

UCB - CS189  
Introduction to Machine Learning  
Fall 2015

Lecture 11: Embedded methods for  
model and feature selection

Isabelle Guyon  
ChaLearn

Come to my office hours...

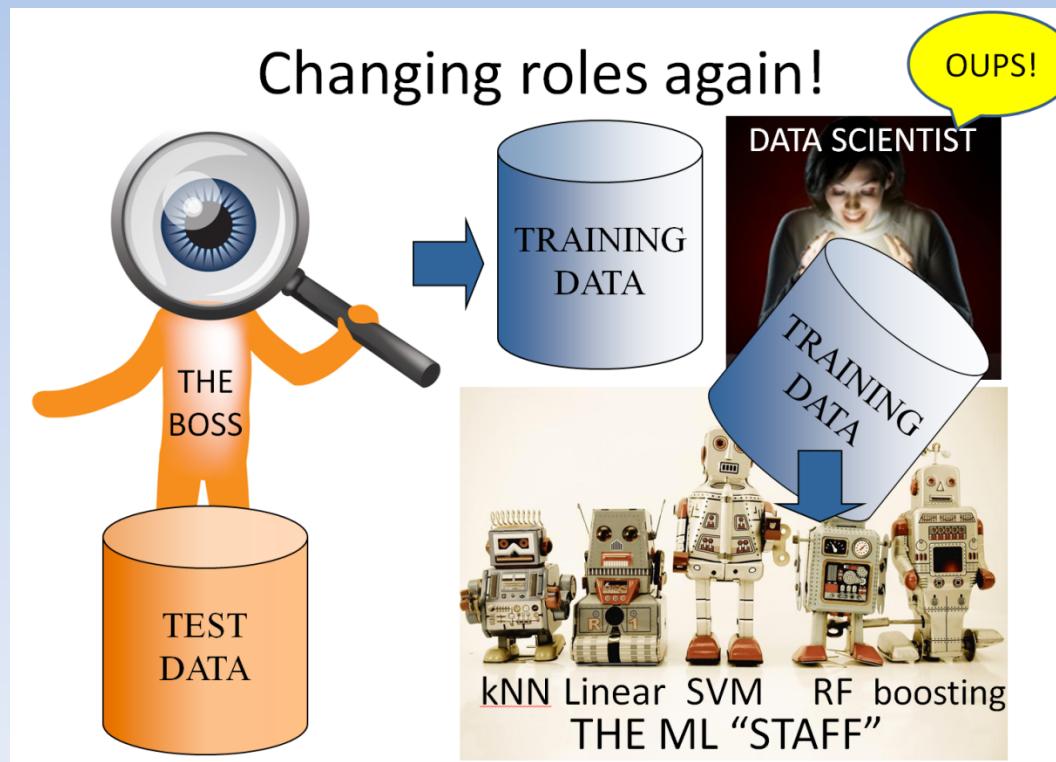
Wed 2:30-4:30 Soda 329

Last time: model search

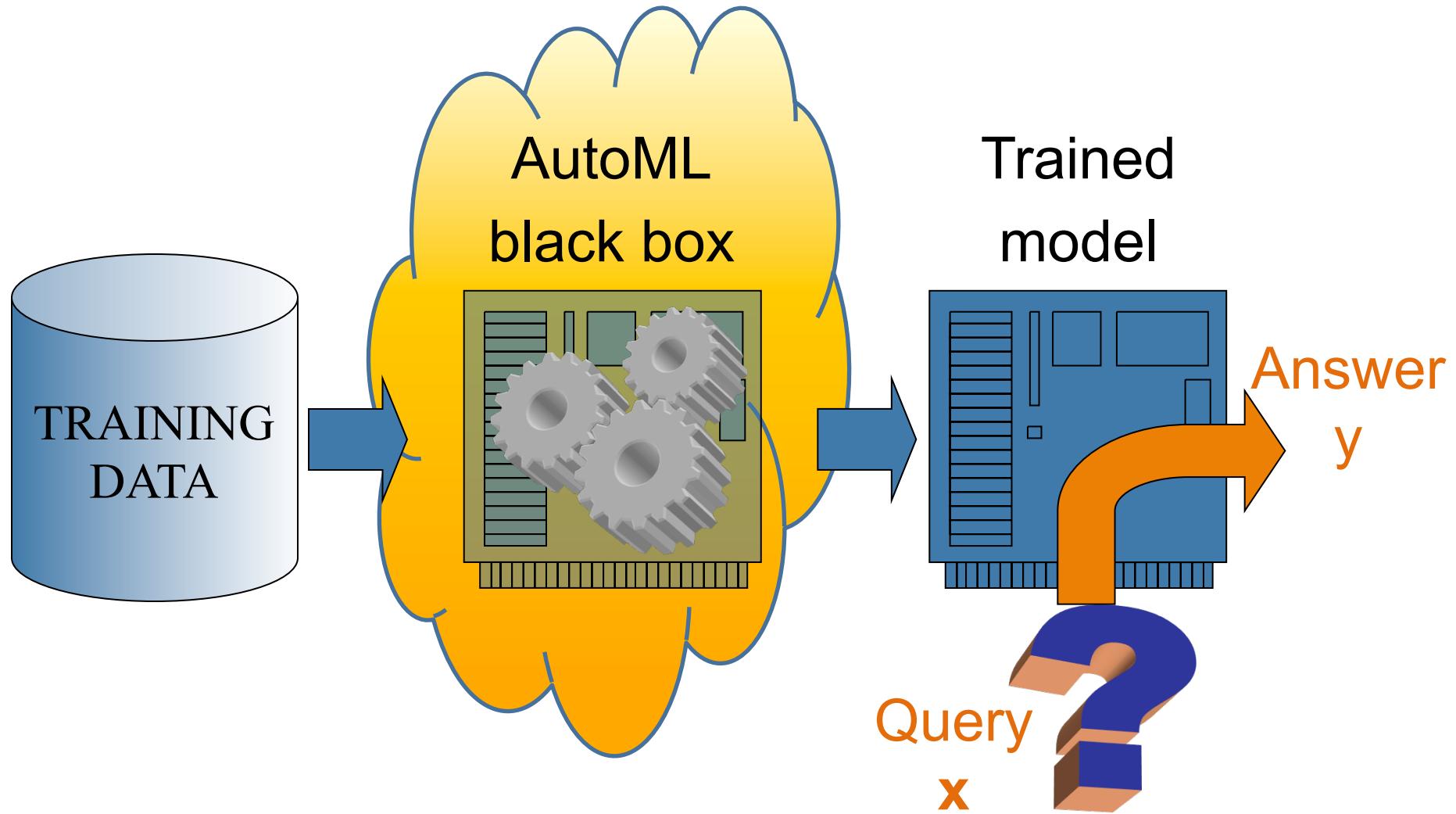


Come to my office hours...  
Wed 2:30-4:30 Soda 329

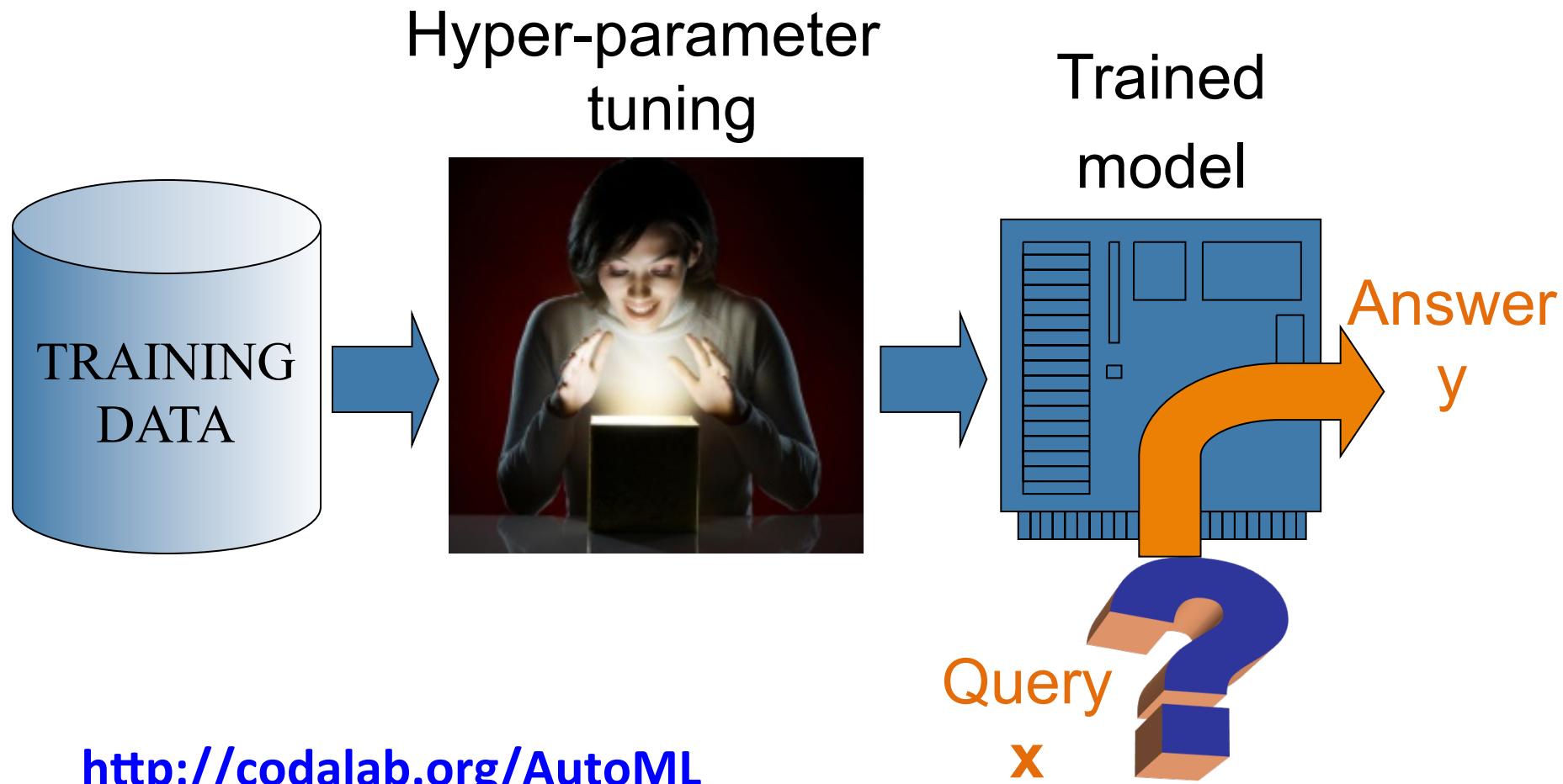
## Today



# Remember: The DREAM

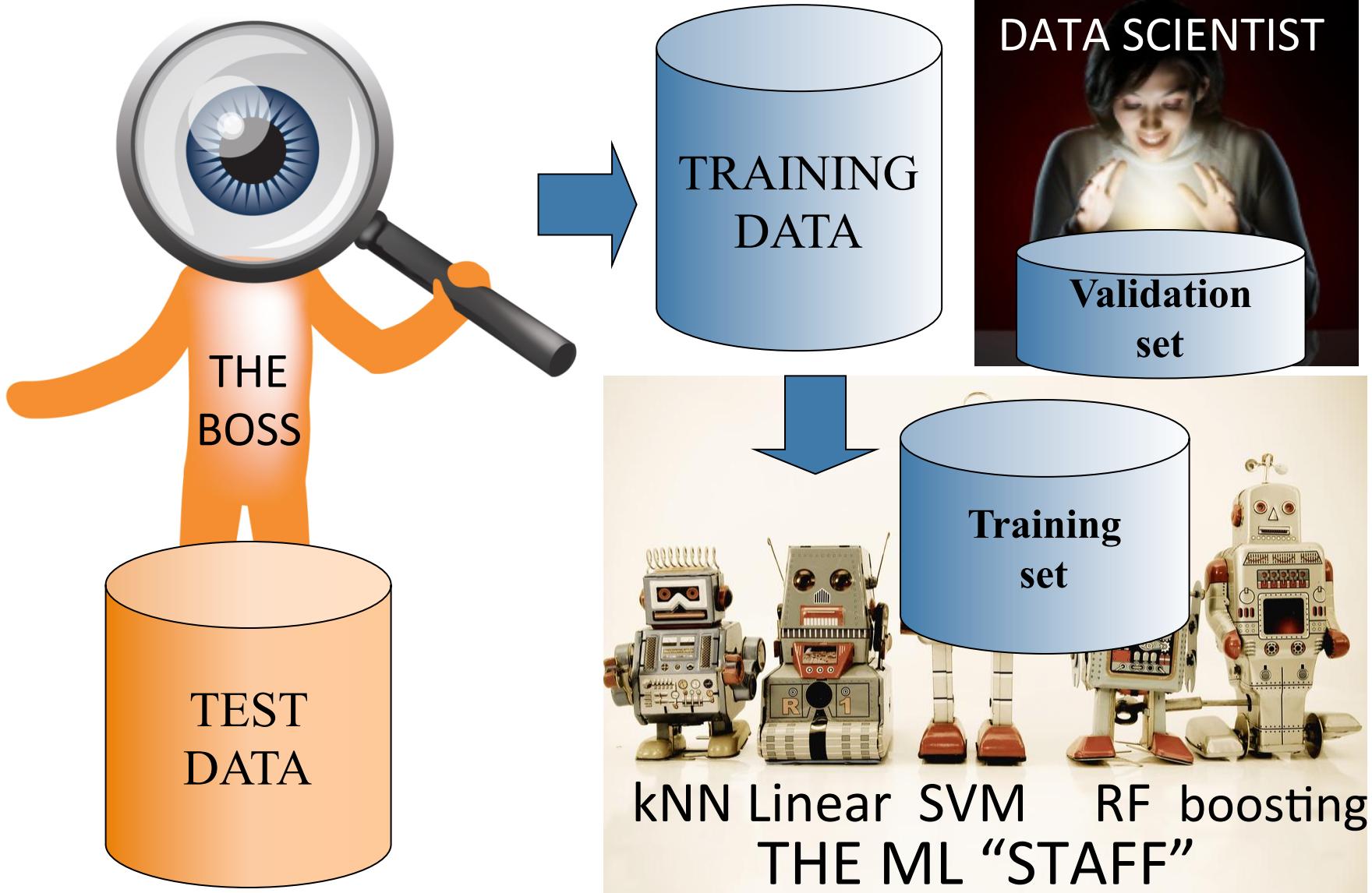


# Remember: The REALITY



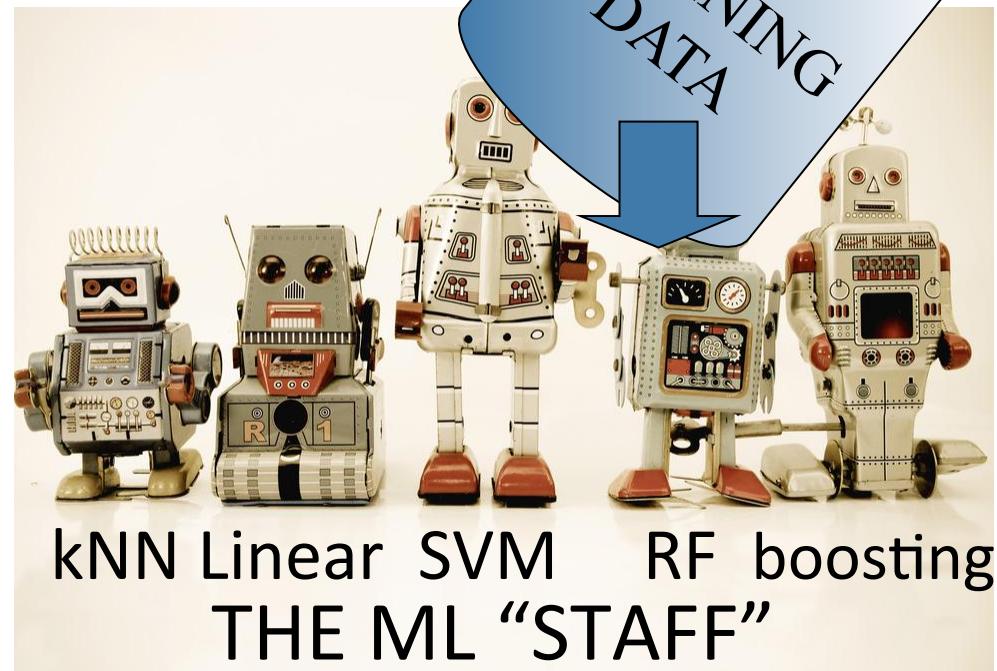
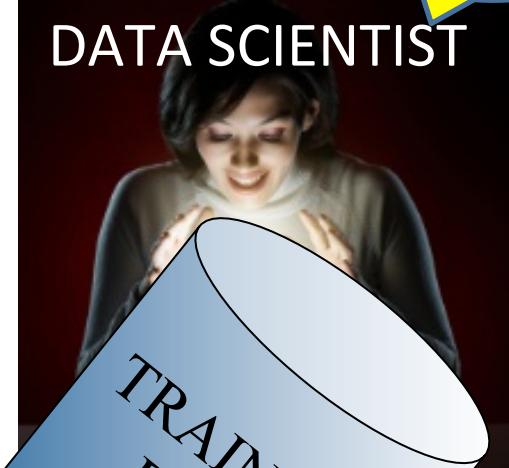
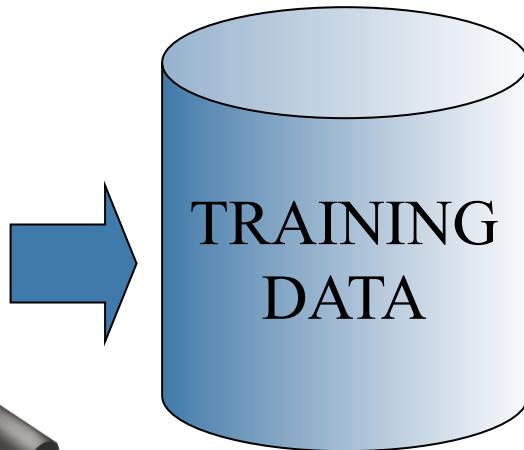
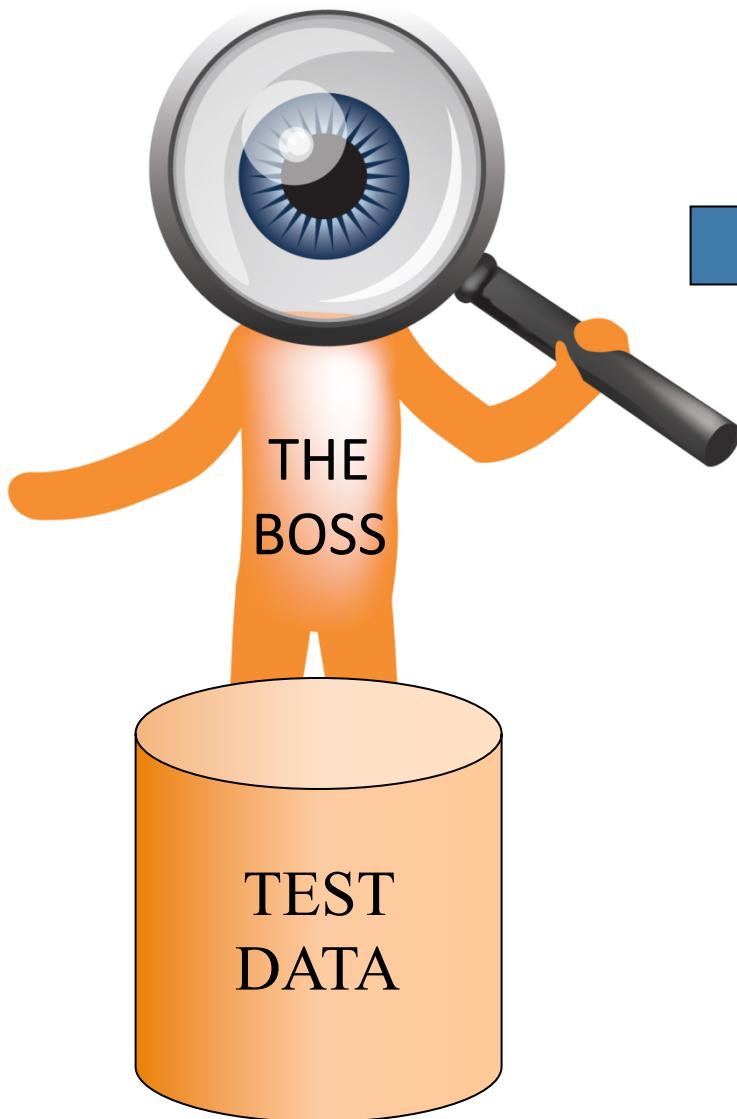
<http://codalab.org/AutoML>

# Last time you were the data scientist



# Changing roles again!

OUPS!



# Remember your options

- 1. Reduce the number of hyper-parameters.**
- 2. Learn how to do “fancy” search.**



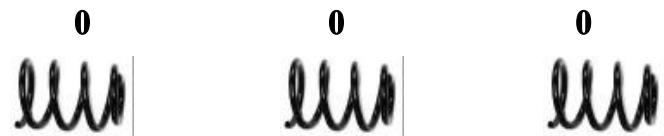
# Remember your options

1. Reduce the number of hyper-parameters.
2. Learn how to do “fancy” search.
3. Push more hyper-parameters to the “lower level”.

# Embedded methods

000  


00      0 0      00  

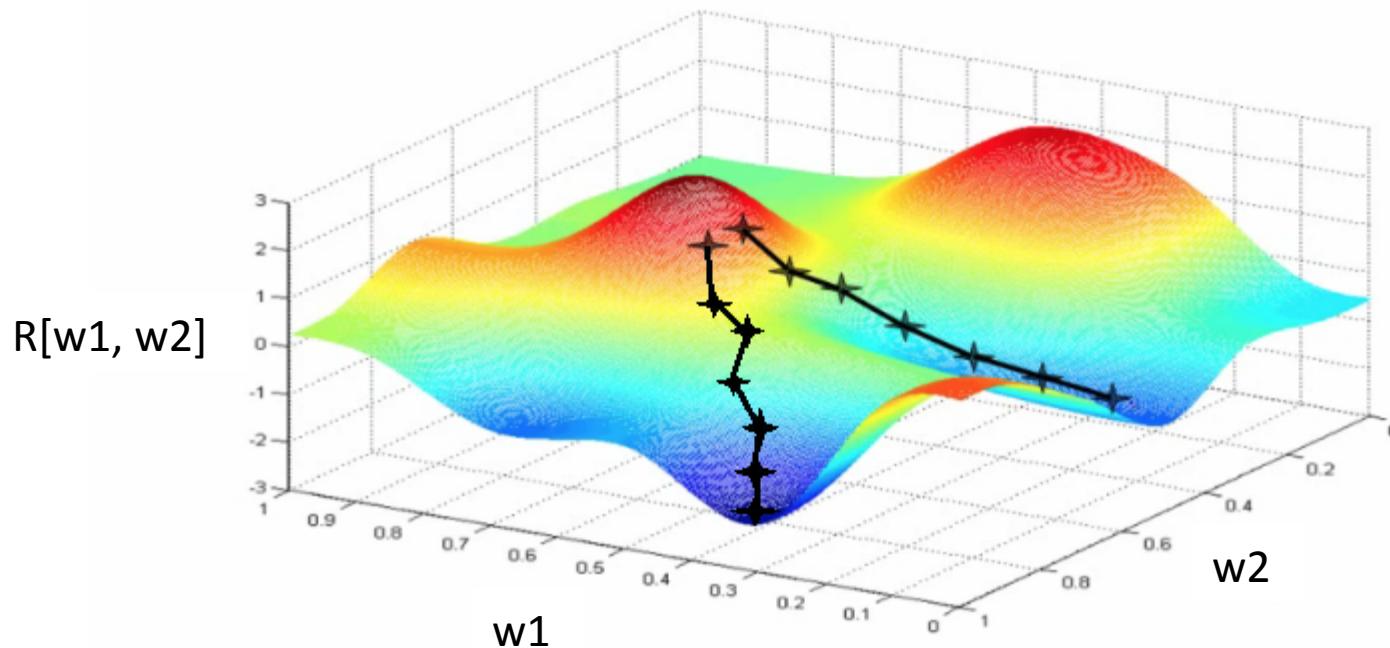

0      0      0  




# Why do we need HP?

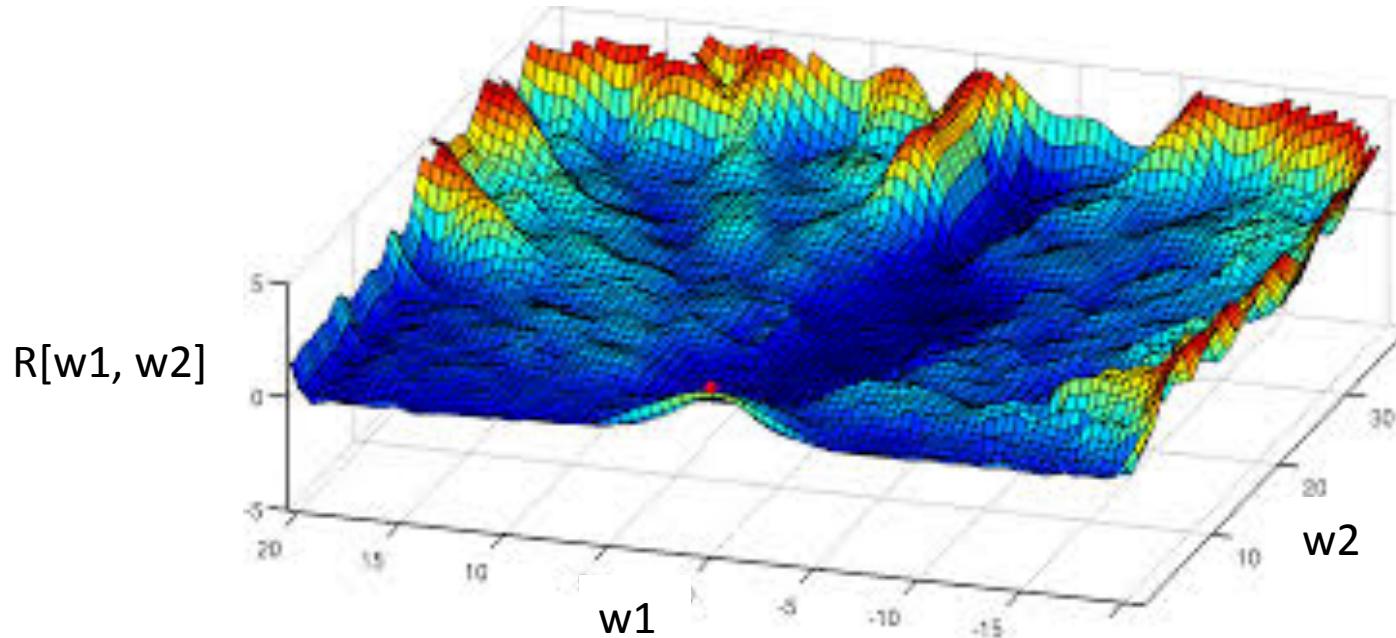
- Could we not just optimize the parameters and hyper-parameters together with the training error?
- No, because:
  - **Computational complexity**: we loose convexity in parameter space.
  - **Statistical complexity**: we get infinite (or zero) capacity.

# Gradient descent falls into local minima...

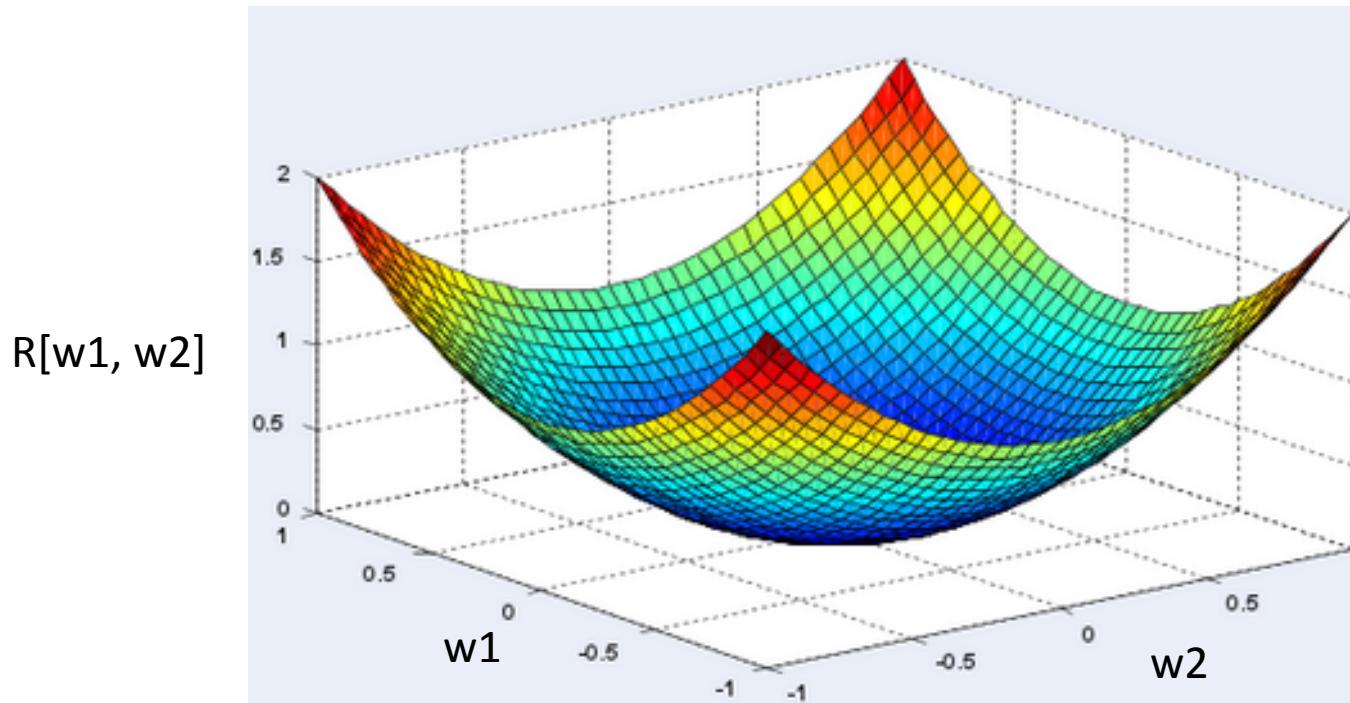


$\eta$  is the learning rate of gradient step.

... finding the global optimum  
can be hard ...



... except if  
the risk functional is convex!



Convexity: Hessian matrix is positive semidefinite.

# Convexity

- $R = RSS = \sum_k (\mathbf{x}^k \mathbf{w}^T - y^k)^2$   
 $= \| X\mathbf{w}^T - \mathbf{y} \|^2$   
 $= \mathbf{w}^T X^T X \mathbf{w} - 2\mathbf{w}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$

- Gradient:

$$\nabla_{\mathbf{w}} R = 2 (X^T X \mathbf{w}^T - X^T \mathbf{y})$$

- Hessian:

$$H(R) = 2 X^T X$$

$X^T X$  is an outer product, therefore it is PSD.

# Kernel trick

- Hessian:

$$H(R) = 2 X^T X$$

- In  $\Phi$ -space:

$$\begin{aligned} H(R) &= 2 \Phi^T \Phi \\ &= 2 K \end{aligned}$$

This is why it is useful that the kernel matrix  $K$  be PSD.

# SVM in $\alpha$ space

- $f(\mathbf{x}) = \sum_k \alpha_k k(\mathbf{x}^k, \mathbf{x})$

$$K = [k(\mathbf{x}^k, \mathbf{x}^h)]$$

## 1) Hinge loss version:

$$\min_{\alpha} \alpha^T \mathbf{y} - (1/2) \alpha^T K \alpha$$

$$0 \leq \alpha_k y_k \leq 1/\lambda \quad \text{box constraint}$$

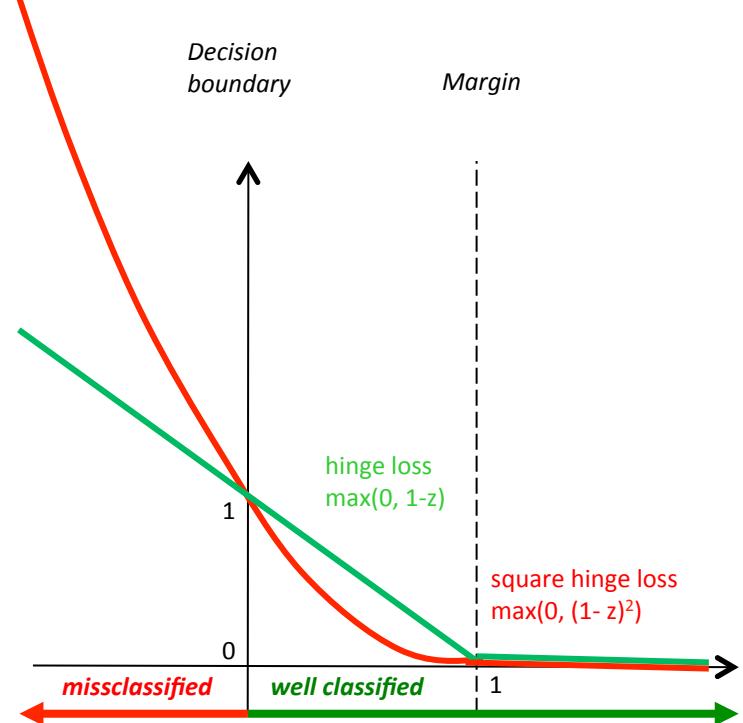
$$\alpha^T \mathbf{1} = 0 \quad \text{bias constraint}$$

## 2) Square hinge loss version

$$\min_{\alpha} \alpha^T \mathbf{y} - (1/2) \alpha^T (K + \lambda I) \alpha$$

$$\alpha_k y_k \geq 0 \quad \text{ridge}$$

$$\alpha^T \mathbf{1} = 0 \quad \text{bias constraint}$$



# Loosing convexity

$$R[\alpha, \theta] = \alpha^T y - (1/2) \alpha^T K(\theta) \alpha$$

$$K(\theta) = [ k(x^k, x^h; \theta) ]_{k=1:N, h=1:N}$$

$$k(x, x'; \theta) = (a + x \cdot x')^b \exp(-c \|x - x'\|^d)$$

$$\theta = [a, b, c, d]$$

# Avoid overfitting...

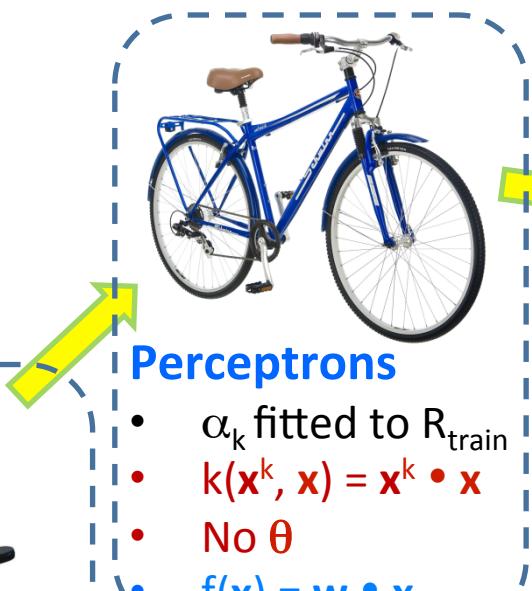
$$f(x) = \left( \sum_k \alpha_k [\Phi(x^k)] \cdot \Phi(x) \right)$$

**w**       **$k(x^k, x; \theta)$**



## Hebb's method

- $\alpha_k = y_k$
- $k(x^k, x) = x^k \cdot x$
- No  $\theta$
- $f(x) = w \cdot x$



## Parzen windows

- $\alpha_k = y_k$
- Any  $k(x^k, x; \theta)$
- $\theta$  fitted to  $R_{CV}$
- $f(x) = \sum_k \alpha_k k(x^k, x; \theta)$



## Kernel methods

- $\alpha_k$  fitted to  $R_{train}$
- Any  $k(x^k, x; \theta)$
- $\theta$  fitted to  $R_{CV}$
- $f(x) = \sum_k \alpha_k k(x^k, x; \theta)$

# Nested subsets

Capacity /complexity



## Hebb's method

- $\alpha_k = y_k$
- $k(x^k, x) = x^k \bullet x$
- No  $\theta$
- $f(x) = w \bullet x$



## Perceptrons

- $\alpha_k$  fitted to  $R_{train}$
- $k(x^k, x) = x^k \bullet x$
- No  $\theta$
- $f(x) = w \bullet x$



## Kernel methods

- $\alpha_k$  fitted to  $R_{train}$
- Any  $k(x^k, x; \theta)$
- $\theta$  fitted to  $R_{CV}$
- $f(x) = \sum_k \alpha_k k(x^k, x; \theta)$

# Nested subsets

Capacity /complexity



## Hebb's method

- $\alpha_k = y_k$
- $k(x^k, x) = x^k \cdot x$
- No  $\theta$
- $f(x) = w \cdot x$

## Parzen windows

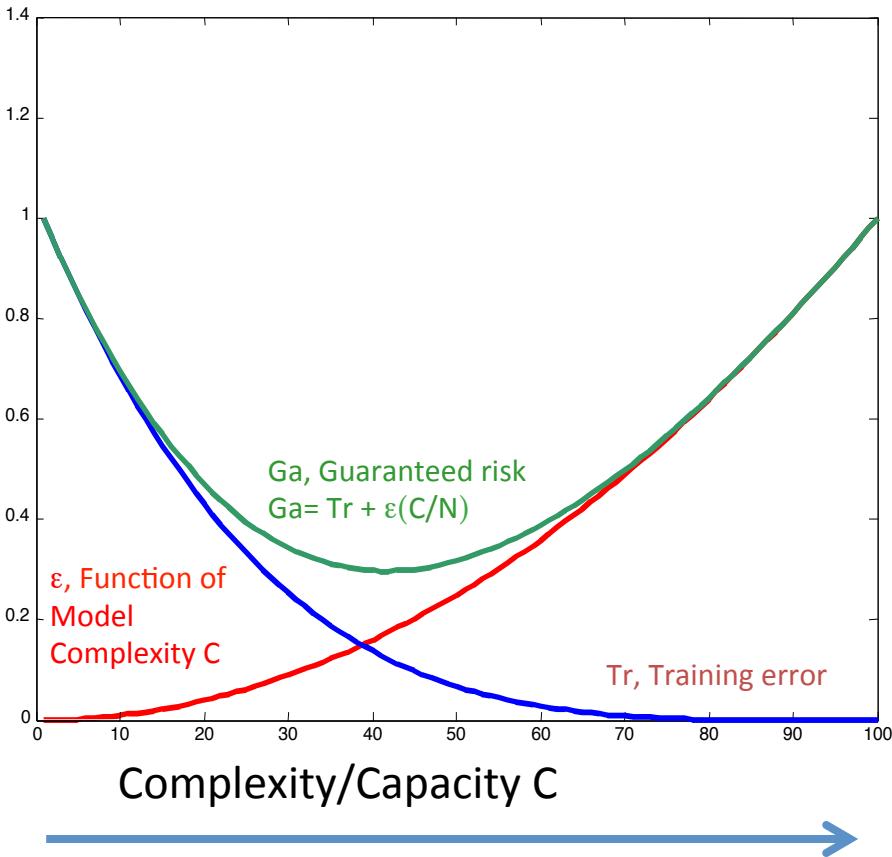
- $\alpha_k = y_k$
- Any  $k(x^k, x; \theta)$
- $\theta$  fitted to  $R_{CV}$
- $f(x) = \sum_k \alpha_k k(x^k, x; \theta)$



## Kernel methods

- $\alpha_k$  fitted to  $R_{train}$
- Any  $k(x^k, x; \theta)$
- $\theta$  fitted to  $R_{CV}$
- $f(x) = \sum_k \alpha_k k(x^k, x; \theta)$

# Guaranteed risk to avoid overfitting



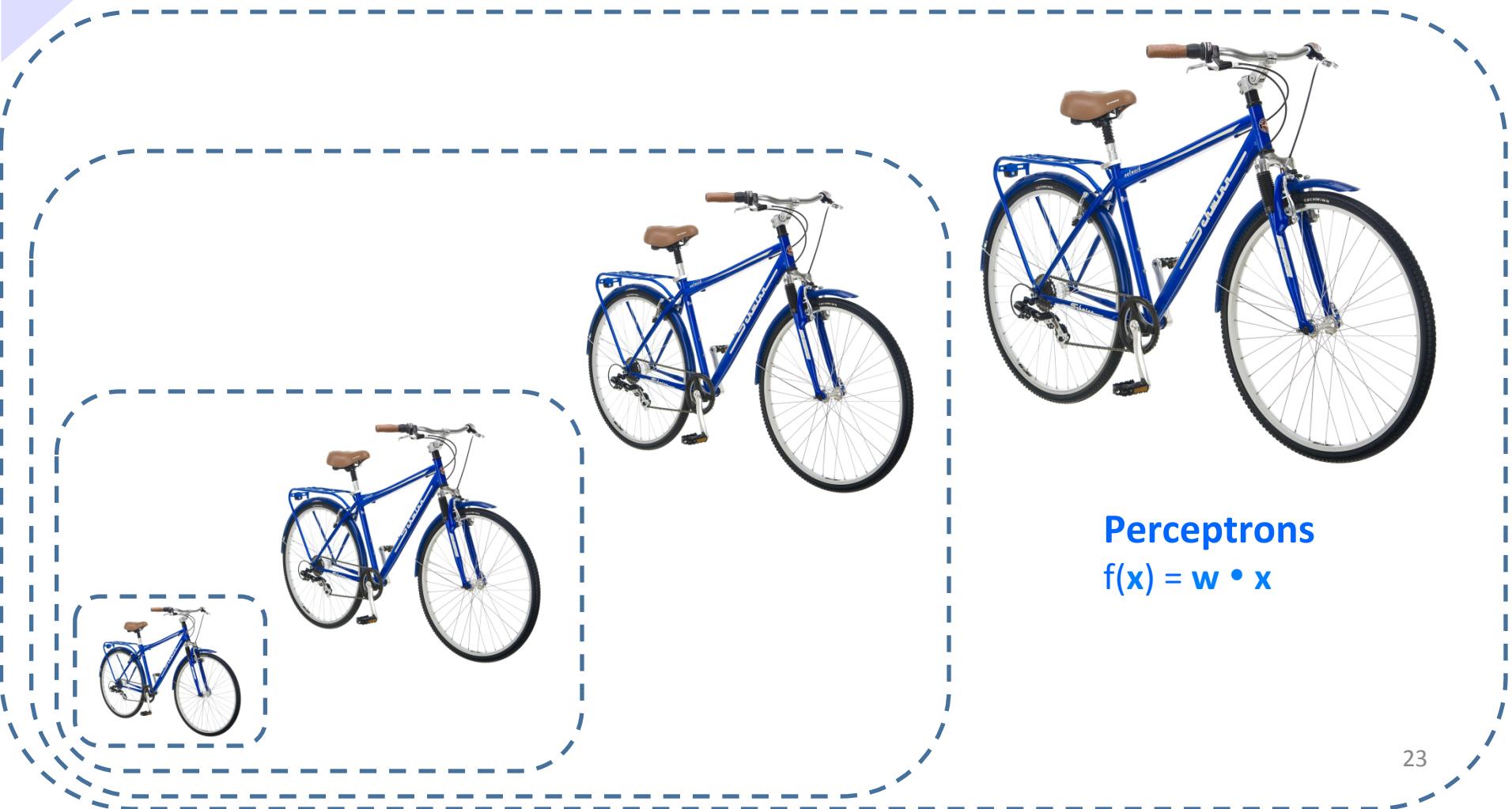
With *high* probability  $(1-\delta)$ ,  
 $0 < \delta \ll 1$

$$R[f] \leq R_{\text{train}}[f] + \varepsilon(\delta, C/N)$$

$R_{\text{gua}}[f]$

$$\text{Shrinkage: } \min_w R_{\text{train}}[w] + \lambda \|w\|^2$$

Capacity =  $d_{\text{effective}}$



Perceptrons  
 $f(x) = w \cdot x$

Shrinkage:  $\min_w R_{\text{train}}[w] + \lambda \|w\|^2$

Capacity =  $d_{\text{effective}}$

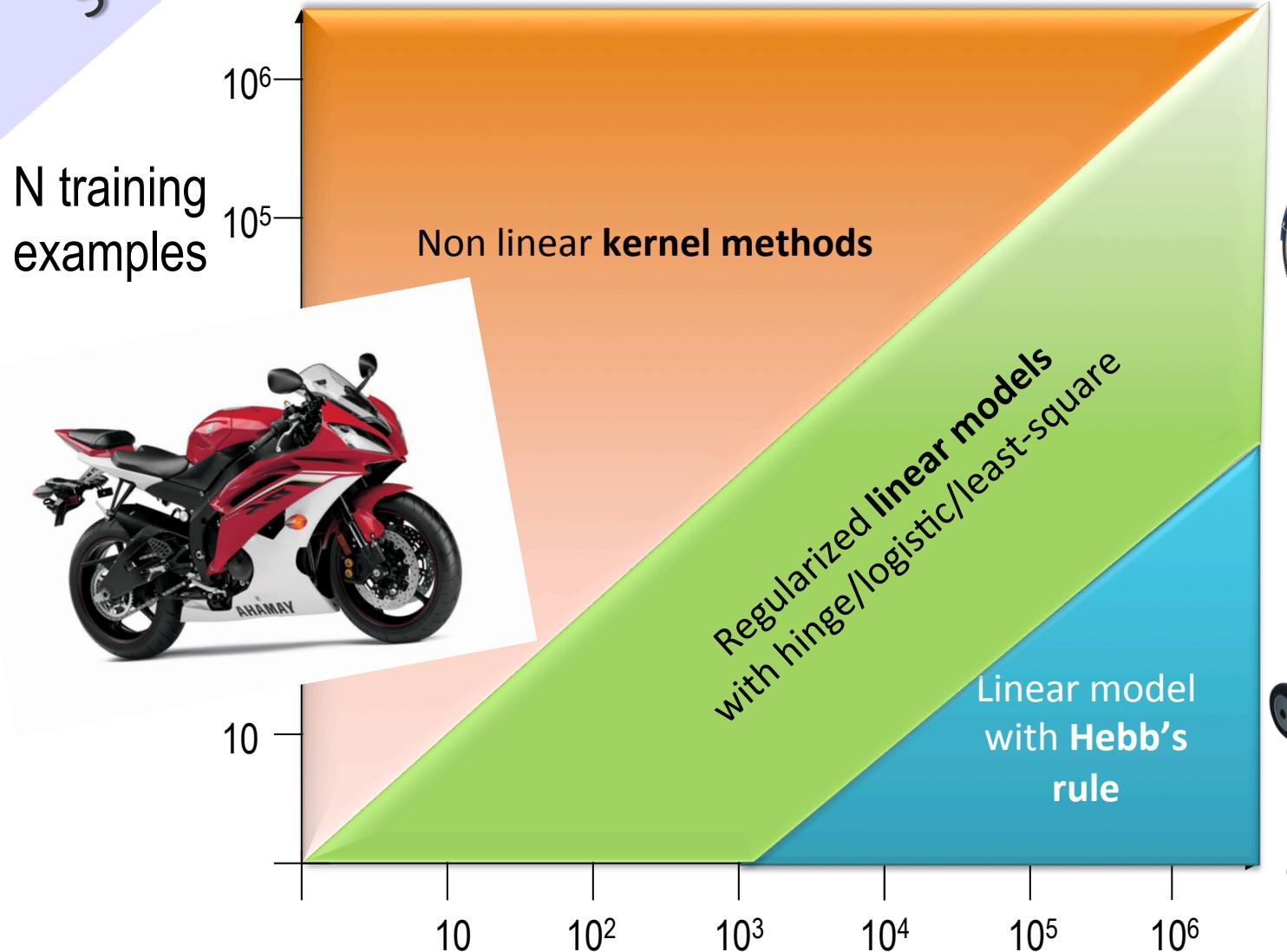


Kernel methods

$$f(x) = \sum_k \alpha_k k(x^k, x; \theta)$$

$$\begin{aligned}\|w\|^2 &= w w^\top \\ &= \alpha^\top \Phi \Phi^\top \alpha \\ &= \alpha^\top K \alpha\end{aligned}$$

# Statistical domains

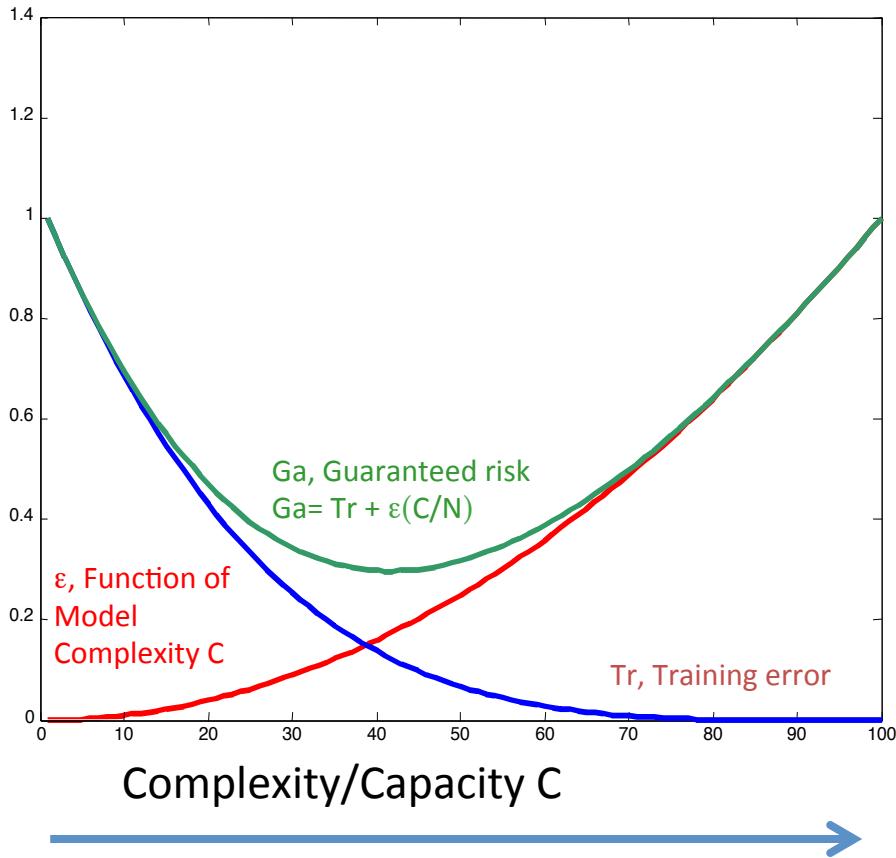


# What you want:

## Keep C/N (or $d_{\text{eff}}/N$ ) under control



# Guaranteed risk to avoid overfitting



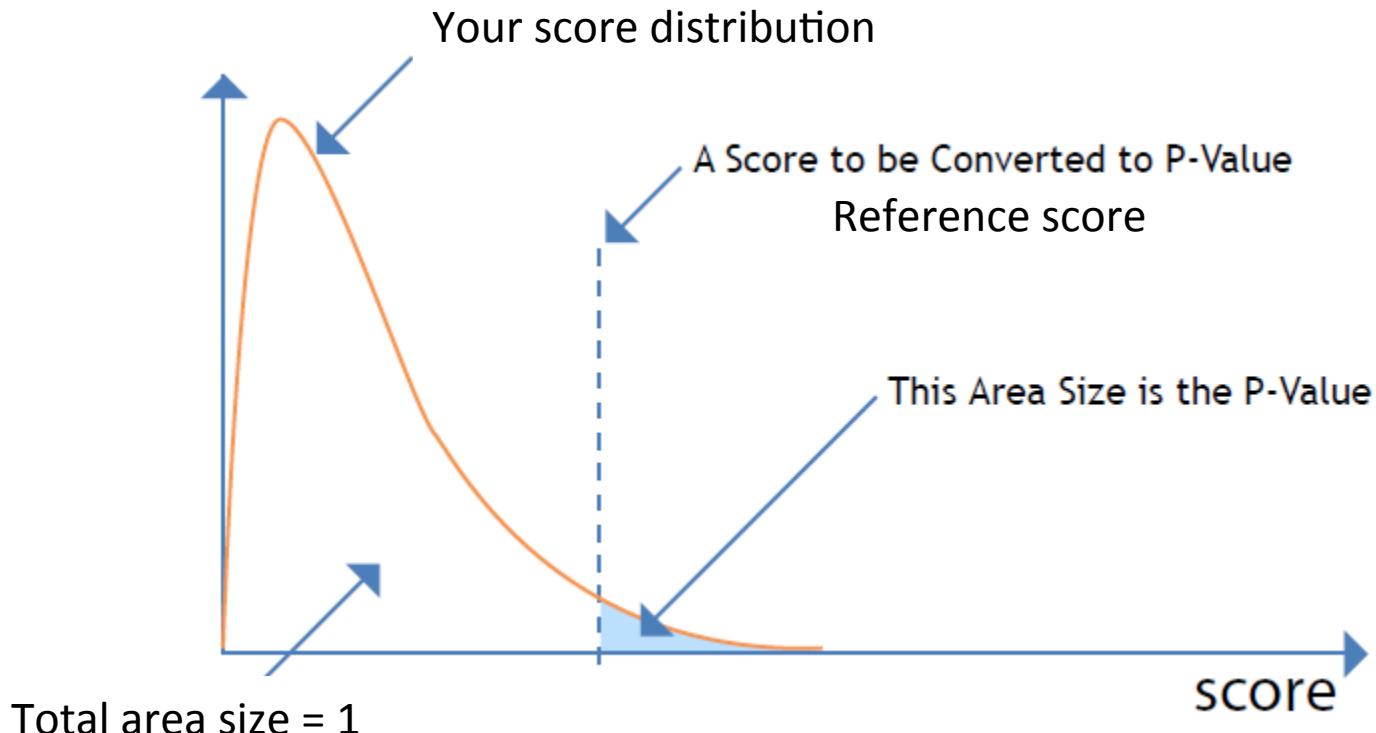
With *high* probability  $(1-\delta)$ ,  
 $0 < \delta \ll 1$

$$R[f] \leq R_{\text{train}}[f] + \varepsilon(\delta, C/N)$$

$R_{\text{gua}}[f]$

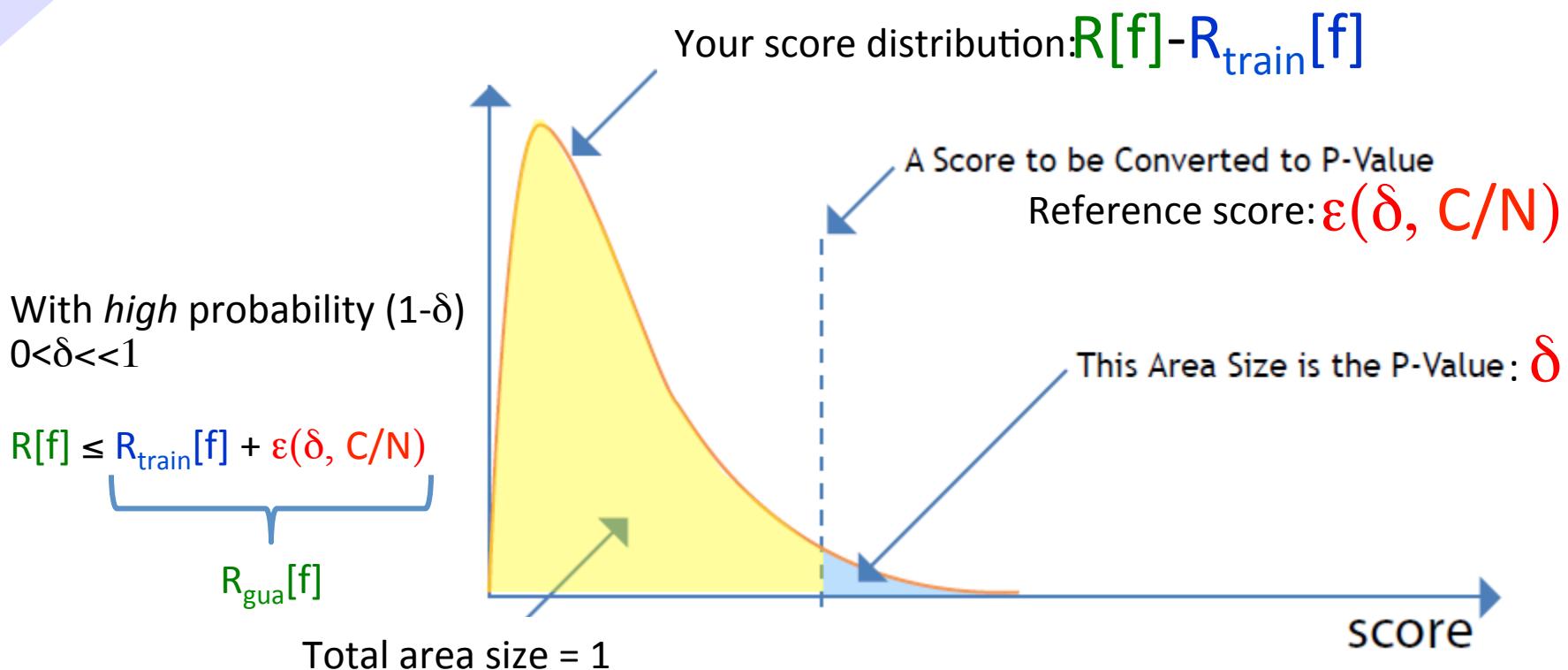
DATA ANALYSIS COMPLETE! www.datavizmentor.net

# The intuition: remember pvalues



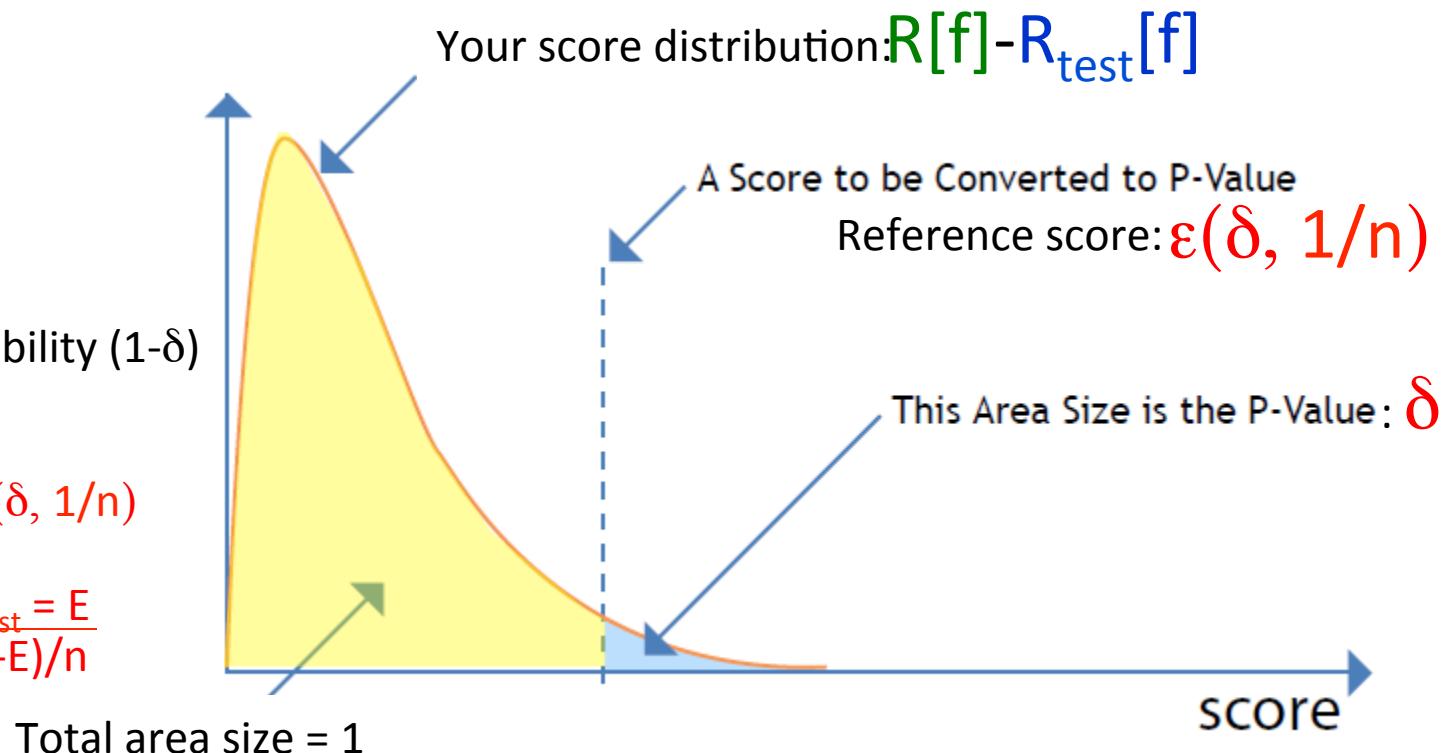
Pvalue = fraction of score values above the reference score.  
**SMALL pvalues are good:** they shed doubt on the “null hypothesis”.

# Guaranteed risk and pvalue



Pvalue = fraction of score values above the reference score.  
**SMALL pvalues are good:** they shed doubt on the “null hypothesis”.

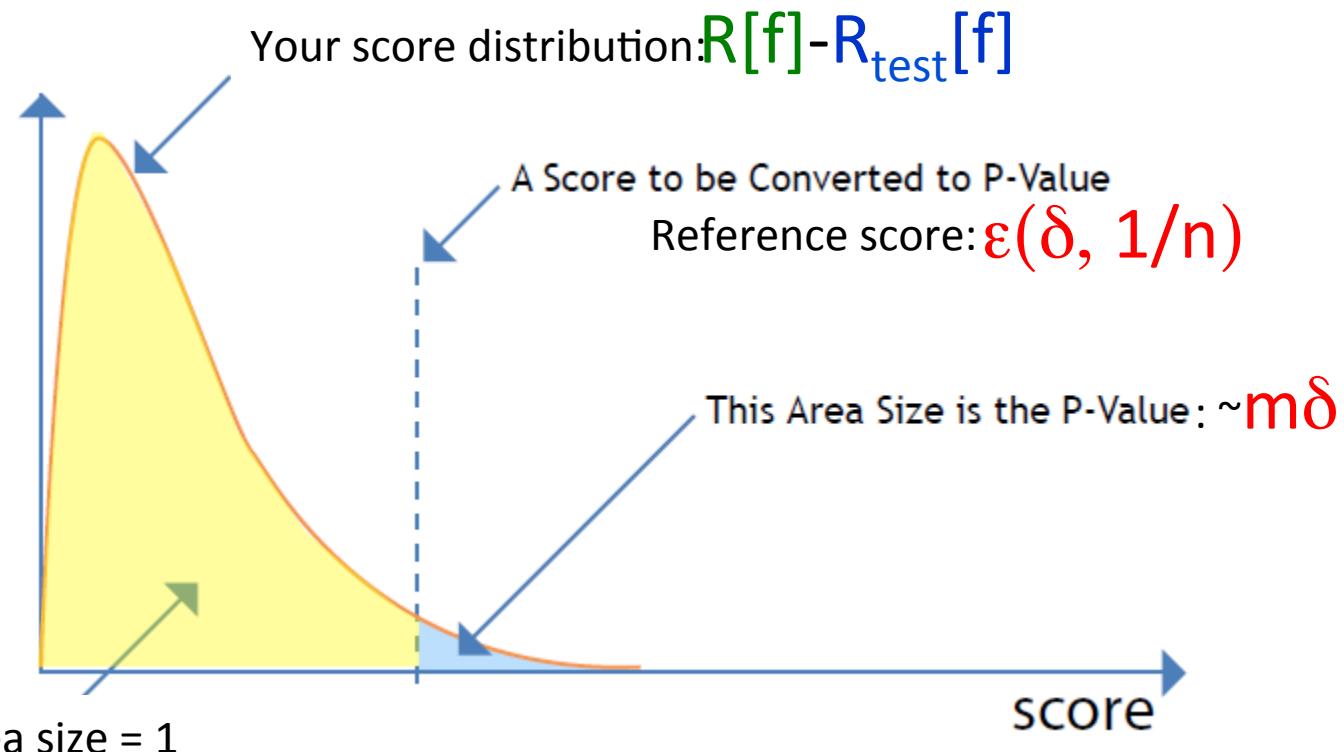
# Test set bounds



Pvalue = fraction of score values above the reference score.  
**SMALL pvalues are good:** they shed doubt on the “null hypothesis”.

# Multiple testing

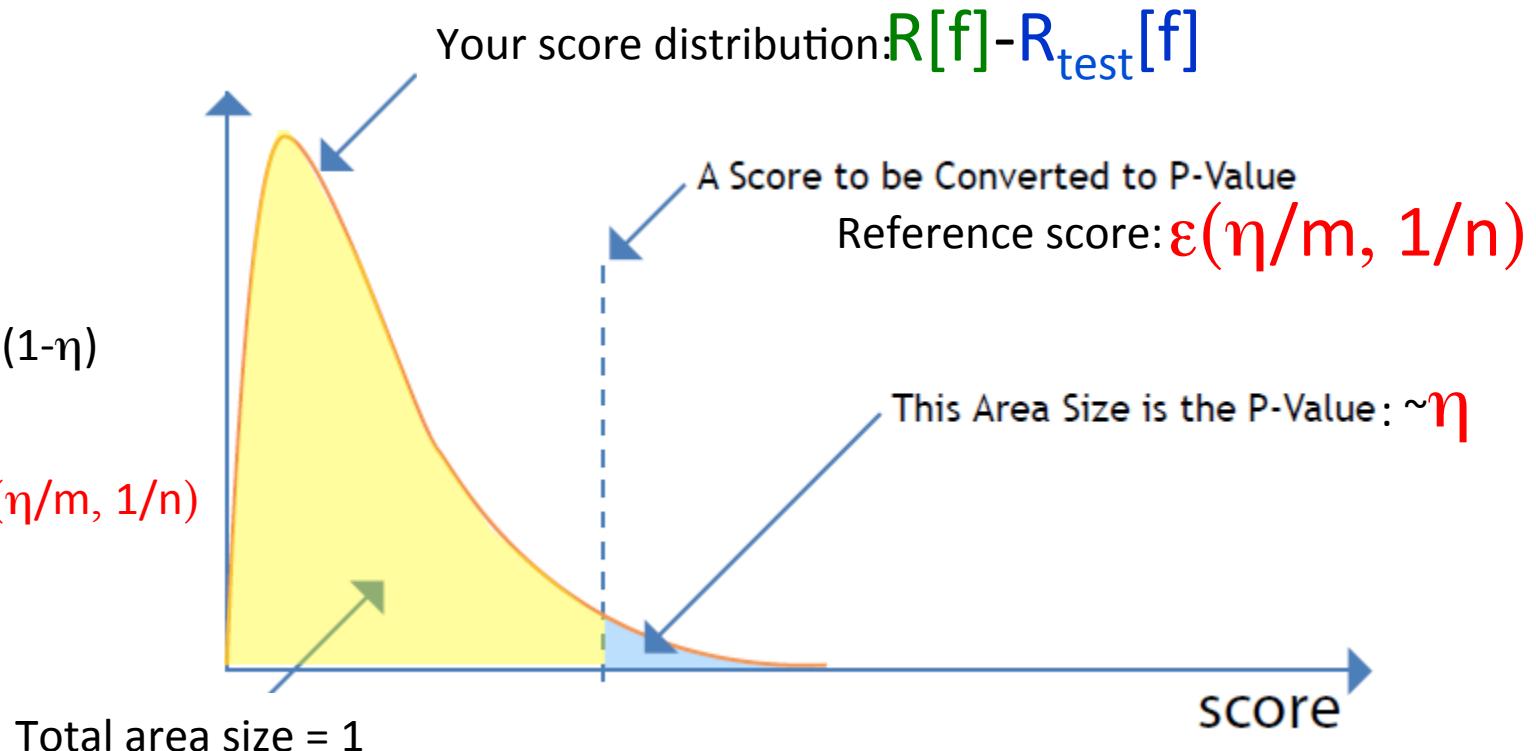
Bonferroni's correction:  $pvalue * m$



Pvalue = fraction of score values above the reference score.  
**SMALL pvalues are good:** they shed doubt on the “null hypothesis”.

# Multiple testing

## Variable change

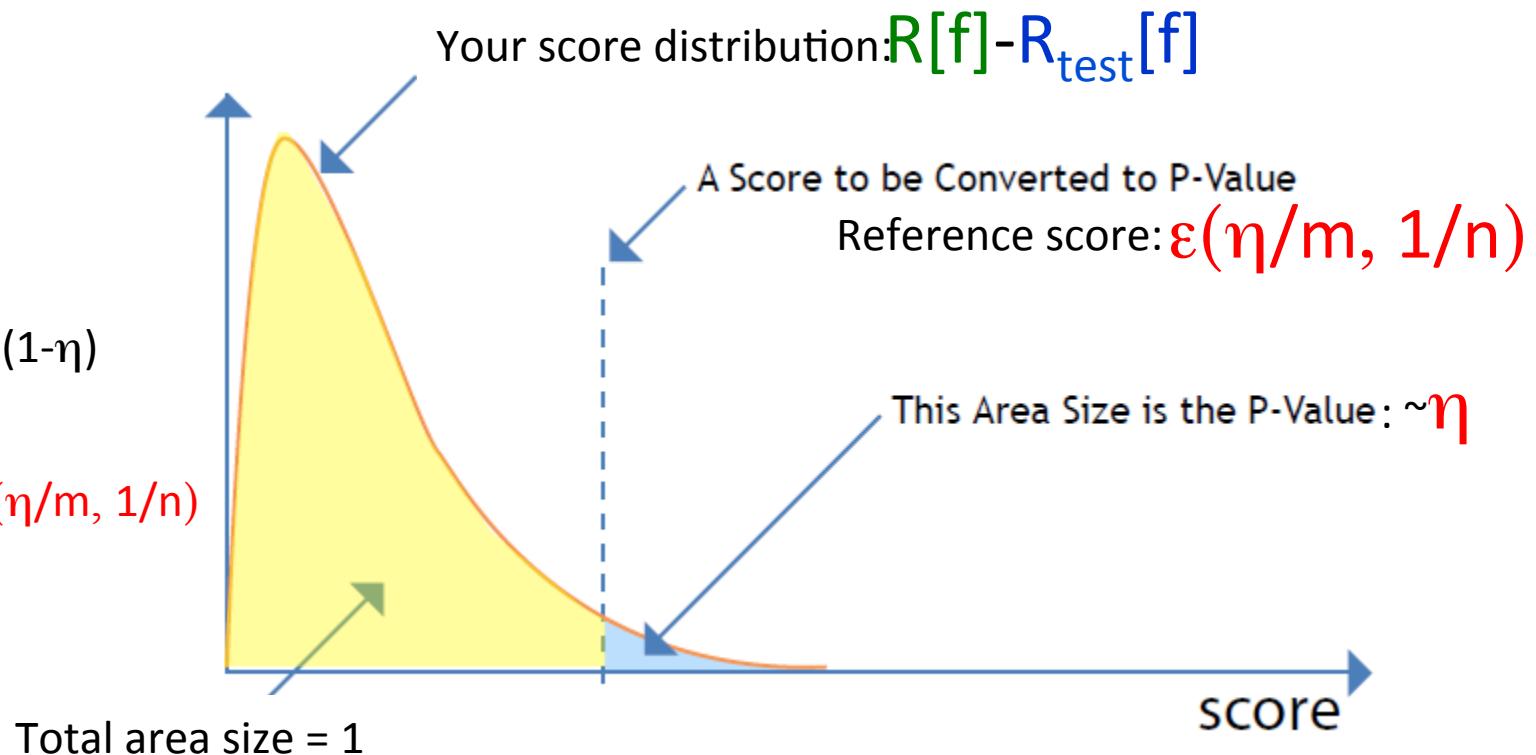


Pvalue = fraction of score values above the reference score.  
**SMALL pvalues are good:** they shed doubt on the “null hypothesis”.

In fact...

$\log(m)$  is a complexity/capacity

Vapnik 1974

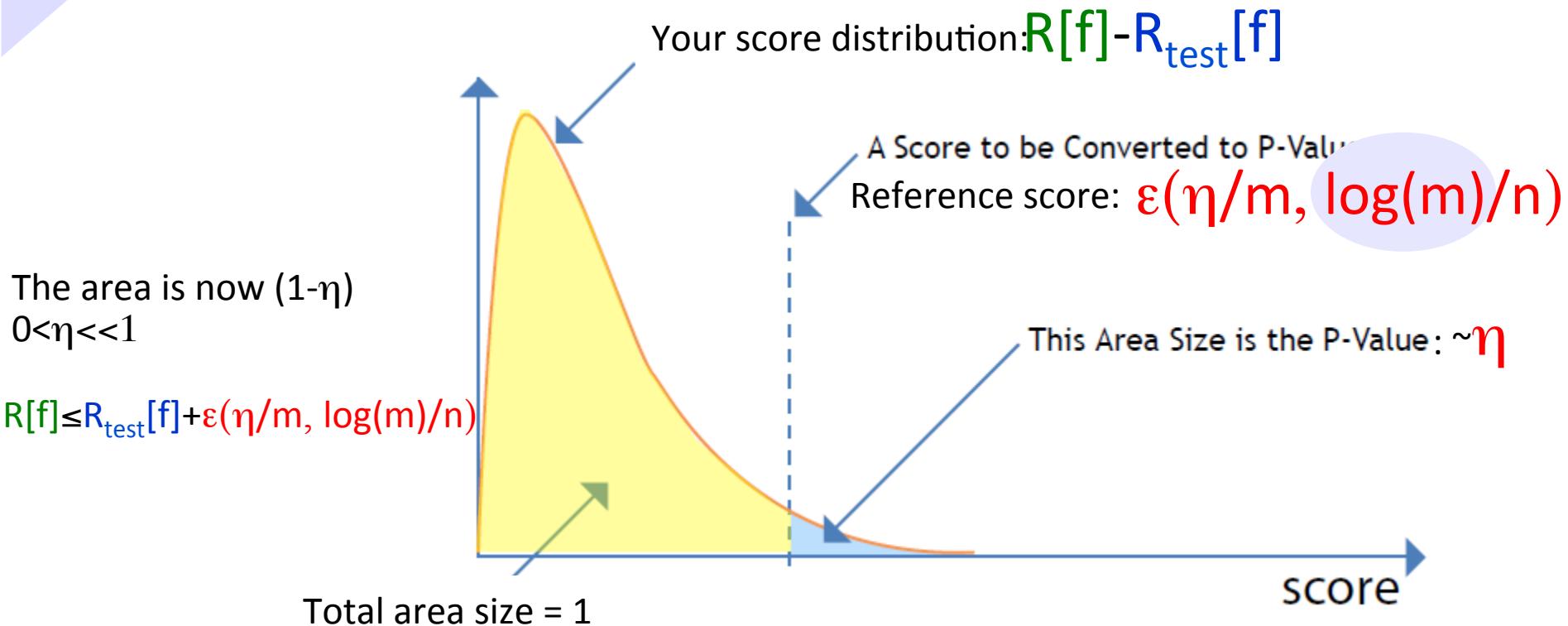


$$\epsilon(\eta/m, 1/n) = \sqrt{\log(2m/\eta)/2n}$$

In fact...

$\log(m)$  is a complexity/capacity

Vapnik 1974

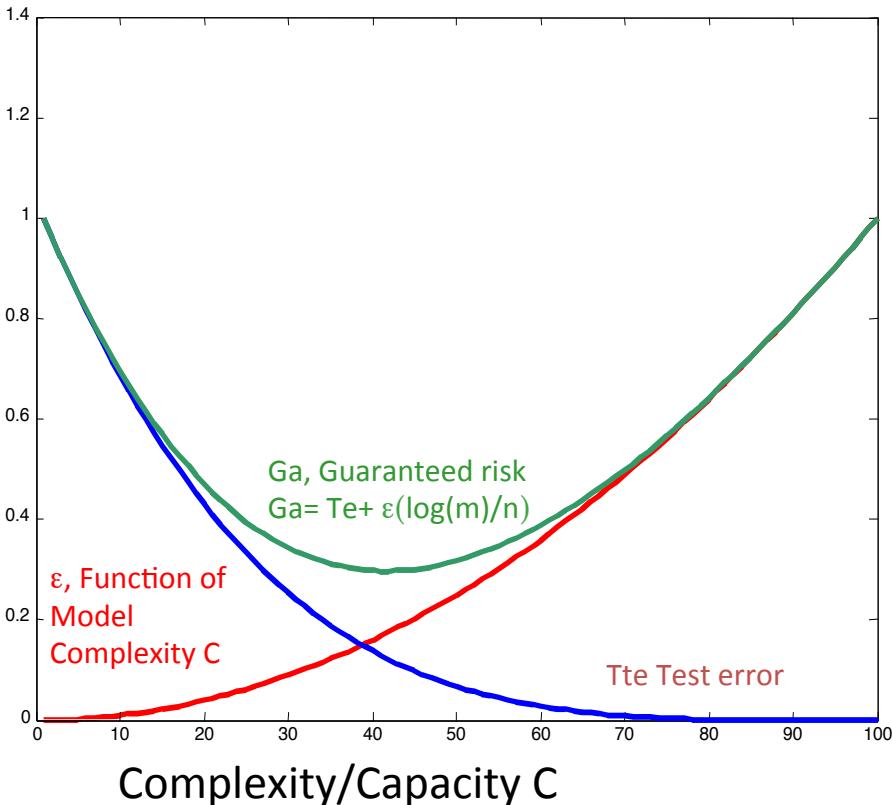


$$\epsilon(\eta/m, 1/n) = \sqrt{\log(2m/\eta)/2n}$$

# Capacity for “multiple testing”

Vapnik 1974

The “test” set became a training set!!!!



With *high* probability  $(1-\delta)$ ,  
 $0 < \delta \ll 1$

$$R[f] \leq R_{\text{test}}[f] + \varepsilon(\delta, \log(m)/N)$$

$C$

R<sub>gua</sub>[f]
C

# Capacity

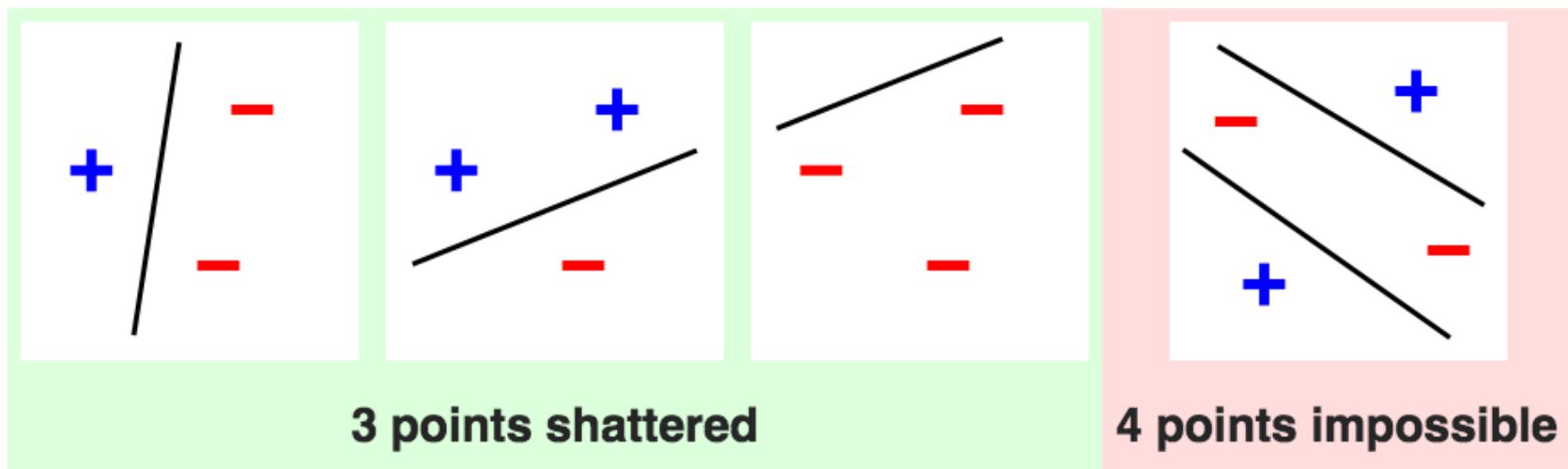
Vapnik-Chervonenkis 1971

- $C = \log(m)$  is the capacity for a FINITE number of hypotheses ( $m$  is number of models tried).
- What is  $C$  for a continuous hypothesis space (parameterized with  $w, \alpha, q, \sigma$ , etc.?)

# VC dimension

Vapnik-Chervonenkis 1971

$C = \text{number of training examples that can be separated, regardless of the assignment of labels.}$



Source: Wikipedia

For  $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$  (linear decision boundary/hyperplane)  
the VC dimension is  $C = d$  (the dimension of input space).

# VC dimension of Perceptrons and kernel methods

$f(\mathbf{x}) = \mathbf{w} \bullet \Phi(\mathbf{x})$ , the VC dimension is D  
(the dimension of FEATURE space).

$$\mathbf{w} = \sum_k \alpha_k \Phi(\mathbf{x}^k) \rightarrow \text{kernel trick}$$

$$f(\mathbf{x}) = \sum_k \alpha_k \Phi(\mathbf{x}^k) \bullet \Phi(\mathbf{x}) = \sum_k \alpha_k k(\mathbf{x}^k, \mathbf{x})$$

$$k(\mathbf{x}^k, \mathbf{x}) = \Phi(\mathbf{x}^k) \bullet \Phi(\mathbf{x})$$

If D is infinite, the VC dimension is infinite!

# VC dimension of kernel methods

Not quite that bad...

Remember the representer theorem:

$f^*(x) = \operatorname{argmin}_f R_{\text{train}}[f] + \lambda \Omega[f]$  admits a representation

$$f^*(x) = \sum_k \alpha_k k(x, x^k)$$

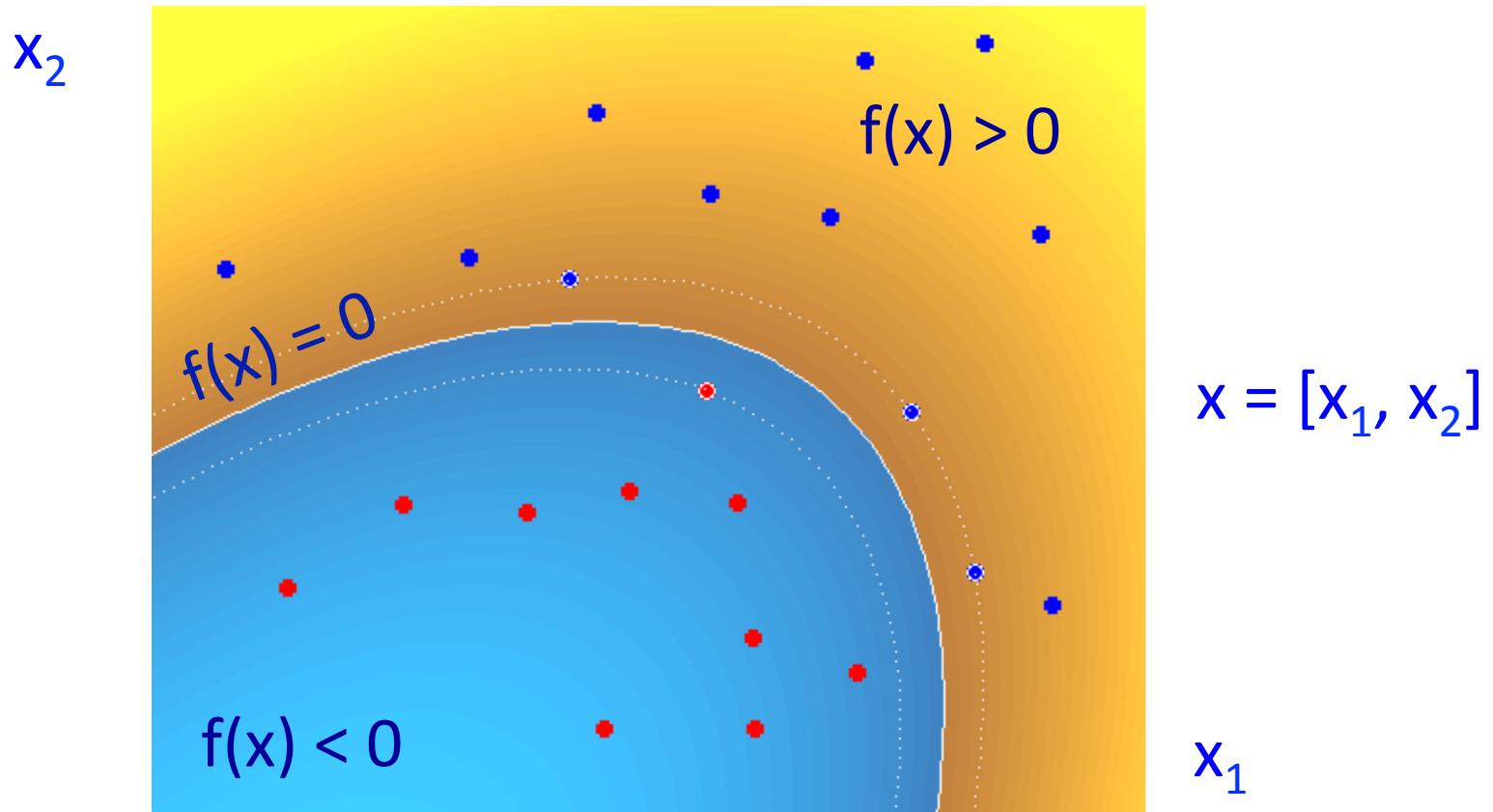
We can consider that  $k(\cdot, x^k)$  are features  $\Phi(\cdot)$  of a Perceptron with  $N$  parameters (the  $\alpha_k$ ). So:

$$C \leq N.$$

For support vector machines the sum runs only over the  $N_S$  support vectors, so:

$$C \leq N_S.$$

# SVM leave-one-out error: Bounded by $N_{SV}/N$



$$f(x) = \sum_k \alpha_k k(x_k, x) + b$$

*SVM, Boser-Guyon-Vapnik, 1992*

# But... the kernels have HP

$$f(\mathbf{x}) = \sum_k \alpha_k k(\mathbf{x}, \mathbf{x}^k; \theta)$$

They are IN FACT universal approximators:

$$C = \infty$$

if we also optimize  $\theta$  with the training set.

So we need to FIX  $\theta$  to enjoy a finite  $C$  when we use the training set. We can then find the best  $\theta$  with a validation set (a second level learning problem).

# For $\lambda$ we also need 2 levels

$$\min_f R_{\text{reg}}[f]$$

$$R_{\text{reg}}[f] = R_{\text{train}}[f] + \lambda \| \mathbf{w} \|^2$$

$$\partial R_{\text{reg}} / \partial \lambda = \| \mathbf{w} \|^2$$

At the optimum,  $\| \mathbf{w} \|^2 = 0$ . Not very helpful:  $C = 0$  if we also optimize  $\lambda$  with the training set.

So we need to FIX  $\lambda$  TOO to enjoy a finite  $C$  when we use the training set. We can then find the best  $\lambda$  with a validation set (another second level learning problem).

# Summary (so far)

- Heuristic HP search (wrapper-style) is computationally expensive, so we want to “push down” some HP to the learning machine.
- But optimizing both parameters and HP with the **training error** is not a good idea because of:
  - **Computational complexity**: we loose **convexity** in parameter space.
  - **Statistical complexity**: to train we get need **finite capacity**  $0 < C < \infty$ .
- Embedded methods solve a **bi-level optimization problem** by direct optimization of the HP by **cross-validation using knowledge about the learning machine to conduct a more efficient search**. But how? Grid search =>  $K * g^d$  learning machine to train. K=number of folds in K-fold, g = number of grid points, d the number of hyper-parameters.
- Idea #1: Approximate the leave-one-out error by “sensitivity analysis” to get a cheap  $R_{CV}$ . Gain: **replace K-fold by training ONCE**.
- Idea #2: Diagonalise  $X^T X$  or  $XX^T$  (or the kernel matrix): makes inversion “cheap” by varying  $\lambda$ . Gain: **replace training g machines by ONE**.

# Idea #1: Virtual leave-one-out

**Support Vector Machines:**

Refinements of  $N_{SV}/N$ .

**Linear regression:**

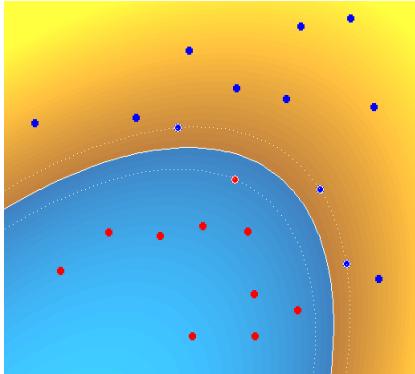
$$r_{LOO}^k = r_{TRAIN}^k / (1 - [XX^+]_{kk})$$

**Kernel methods including SVM, regression, logistic regression:**

Efficient approximate leave-one-out cross-validation for kernel logistic regression (Cawley-Talbot, 2008)

**Neural Networks:**

Local Linear Least Squares: Performing Leave-one-out Without Leaving Anything Out (Monari-Dreyfus, 2002)

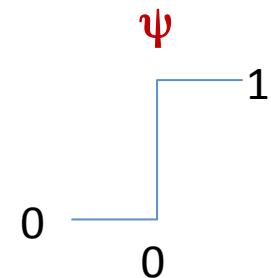


# Example 1: SVM (hard margin)

$$R_{\text{LOO}} = (1/N) \sum_{k=1:N} \Psi(-y_k f^{(-k)}(\mathbf{x}^k))$$

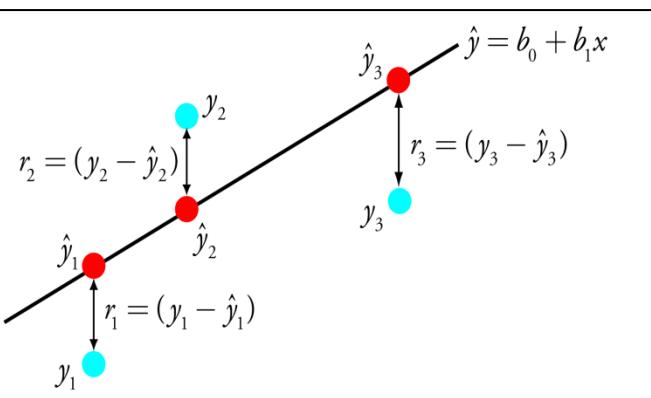
- LOO functional margin

$$R_{\text{LOO}} \leq (1/N) \sum_{k \in \{\text{SV}\}} \Psi(U^k - 1)$$



- $U^k$  is a bound on the change in functional margin for SV
- $$0 \leq y_k f^{(0)}(\mathbf{x}^k) - y_k f^{(-k)}(\mathbf{x}^k) \leq U^k$$
- Simplest bound  $U^k = \infty \Rightarrow R_{\text{LOO}} = N_{\text{SV}}/N.$
  - Jaakkola-Haussler (1999):  $U^k = \alpha_k k(\mathbf{x}^k, \mathbf{x}^k)$
  - Opper-Winther (2000):  $U^k = \alpha_k / [K_{\text{SV}}^{-1}]_{kk}$
  - Span bound Vapnik-Chapelle (2000):  $U^k = \alpha_k S_k^2$   
where  $S_k$  is the distance of  $\phi(\mathbf{x}^k)$  to the span of other SV.

## Example 2: Linear regression



$$R_{\text{LOO}} = (1/N) \sum_{k=1:N} (r_{\text{LOO}}^k)^2 \quad \text{MSE}$$

- $r_{\text{LOO}}^k = y_k - f^{(-k)}(\mathbf{x}^k)$  residual of example k for  $f^{(-k)}$
- $r_{\text{TRAIN}}^k = y_k - f(\mathbf{x}^k)$  residual of example k for  $f$
- Efron-Tibshirani, 1993

$$r_{\text{LOO}}^k = r_{\text{TRAIN}}^k / (1 - [\mathbf{X}\mathbf{X}^+]_{kk})$$

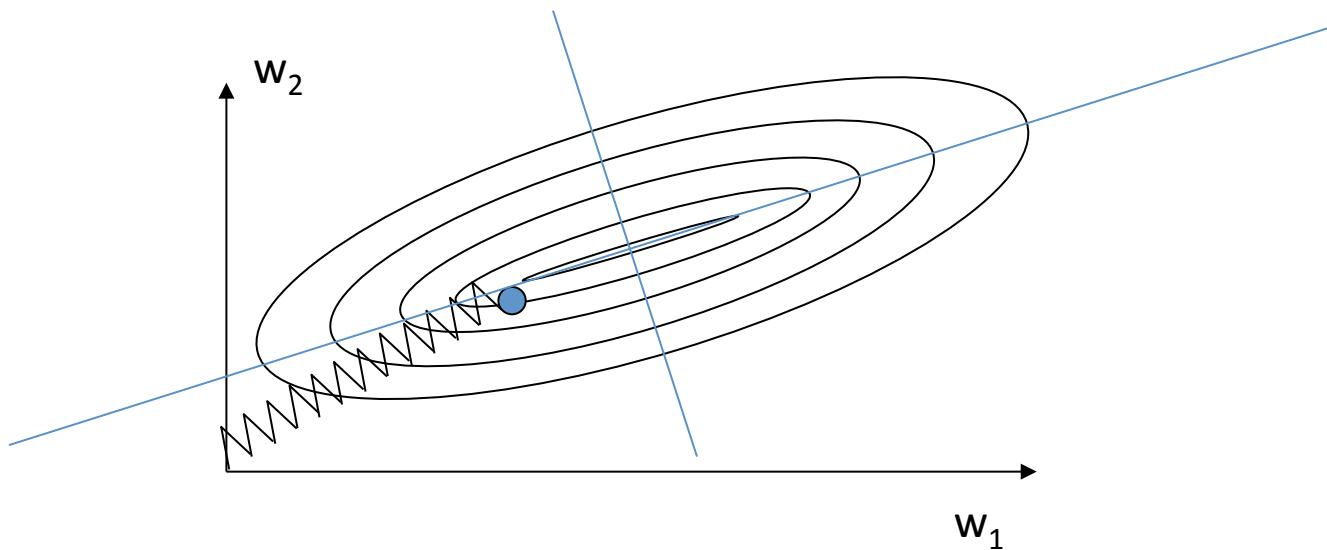
# Idea #2: Matrix factorization

- We need to compute:  $(XX^T + \lambda I)^{-1}$  or  $(X^TX + \lambda I)^{-1}$ , whichever is smallest:  $\dim(N, N)$  or  $(d, d)$ .  
So  $(K + \lambda I)^{-1}$  in general.
- We want to factorize a symmetric matrix  $K$  as:  
 $K = U^T D U$ , where  $D$  is a diagonal matrix  
(of eigenvalues) and  $U$  is an orthogonal matrix  
of eigenvectors:  $U^T U = I$  (and  $U^T(\lambda I)U = \lambda I$ )
- Then  $(K + \lambda I)^{-1} = (U^T D U + U^T(\lambda I)U)^{-1}$   
 $= (U^T(D + \lambda I)U)^{-1}$

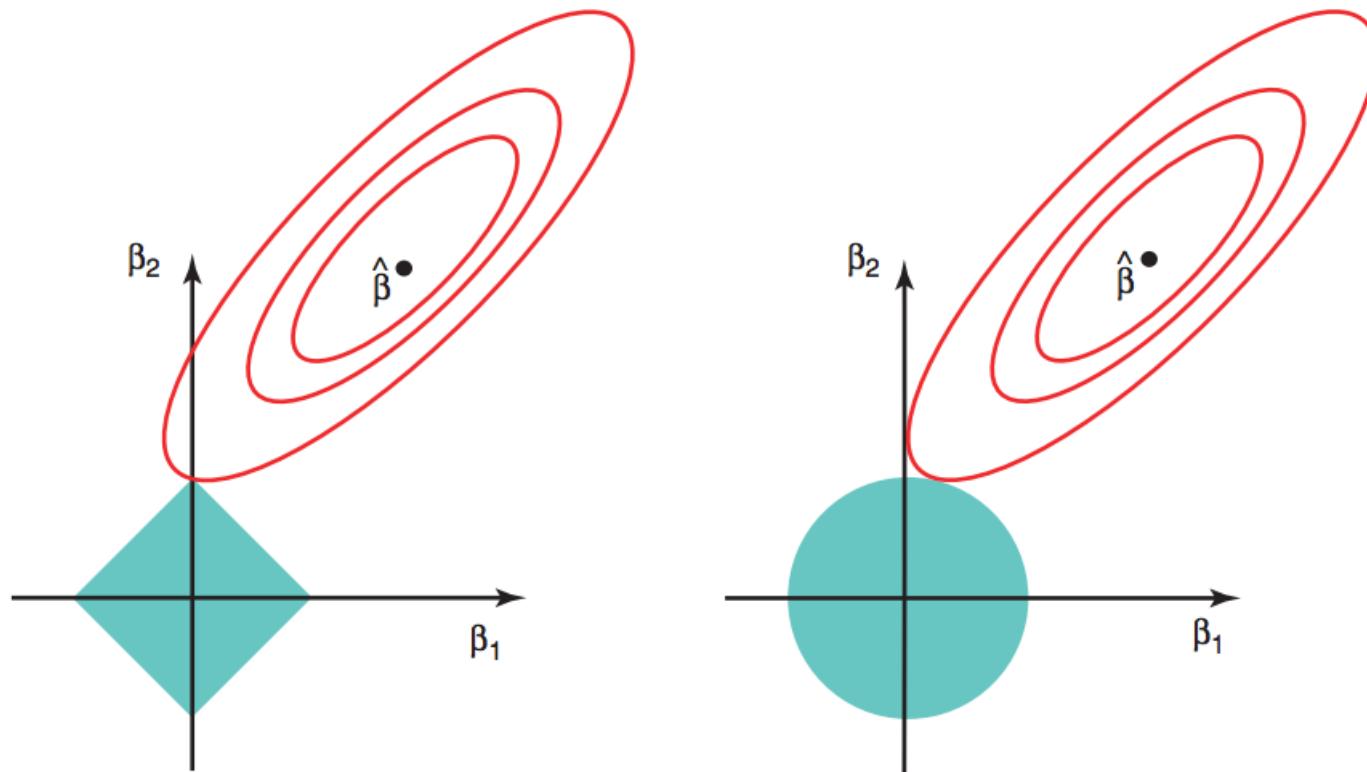
$$(K + \lambda I)^{-1} = U(D + \lambda I)^{-1} U^T$$

# How does this relate to dimensionality reduction?

- Reminder: For  $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$  the VC dimension is  $C = d$  (the dimension of input space).
- $R_{\text{reg}} = R_{\text{train}} + \lambda \|\mathbf{w}\|^2$

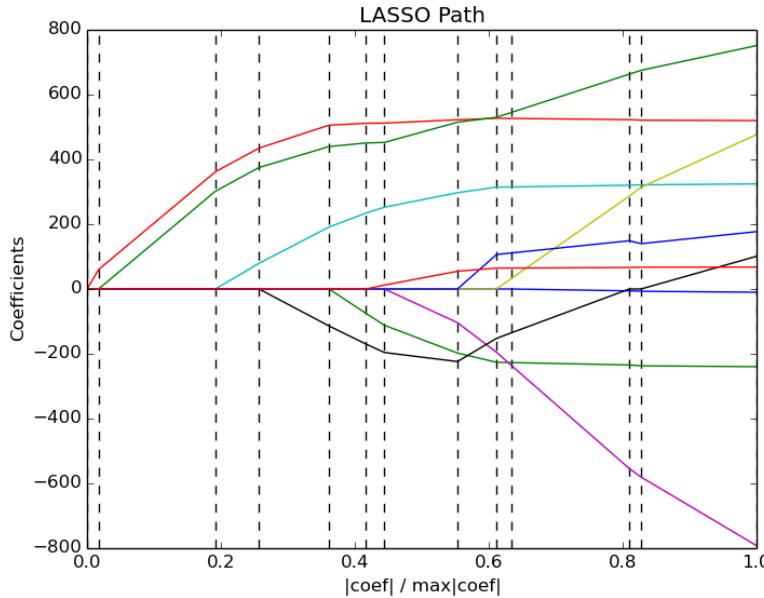


# Different regularizers



**FIGURE 6.7.** Contours of the error and constraint functions for the lasso (left) and ridge regression (right). The solid blue areas are the constraint regions,  $|\beta_1| + |\beta_2| \leq s$  and  $\beta_1^2 + \beta_2^2 \leq s$ , while the red ellipses are the contours of the RSS.

# Regularization path



See also:

[http://scikit-learn.org/stable/auto\\_examples/linear\\_model/plot\\_lasso\\_lars.html](http://scikit-learn.org/stable/auto_examples/linear_model/plot_lasso_lars.html)

- Hastie et al 2004:

<http://www.jmlr.org/papers/volume5/hastie04a/hastie04a.pdf>

- Efron 2004:

<http://statweb.stanford.edu/~tibs/ftp/lars.pdf>

- Friedman et al 2010:

<http://core.ac.uk/download/pdf/6287975.pdf>

# Summary

- Heuristic HP search (wrapper-style) is computationally expensive, so we want to “push down” some HP to the learning machine.
- But optimizing both parameters and HP with the training error is not a good idea because of:
  - **Computational complexity:** we loose convexity in parameter space.
  - **Statistical complexity:** we get infinite capacity.
- Embedded methods solve a **bi-level optimization problem** by direct optimization of the HP with the CV error, e.g. by gradient descent.
- Computational tricks for least square regression:
  - Leave-one-out error computed “cheaply” (**virtual LOO**).
  - Matrix factorization to vary the ridge value inexpensively.
- Shrinkage ( $\|\mathbf{w}\|_q$  regularizers) perform space dimensionality reduction.

Come to my office hours...  
Wed 2:30-4:30 Soda 329

**Next time:**  
**Gaussian classifier**

