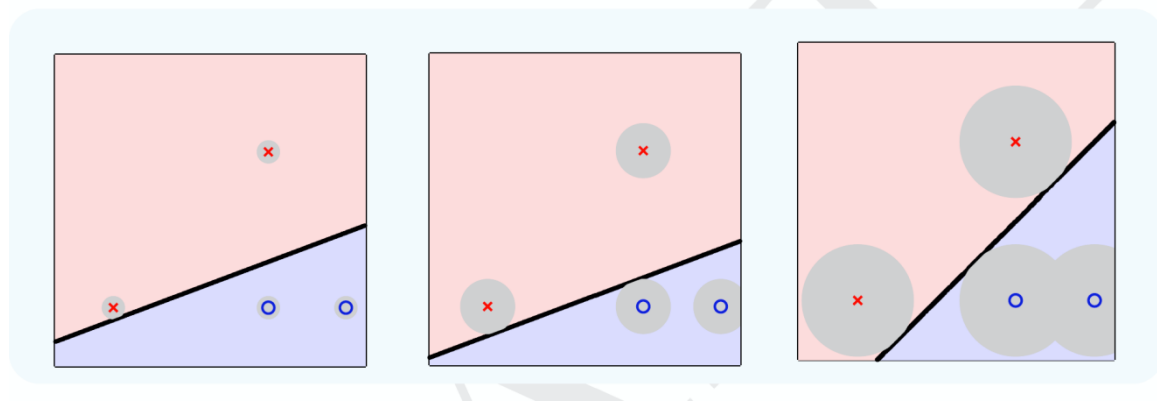
# The Optimal Hyperplane

- Let us revisit the perceptron model



Three possible separating hyperplanes, all three with  $E_{in}=0$  and the same VC- dimension (vector dimensionality)

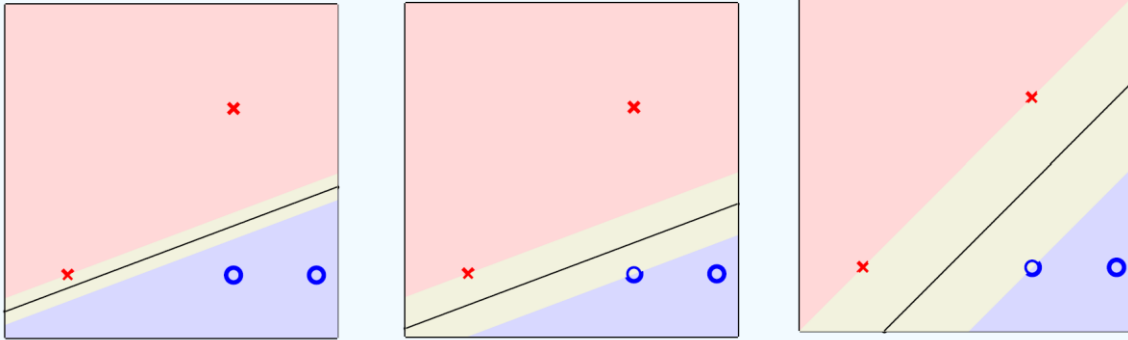
According to the VC-theory all have the same generalization bound



Question: Could we order our preferences?

Let's use noise tolerance !

# Fat separator



The robustness to the noise can be quantified by a fat separator.

The thickness reflect the amount of noise the separator can tolerate.

Margin: the maximum thickness possible for a separator

Three important questions:

1. Can we efficiently find the fattest separator? (Algorithm)
2. Why is a fat separator better than a thin one ? ( $E_{\text{out}}$ )
3. What should we do if the data is not separable? ( $E_{\text{in}}$ )

# Finding the Fattest Separating Hyperplane

- Notation: Now go back to denote an hyperplane in  $\mathbb{R}^d$  as:  $\mathbf{w}^T \mathbf{x} + b = 0$ 
  - $\mathbf{w}$  and  $b$  are estimated in different way.
  - $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ ,  $\mathbf{w} \in \mathbb{R}^d, \mathbf{x} \in \mathbb{R}^d, b \in \mathbb{R}$

- The hyperplane  $(\mathbf{w}, b)$  separates the data if y only if for  $i=1, \dots, N$

$$y_n(\mathbf{w}^T \mathbf{x}_n + b) > 0$$

- The magnitud depends on a free scale that must be normalized

$$\rho = \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b)$$

$$\min_{n=1, \dots, N} y_n \left( \frac{\mathbf{w}^T}{\rho} \mathbf{x}_n + \frac{b}{\rho} \right) = \frac{1}{\rho} \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = \frac{\rho}{\rho} = 1$$

- Thus, for any separating hyperplane , it is always possible to choose weights so that all the signals  $y_n(\mathbf{w}^T \mathbf{x}_n + b)$  are of magnitud greater or equal to 1, with equality satisfied by at least one  $(\mathbf{x}_n, y_n)$

# Finding the Fattest Separating Hyperplane

- **Definition** (Separating Hyperplane): The hyperplane  $h$  separates the data **if y only if** it can be represented by weights  $(\mathbf{w}, b)$  that satisfy  $\min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
- **Margin of a hyperplane:**
  - Consider the distance from a point to a hyperplane  $(\mathbf{w}, b)$ :

$$d(\mathbf{x}_n, h) = \frac{|\mathbf{w}^T \mathbf{x}_n + b|}{\|\mathbf{w}\|} = \frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|}$$

- Then, we can compute the distance from the separating hyperplane to the nearest point

$$\min_{n=1, \dots, N} \text{dist}(\mathbf{x}_n, h) = \frac{1}{\|\mathbf{w}\|} \cdot \min_{n=1, \dots, N} y_n(\mathbf{w}^T \mathbf{x}_n + b) = \frac{1}{\|\mathbf{w}\|}$$

- This simple expression for the distance of the the nearest point to the hyperplane is a consequence of above normalization
- Clearly, to maximize  $\frac{1}{\|\mathbf{w}\|}$  is equivalent to minimize  $\mathbf{w}^T \mathbf{w}$

# Finding the Optimal Hyperplane

To estimate the maximum-margin separating hyperplane we need to solve the following optimization problem:

$$\begin{aligned} & \underset{b, \mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{Subject to: } \min_n y_n (\mathbf{w}^T \mathbf{x}_n + b) = 1, \quad (n = 1, \dots, N) \end{aligned}$$

To make the optimization problem easier to solve, we approach

**PRIMAL**

$$\begin{aligned} & \underset{b, \mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ & \text{Subject to: } y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1, \quad (n = 1, \dots, N) \end{aligned}$$

It can be proved that the solution verify the constraints in equality

This problem can be written as:

$$\underset{\mathbf{u} \in \mathbb{R}^L}{\text{minimize}} \quad \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \mathbf{p}^T \mathbf{u}$$

$$\text{Subject to: } \mathbf{A} \mathbf{u} \geq \mathbf{c}$$

Quadratic Programming with Lineal Constraint

$$\mathbf{Q} = \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & \mathbf{I}_d \end{bmatrix}, \quad \mathbf{p} = \mathbf{0}_{d+1}$$

$$[y_n \ y_n \mathbf{x}_n^T] \mathbf{u} \geq 1$$

# Linear Hard-Margin SVM with QP

## Linear Hard-Margin SVM with QP

- 1: Let  $\mathbf{p} = \mathbf{0}_{d+1}$  ( $(d + 1)$ -dimensional zero vector) and  $\mathbf{c} = \mathbf{1}_N$  ( $N$ -dimensional vector of ones). Construct matrices  $\mathbf{Q}$  and  $\mathbf{A}$ , where

$$\mathbf{Q} = \begin{bmatrix} 0 & \mathbf{0}_d^T \\ \mathbf{0}_d & \mathbf{I}_d \end{bmatrix}, \quad \mathbf{A} = \underbrace{\begin{bmatrix} y_1 & -y_1 \mathbf{x}_1^T \\ \vdots & \vdots \\ y_N & -y_N \mathbf{x}_N^T \end{bmatrix}}_{\text{signed data matrix}}.$$

- 2: Calculate  $\begin{bmatrix} b^* \\ \mathbf{w}^* \end{bmatrix} = \mathbf{u}^* \leftarrow \text{QP}(\mathbf{Q}, \mathbf{p}, \mathbf{A}, \mathbf{c})$ .

- 3: Return the hypothesis  $g(\mathbf{x}) = \text{sign}(\mathbf{w}^{*\top} \mathbf{x} + b^*)$ .

- The solution depends only on a few data points, called “support vectors”.
- These points are those that verify the restriction in equality



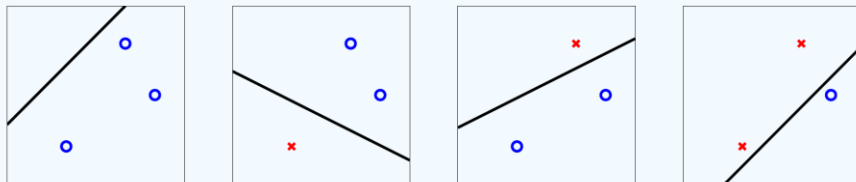
# Is a Fat Separator Better ?

- The SVM (optimal hyperplane) is a linear model, therefore  $d_{VC} = d + 1$ .
- The question is : does the support vector machine gain any more generalization ability by maximizing the margin ?
- Compare regularization with the optimal hyperplane problem

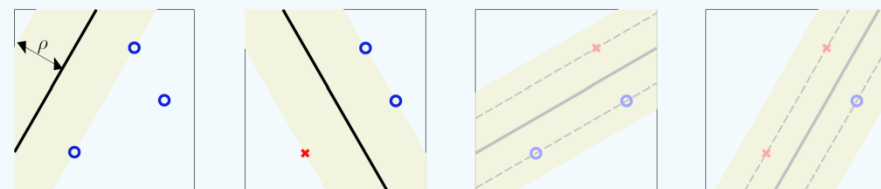
	optimal hyperplane	regularization
minimize:	$\mathbf{w}^T \mathbf{w}$	$E_{in}$
subject to:	$E_{in} = 0$	$\mathbf{w}^T \mathbf{w} \leq C$

- In SVM the error is fixed (to zero) and minimize the complexity of the model.
- There exist possibilities of improving in terms of VC dimension

# Fat Hyperplanes Shatter Fewer Points



Thin hyperplanes can implement all 8 dichotomies



Only 4 of the 8 dichotomies can be separated by hyperplanes with **thickness  $\rho$**

What matters is the thickness of the hyperplane relative to the spacing of the data

**MAIN RESULT: ( VC dimension of  $\rho$ -fat hyperplanes)**

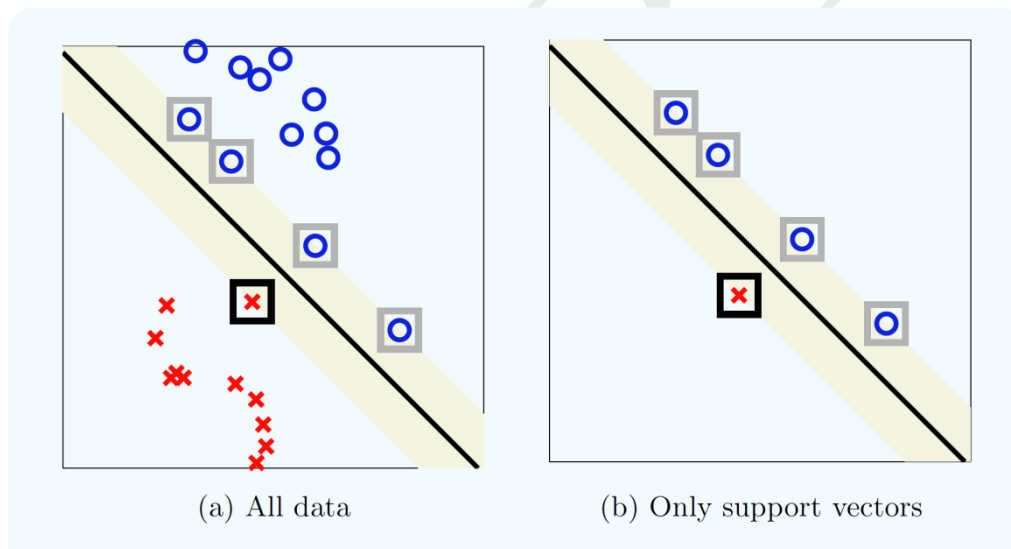
Suppose the input space is a ball of radius  $R$  in  $\mathbb{R}^d$ , so  $\|\mathbf{x}\| \leq R$ . Then

$$d_{\text{vc}}(\rho) \leq \left\lceil R^2 / \rho^2 \right\rceil + 1$$

where  $\left\lceil R^2 / \rho^2 \right\rceil$  is the smallest integer greater than or equal to  $R^2 / \rho^2$ . Combining results

$$d_{\text{vc}}(\rho) \leq \min \left( \left\lceil R^2 / \rho^2 \right\rceil, d \right) + 1$$

# Bounding the Cross Validation Error



$$E_{cv} = \frac{1}{N} \sum_{n=1}^N e_n \quad (\text{LOO})$$

But  $e_n = 0$  for any point not support vector

For the support vectors  $e_n \leq 1$   
(binary classification)

$$E_{out}(\text{SVM}) \leq E_{cv}(\text{SVM}) = \frac{1}{N} \sum_{n=1}^N e_n \leq \frac{\# \text{ support vectors}}{N}$$

**Clearly removing points not support vectors do not influence the solution**

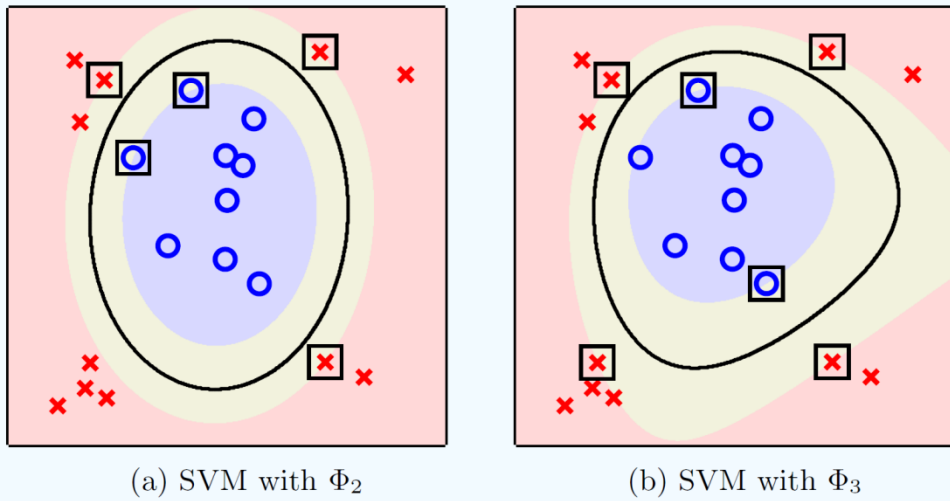
# Non-Separable data

- Let's assume the data are separable after some non-linear data transformation

$$\mathbf{z}_n = \Phi(\mathbf{x}_n)$$

- We solve the hard-margin SVM in the  $Z$ -space:

$$\Phi_2(\mathbf{x}) = (x_1, x_2, x_1^2, x_1x_2, x_2^2)$$



$$\begin{aligned} & \underset{\tilde{b}, \tilde{\mathbf{w}}}{\text{minimize}} \quad \frac{1}{2} \tilde{\mathbf{w}}^T \tilde{\mathbf{w}} \\ & \text{Subject to: } y_n (\tilde{\mathbf{w}}^T \mathbf{z}_n + \tilde{b}) \geq 1, \\ & \quad (n = 1, \dots, N) \end{aligned}$$

$$g(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^{*T} \Phi(\mathbf{x}) + \tilde{b}^*)$$

- NLT is used to lower  $E_{\text{in}}$ . However, you pay a price to get this sophisticated boundary in terms of a larger  $d_{\text{vc}}$  and a tendency to overfit.
- Nevertheless, only a light increasing of the boundary complexity is observed from  $\Phi_2$  to  $\Phi_3$

# Lagrange Dual for a QP-Problem

$$\underset{\mathbf{u} \in \mathbb{R}^L}{\text{minimize}} \quad \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \mathbf{p}^T \mathbf{u}$$

Subject to:  $\mathbf{a}^T \mathbf{u} \geq c$



$$\underset{\mathbf{u} \in \mathbb{R}^L}{\text{minimize}} \quad \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \mathbf{p}^T \mathbf{u} + \max_{\alpha \geq 0} \alpha (c - \mathbf{a}^T \mathbf{u})$$

$$\mathcal{L}(\mathbf{u}, \alpha) = \frac{1}{2} \mathbf{u}^T \mathbf{Q} \mathbf{u} + \mathbf{p}^T \mathbf{u} + \alpha (c - \mathbf{a}^T \mathbf{u})$$



$$\min_{\mathbf{u}} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{u}, \alpha)$$

- For convex quadratic programming when  $\mathcal{L}(\mathbf{u}, \alpha)$  has the same form that here and there exists  $\mathbf{u}$  such that  $c - \mathbf{a}^T \mathbf{u} \leq 0$ , then can be proved

$$\min_{\mathbf{u}} \max_{\alpha \geq 0} \mathcal{L}(\mathbf{u}, \alpha) = \max_{\alpha \geq 0} \min_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \alpha)$$

- The optimization problem on the RHS is easier to solve since the minimum in  $\mathbf{u}$  can be calculated analytically

# Dual of the Hard-Margin SVM

$$\underset{b, \mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{Subject to: } y_n(\mathbf{w}^T \mathbf{z}_n + b) \geq 1, \quad (n = 1, \dots, N)$$

$$\begin{aligned} \mathcal{L}(b, \mathbf{w}, \boldsymbol{\alpha}) &= \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{n=1}^N \alpha_n \left( 1 - y_n(\mathbf{w}^T \mathbf{z}_n + b) \right) \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{w}^T \mathbf{z}_n - b \sum_{n=1}^N \alpha_n y_n + \sum_{n=1}^N \alpha_n \end{aligned}$$

Very important to realize that we must first minimize  $\mathcal{L}$  respecto to  $(b, \mathbf{w})$  and then maximize with respect to  $\alpha \geq 0$  **(this encourage to satisfy the constraints)**

We need derivatives of  $\mathcal{L}$ :

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{n=1}^N \alpha_n y_n \quad \text{and} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n$$

Setting these derivatives to zero,

$$\sum_{n=1}^N \alpha_n y_n = 0$$

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{z}_n$$

# Dual of the Hard-Margin SVM

- Plugging the above stationary conditions back into the Lagrangian we get

$$\mathcal{L}(\alpha) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{z}_n^T \mathbf{z}_m + \sum_{n=1}^N \alpha_n$$

We must maximize  $\mathcal{L}(\alpha)$  subject to  $\alpha \geq \mathbf{0}$ , and the condition  $\sum_{n=1}^N \alpha_n y_n = 0$

We formulate this as a equivalent minimization

$$\underset{\alpha \in \mathbb{R}^N}{\text{minimize}} \quad \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{z}_n^T \mathbf{z}_m - \sum_{n=1}^N \alpha_n$$

$$\begin{aligned} \text{Subject to: } & \sum_{n=1}^N \alpha_n y_n = 0 \\ & \alpha_n \geq 0 \quad (n = 1, \dots, N) \end{aligned}$$

This is a new  
QP-problem  
in  $\alpha$

# Optimal solution conditions(KKT)

- The Karush-Khun-Tucker (KKT) conditions characterize the optimal solution  $(\mathbf{u}^*, \boldsymbol{\alpha}^*)$  of the primal and dual formulations of a Lagragian QP-problem,

$$\min_{\mathbf{u}} \max_{\boldsymbol{\alpha} \geq 0} \mathcal{L}(\mathbf{u}, \boldsymbol{\alpha}) = \max_{\boldsymbol{\alpha} \geq 0} \min_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \boldsymbol{\alpha})$$

1. Primal and dual constraints :  $\mathbf{a}_m^T \mathbf{u}^* \geq \mathbf{c}_m, \alpha_m \geq 0, m = 1, 2, \dots, M$
2. Complementary slackness:  $\alpha_m^* (y_m (\mathbf{z}_m^T \mathbf{w}^* + b^*) - 1) = 0 \quad (m = 1, \dots, M)$
3. Stationarity respect to  $\mathbf{u}$ :  $\nabla_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \boldsymbol{\alpha})|_{\mathbf{u}=\mathbf{u}^*, \boldsymbol{\alpha}=\boldsymbol{\alpha}^*} = \mathbf{0}$

- The consequence is
  1. The primal problem can be used to discover relations between the true unknowns and its Lagrange multipliers.
  2. These relations allow reformulate the primal problem as an equivalent optimization problem on the Lagrange multipliers (dual)
  3. The solution of the dual problem give us the solution of the primal one.



# ALGORITHM

## Hard-Margin SVM with Dual QP

1: Construct  $Q_D$  and  $A$  from the QP-Dual formulation

$$Q_D = \begin{bmatrix} y_1 y_1 \mathbf{z}_1^T \mathbf{z}_1 & \dots & y_1 y_N \mathbf{z}_1^T \mathbf{z}_N \\ y_2 y_1 \mathbf{z}_2^T \mathbf{z}_1 & \dots & y_2 y_N \mathbf{z}_2^T \mathbf{z}_N \\ \vdots & \vdots & \vdots \\ y_N y_1 \mathbf{z}_N^T \mathbf{z}_1 & \dots & y_N y_N \mathbf{z}_N^T \mathbf{z}_N \end{bmatrix} \quad \text{and} \quad A_D = \begin{bmatrix} \mathbf{y}^T \\ -\mathbf{y}^T \\ \mathbf{I}_{N \times N} \end{bmatrix}.$$

2: Use a QP-solver to optimize the dual problem:

$$\boldsymbol{\alpha}^* \leftarrow \text{QP}(Q_D, -\mathbf{1}_N, A_D, \mathbf{0}_{N+2}).$$

3: Let  $s$  be a support vector for which  $\alpha_s^* > 0$ . Compute  $b^*$ ,

$$b^* = y_s - \sum_{\alpha_n^* > 0} y_n \alpha_n^* \mathbf{z}_n^T \mathbf{z}_s.$$

4: Return the final hypothesis

$$g(\mathbf{x}) = \text{sign} \left( \sum_{\alpha_n^* > 0} y_n \alpha_n^* \mathbf{z}_n^T \mathbf{z} + b^* \right).$$

# Kernel Trick via Dual SVM

- Let's consider again the problem of solving nonlinear SVM after a nonlinear transform  $\Phi: \mathcal{X} \rightarrow \mathcal{Z}, \Phi(\mathbf{x}) = \mathbf{z}$
- We can observe that the only and important change in the dual formulation respect to the linear one is the needed of computing the  $\mathcal{Z}$ - space inner product

$$\Phi(\mathbf{x})^T \Phi(\mathbf{x}')$$

This inner product is needed to create  $\mathbf{Q}_D$  and to compute  $g(\mathbf{x})$   
When the dimension  $d'$  of  $\Phi(\mathbf{x})$  is very large (or infinity) the inner product computation is inefficient or non-viable.

The **kernel trick** defines a mechanism to compute the inner product without explicitly compute the function  $\Phi(\mathbf{x})$ .

Let's define a function that combines both the transform and inner product

$$K_{\Phi}(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}') = \mathbf{z}^T \mathbf{z}' \quad \text{Kernel Function}$$

# Hard-Margin SVM with Kernel

## Hard-Margin SVM with Kernel

- 1: Construct  $Q_D$  from the kernel  $K$ , and  $A_D$ :

$$Q_D = \begin{bmatrix} y_1 y_1 K_{11} & \dots & y_1 y_N K_{1N} \\ y_2 y_1 K_{21} & \dots & y_2 y_N K_{2N} \\ \vdots & \vdots & \vdots \\ y_N y_1 K_{N1} & \dots & y_N y_N K_{NN} \end{bmatrix} \quad \text{and} \quad A_D = \begin{bmatrix} \mathbf{y}^T \\ -\mathbf{y}^T \\ \mathbf{I}_{N \times N} \end{bmatrix},$$

where  $K_{mn} = K(\mathbf{x}_m, \mathbf{x}_n)$ . ( $K$  is called the *Gram* matrix.)

- 2: Use a QP-solver to optimize the dual problem:

$$\alpha^* \leftarrow \text{QP}(Q_D, -\mathbf{1}_N, A_D, \mathbf{0}_{N+2}).$$

- 3: Let  $s$  be any support vector for which  $\alpha_s^* > 0$ . Compute

$$b^* = y_s - \sum_{\alpha_n^* > 0} y_n \alpha_n^* K(\mathbf{x}_n, \mathbf{x}_s).$$

- 4: Return the final hypothesis

$$g(\mathbf{x}) = \text{sign} \left( \sum_{\alpha_n^* > 0} y_n \alpha_n^* K(\mathbf{x}_n, \mathbf{x}) + b^* \right).$$

The algorithm is equivalent to the lineal one. (HM SVM with Dual QP), but needing the kernel values

First of all we needs to characterize the class of function and its properties

# Polynomial kernel

Order two transformation

$$\Phi_2(\mathbf{x}) = (1, x_1, x_2, \dots, x_d, x_1x_1, x_1x_2, \dots, x_dx_d),$$

$$\Phi_2(\mathbf{x})^\top \Phi_2(\mathbf{x}') = 1 + \sum_{i=1}^d x_i x'_i + \sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j$$

$$\sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j = \sum_{i=1}^d x_i x'_i \times \sum_{j=1}^d x_j x'_j = (\mathbf{x}^\top \mathbf{x}')^2$$

$$K(\mathbf{x}, \mathbf{x}') = 1 + (\mathbf{x}^\top \mathbf{x}') + (\mathbf{x}^\top \mathbf{x}')^2$$

- This function is computed in time  $\mathcal{O}(d)$  instead of  $\mathcal{O}(\tilde{d})$ , where  $\tilde{d}$  is the dimension of  $\Phi_2(\mathbf{x})$ .
- This shows that for polynomial transformations, the function  $K$  represents a shortcut to combine the expansion and inner product.

# Example: $K(x, x') = (\zeta + \gamma x^T x')^Q$

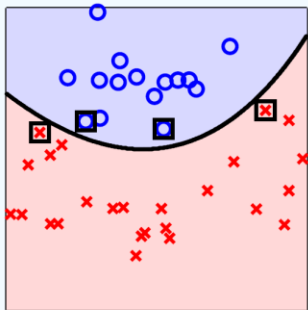
- This kernel is also called the *polynomial kernel*

Now  $K_\Phi(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}')$ , but coefficients are added  
 Let be  $\Phi(\mathbf{x}) = \mathbf{c}^T \Phi_2^d(\mathbf{x})$  where

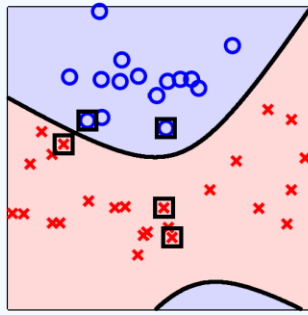
$\mathbf{c}^T = (1, \sqrt{2\gamma\zeta}, \sqrt{2\gamma\zeta}, \dots, \sqrt{2\gamma\zeta}, \gamma, \gamma, \dots, \gamma)$  coefficients

$\Phi_2^d(\mathbf{x}) = (1, x_1, x_2, \dots, x_d, x_1x_1, x_1x_2, \dots, x_dx_d)$  NLT

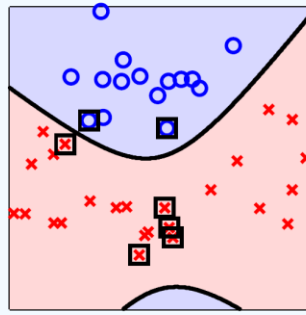
$$K(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \mathbf{x}^T \mathbf{x}')^2, \quad \zeta > 0 \quad \gamma > 0 \quad \mathbf{x} \in \mathbb{R}^d$$



$$(1 + 0.001 \mathbf{x}^T \mathbf{x}')^2$$



$$1 + (\mathbf{x}^T \mathbf{x}') + (\mathbf{x}^T \mathbf{x}')^2$$

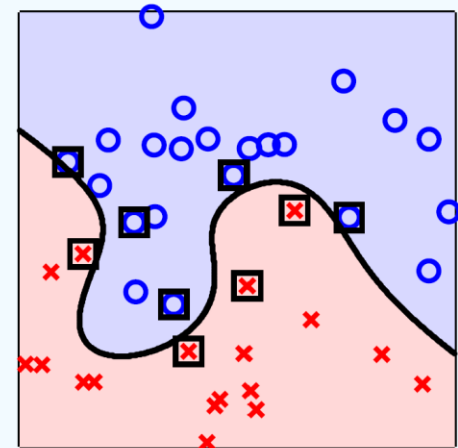


$$(1 + 1000 \mathbf{x}^T \mathbf{x}')^2$$

Examples of different coefficient values

Observe that for  $\zeta=1$  values higher than 1 for  $\gamma$  appears to have a low influence in the solution.

$$K(x, x') = (\zeta + \gamma x^T x')^{10}$$

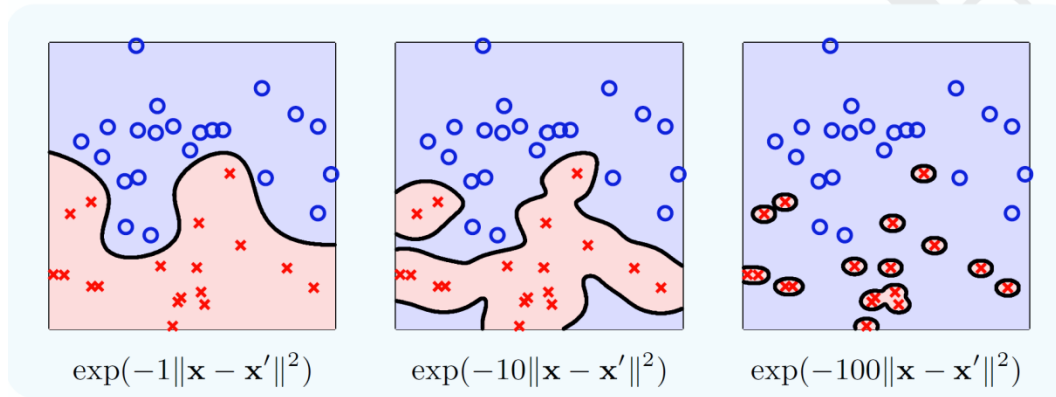


We can efficiently use high-dimensional kernels.

The model complexity is controlled by maximizing the margin

# Example: $K(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$

- This kernel is called the *Gaussian-RBF Kernel*, where  $\gamma > 0$ ,  $\mathbf{x} \in \mathbb{R}^d$



- To understand the implicit NLT let consider the case  $\gamma = 1$ ,  $\mathbf{x} \in \mathbb{R}$

$$\begin{aligned} K(x, x') &= \exp(-x^2) \cdot \exp(2xx') \cdot \exp(-x'^2) \\ &= \exp(-x^2) \cdot \exp\left(\sum_{k=0}^{\infty} \frac{2^k x^k x'^k}{k!}\right) \cdot \exp(-x'^2) \end{aligned}$$

$$\Phi(x) = \exp(-x^2) \cdot \left(1, \sqrt{\frac{2^1}{1!}}x, \sqrt{\frac{2^1}{1!}}x^2, \sqrt{\frac{2^1}{1!}}x^3, \dots\right)$$

Note that in this case  $\Phi(x)$  is an infinity NLT

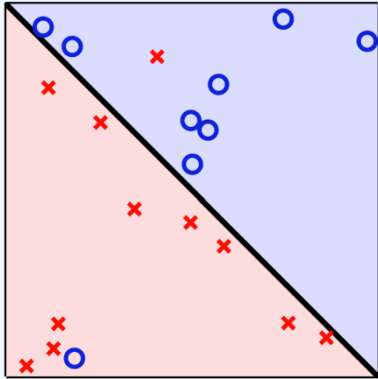
# Choice of Kernels

- Three kernels are popular in practice:
  - Linear: Polynomial con  $Q=1, \zeta=0, \gamma=1$
  - Polynomial:  $Q < 10$  and appropriated  $\zeta$  and  $\gamma$  ( not easy of selecting)
  - Gaussian-RBF: only one parameter  $\gamma \in [0,1]$ .
- Many other kernels can be proposed as combinations of simpler kernels.
- New kernels can be proposed but it is not an easy task since the function  $K(x, x')$  must accomplish the **Mercer's condition**, that is its Gram Matrix **K** is **positive semidefinite** for any given sample  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$

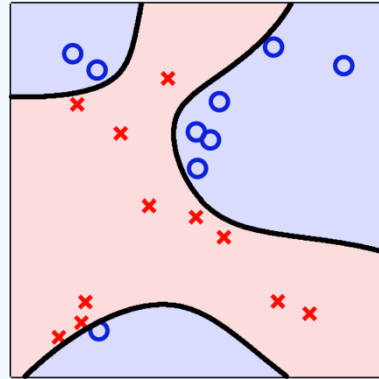
$$K = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_N, \mathbf{x}_1) & \cdots & K(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

# Soft-margin SVM

- The hard-margin SVM assumes that the data are separable in the  $\mathcal{Z}$  space



(a) linear classifier



(b) Gaussian-RBF kernel

But according to the type of noise, overfitting can appear after the transformation.

A best strategy is to assume some noise in the labels using a “soft” fomulation

*Soft-Margin SVM:* Let introduce an amount of margin violation  $\xi_n \geq 0$  for each point  $(\mathbf{x}_n, y_n)$



$$y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n$$

$$\underset{\mathbf{b}, \mathbf{w}, \xi}{\text{minimize}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{n=1}^N \xi_n$$

$$\text{Subject to: } y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \text{ for } n = 1, 2, \dots, N$$
$$\xi_n \geq 0 \text{ for } n = 1, 2, \dots, N$$

This problem can be solved by quadratic programming using the dual problem



# Soft-margin SVM

## DUAL

$$\begin{aligned} & \underset{\alpha}{\text{minimize}} \quad \frac{1}{2} \alpha^T Q_D \alpha + \mathbf{1}^T \alpha \\ & \text{Subject to: } \mathbf{y}^T \alpha = 0 \\ & \quad \mathbf{0} \leq \alpha \leq C \cdot \mathbf{1} \end{aligned}$$

The main difference with the above dual problem is the upper bound of  $\alpha$

**Solution:** similar to hard-SVM

Support vectors  $0 \leq \alpha_n^* \leq C$

Those with  $0 < \alpha_n^* < C$  (Margin S.V.)

Those with  $\alpha_n^* = C$  (Non-Margin S.V.)

$$g(x) = \text{sign}(\mathbf{w}^{*T} \mathbf{x} + b^*) = \text{sign} \left( \sum_{\alpha_n^* > 0} y_n \alpha_n^* \mathbf{x}_n^T \mathbf{x} + b^* \right)$$

The kernel trick is still applicable

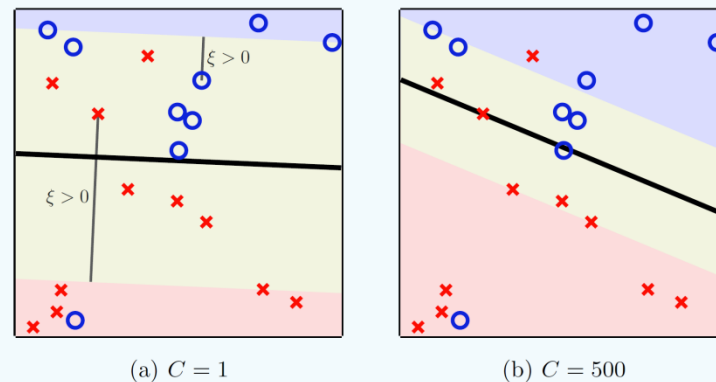
$b^*$  can be estimated from any margin S.V.

$$y_n (w^{*T} x_n + b^*) = 1$$

In other cases there are a range of solutions

The KKT conditions now are:

$$\begin{aligned} \alpha_n^* \cdot (y_n (w^{*T} x_n + b^*) - 1 + \xi_n^*) &= 0 \\ \beta_n^* \cdot \xi_n^* &= (C - \alpha_n^*) \cdot \xi_n^* = 0 \end{aligned}$$



$$E_{cv} \leq \frac{\#[\alpha_n^* > 0]}{N}$$

# Soft-margin SVM & Regularization

- Let's take a closer look at the parameter  $C$ 
  - If  $C$  is large, it means that we do want all errors  $\xi_n$  to be as small as possible with the trade-off of higher complexity hypothesis
  - If  $C$  is small, we will tolerate some amounts of errors with less complicated hypothesis

Such trade-off was found when we studied regularization

Let's define

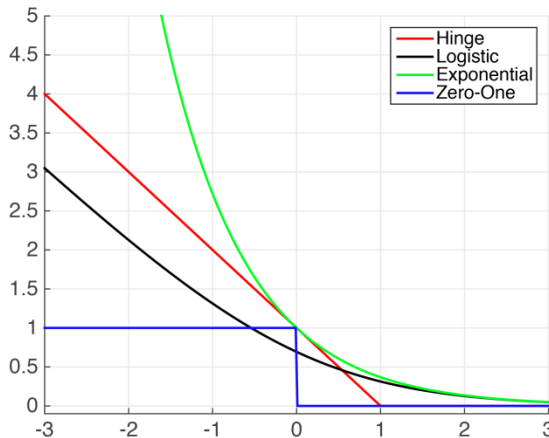
$$E_{SVM}(b, \mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \max(1 - y_n(\mathbf{w}^T \mathbf{x}_n + b), 0)$$

Then the soft-margin SVM can be written as the following optimization problem:

$$\min_{b, \mathbf{w}} \lambda \mathbf{w}^T \mathbf{w} + E_{SVM}(b, \mathbf{w}), \quad \lambda = 1/2CN$$

Soft-margin SVM can be viewed as a special case of regularized classification with  $E_{SVM}(b, \mathbf{w})$  as a surrogate for the in-sample error and  $\mathbf{w}^T \mathbf{w}$  as the regularizer

# Soft-SVM and SGD



$$\min_{b, \mathbf{w}} \lambda \mathbf{w}^T \mathbf{w} + E_{SVM}(b, \mathbf{w}), \quad \lambda = 1/2CN$$



$$\min_{b, \mathbf{w}} \lambda R(\mathbf{w}) + L(\mathbf{x}, y, b, \mathbf{w}),$$

**SVM** : Hinge Loss

**Perceptron**: Zero-One

**Logistic Regression**: Logistic

**Boosting**: Exponential

A strong connection between different binary classifiers appears as consequence that all of them solve the same problem minimizing a different surrogate for the zero-one loss function, and using a different algorithm to reach its solution.

**For separable or almost separable data set the soft-SVM provides in general the best option **but** for overlapping classes LR can provide a better performance.**

SVM is not the only classifier that take advantage of the kernel trick, LR and many others can be formulated to be used with kernels.

# Soft-Margin SVM: Summary

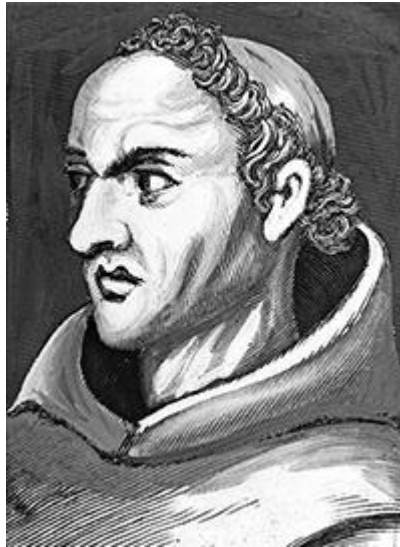
- It can deliver a large-margin hyperplane, and in so doing, it can control the effective model complexity.
- It can deal with high- or infinite-dimensional transforms using the kernel trick.
- It expresses the final hypothesis  $g(x)$  using only a few support vectors, their corresponding Lagrange multipliers, and the kernel.
- It can control the sensitivity to outliers and regularize the solution by setting  $C$  appropriately
- When  $C$  y the kernel are chosen properly, the soft-margin SVM enjoy a low  $E_{\text{out}}$  and define one of the most useful classification models.

# Three Learning Principles

# We will discuss.....

- **Occam's razor**: pick a model carefully
- **Sample Bias**: generate data carefully
- **Data Snooping**: handle the data carefully

# Occam's razor



Entities should not be multiplied beyond necessity “Occam’s razor”  
principle attributed to William of Occam c. 1280–1349

We should seek simpler models over complex ones and optimize the tradeoff between model complexity and the accuracy of model’s description of the training data

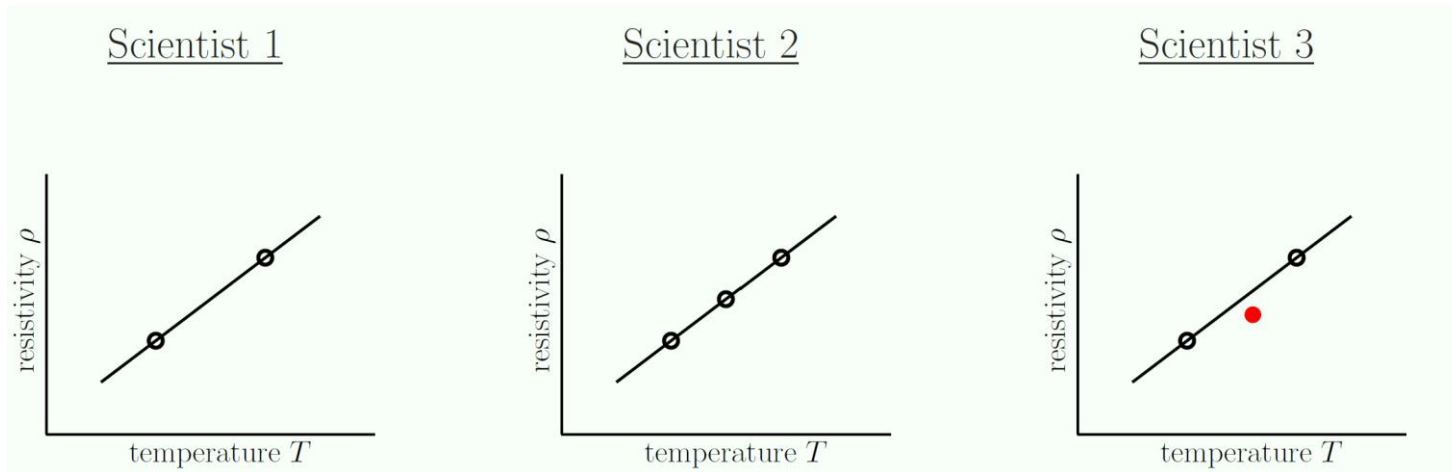
# Why Simpler is Better... ?

- Mathematically : many arguments ( lower VC dimension, less capability to fit noise, etc,etc)

Simple is better because you will be **more surprised** when you fit the data

If something unlikely happens, it is very significant when it happens.

A scientific experiment

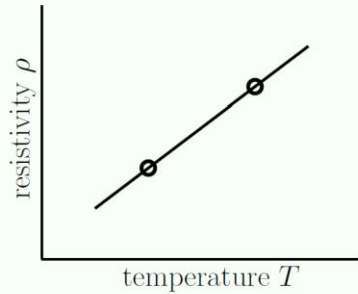


Who provides most evidence to the hypothesis “ $\rho$  is linear in  $T$ ” ?



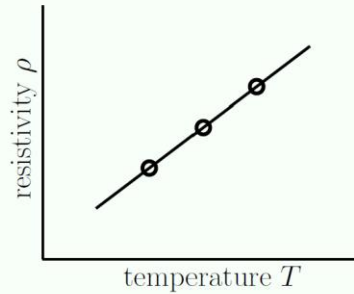
# Scientific Experiment

Scientist 1



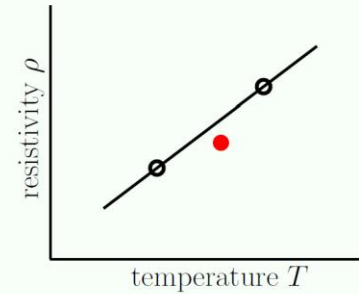
no evidence

Scientist 2



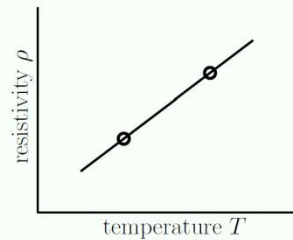
Very convincing

Scientist 3



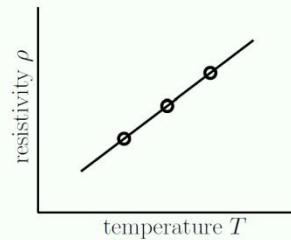
some evidence?

Scientist 1



no evidence

Scientist 2



very convincing

## Axiom of Non-Falsifiability.

*If an experiment has no chance of falsifying a hypothesis, then the result of that experiment provides no evidence one way or the other for the hypothesis.*

# Falsification and $m_{\mathcal{H}}(N)$

- If  $\mathcal{H}$  shatter  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$   
... don't be surprised if you fit the data
- If  $\mathcal{H}$  doesn't shatter  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ , and the target values are uniformly distributed,

$$\mathbb{P}[\text{falsification}] \geq 1 - \frac{m_{\mathcal{H}}(N)}{2^N}$$

A good fit is surprising with simpler  $\mathcal{H}$  (small  $m_{\mathcal{H}}(N)$ ) and hence more significant.

## GOING BEYOND OCCAM'S RAZOR . . .

We may opt for 'a simpler fit than possible', namely an imperfect fit of the data using a simple model over a perfect fit using a more complex one. The reason is that the price we pay for a perfect fit in terms of the penalty for model complexity may be too much in comparison to the benefit of the better fit.

# Sampling Bias in Learning

If the data is sampled in a biased way, learning will produce a similarly biased outcome.

. . . or, make sure the training and test distributions are the same.

**You cannot draw a sample from one distribution and make claims about another distribution**

Example.1: Consider the data set of customer loan from a Bank. ¿where is the bias when used to build a rule for new credit appointment ?

Example.2: Consider the data set of historical values of the current trading companies to select companies in which to invest ¿where is the bias?

# Data Snooping

If a data set has affected any step in the learning process, it cannot be fully trusted in assessing the outcome.

- ... or, estimate performance with a completely uncontaminated test set
- ... and, **choose  $\mathcal{H}$  before** looking at the data

**Example :** Consider a 8 years data set of the daily changes USD/EURO. We want predict UP/DOWN.

Before fitting a model:

- 1.- We normalize the fully data set.
- 2.- We separate the last two years from the data set as Test Set.
- 3.- We fit a model
- 4.- The fitted model works very well on the data set,
- 5.- But when used with new data the prediction error was very high,  
is there any rational explanation ?

# Data Snooping is a Subtle Happy Hell

- The data looks linear, so I will use a linear model, and it worked.
  - If the data were different and didn't look linear, would you do something different?
- Try linear, it fails; try circles it works.
  - If you torture the data enough, it will confess.
- Try linear, it works; so I don't need to try circles.
  - Would you have tried circles if the data were different?
- Read papers, see what others did on the data. Modify and improve on that.
  - If the data were different, would that modify what others did and hence what you did?
  - the data snooping can happen all at once or sequentially by different people
- Input normalization: normalize the data, now set aside the test set.
  - Since the test set was involved in the normalization, wouldn't your  $g$  change if the test set changed?