

VC-DIMENSION

- Measure of complexity of the model (degree of freedom of the model)
(flexibility to fit data)
 - by VC-dimension as we can measure # of parameters of polynomial
 - We can use R_{emp} to approximate the R (general rule) \Rightarrow yes, considering
 - \Rightarrow VC-BOUND: holds with probability $1 - \delta$ that:
- $$R \leq R_{\text{emp}} + \epsilon \left(\frac{1}{l}, \text{VC}, \frac{1}{\delta} \right)$$

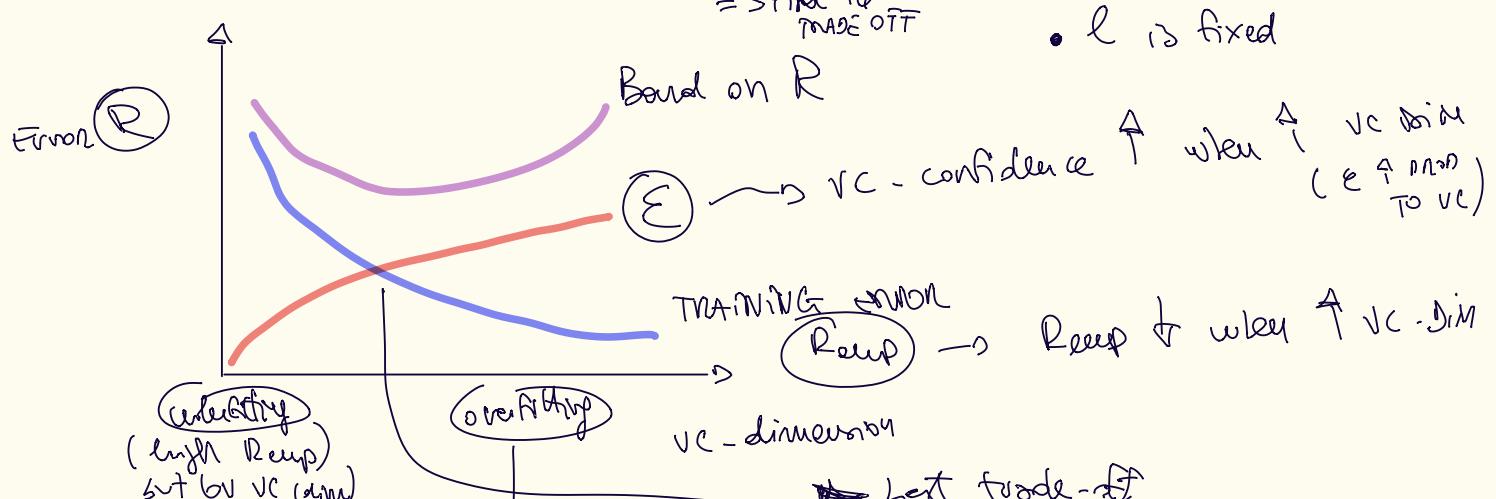
$\xrightarrow{\text{VC-bound}}$
- R_{emp}
- δ
- VC-confidence
- \Rightarrow it's a guaranteed risk: measure this quantity can measure the upper bound
 - \Rightarrow we never measure R \rightsquigarrow because you need to know the entire distribution of the data
 - ϵ is a function that grows with VC-dim
 $(\text{VC} \uparrow \Rightarrow \epsilon \uparrow)$
 $(l \uparrow; \delta \uparrow \Rightarrow \epsilon \downarrow)$
 - δ is the confidence
 \Rightarrow rules the bound holds
 (prob)
 \Rightarrow rule the confidence
 - R_{emp} decreases using complex models \Rightarrow high VC-dimension
 - if $\delta = 0,01 \Rightarrow$ bound hold with prob 0,99 \Rightarrow prob the bound is true
 - \Rightarrow NB: we can overestimate the generalization error R with prob $1 - \delta$
where ingredients are: R_{emp} (can measure)
VC-dim (you can control, e.g. # data)
 \Rightarrow correlation b/w R and R_{emp}
 - if # data can still have a high VC-dim but still have a good control
because $\uparrow l$, $\uparrow \text{VC} \Rightarrow$ w.t. $\frac{1}{l} \cdot \text{VC}$
 \Rightarrow increase $l \Rightarrow$ increase VC
 $\Rightarrow \uparrow l$
 \Rightarrow higher VC
 - if for 1 point, take it 100 times
 \Rightarrow the bound \Rightarrow P of taking 100 times (randomly) the same point is very low and belongs to $1 - \delta$ zone
 - NB: fixing δ we can explain underfitting/overfitting and how control them (when they happens on theoretical POV) (to cut-off steamp cases)

\Rightarrow higher l (# data) \Rightarrow lower VC confidence and band closer to R

- too simple model (low VC-dimension) can be not sufficient
 \Rightarrow due to high R_{emp} \Rightarrow UNDERFITTING

- higher VC-dimension (fix l) \Rightarrow lower R_{emp}
 \Rightarrow but VC-confidence & R may increase
 \Rightarrow OVERFITTING

- STRUCTURAL RISK MINIMIZATION: minimize the bound
 \Rightarrow find the trade-off



Blue + Red: take the sum \Rightarrow purple \Rightarrow U-shape

↓
low R_{emp}

VC-dim ↑, VC-confidence ↑, red line ↑ \Rightarrow overfitting

$$R_{emp} = \frac{1}{l} \sum_{i=1}^l L(h(x_i), y_i)$$

LEARNING ALGO 4 Reg./Class. using linear model (Based on LMS) CLASSIFI.

Learning problem: Find w s.t. minimize the loss on Training data

$$R_{\text{emp}} = \frac{1}{l} \sum_{p=1}^l L(h_w(x_p), y_p)$$

ASSUME USE LEAST SQUARE (or M REGRESSION)

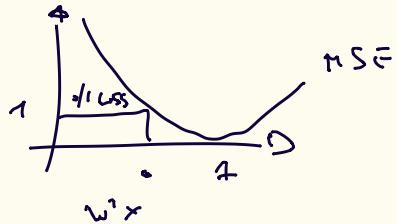
o) Replace original objective function

$L(0/1 \text{ loss})$ by a smooth differentiable one \Rightarrow MSE (loss)

$$\cdot h(x) = 1 \text{ if } w^T x > 0$$

minimum error is for $w^T x > 0$

\Rightarrow hence \rightarrow classification error



MATRIX OF DATA

$$\text{Same as regression } E(w) = \sum_{p=1}^l (y_p - x_p^T w)^2 = \|y - Xw\|^2$$

(find w to minimize sum of square) \rightarrow if target is positive, we are going to $w^T x \approx 1$

here we CANNOT use $h(x)$ instead

$$\text{because } h(x) = \text{sign}(w^T x + w_0)$$

$\Rightarrow \Rightarrow$ Discrete function

WITH a THRESHOLD (THIS IS NOT DIFFERENTIABLE)

\Rightarrow we DO the derivative for $w = 0, \dots, m$

(we just w_0 and w_1)

imposing gradient \Rightarrow

$$\begin{aligned} \frac{\partial E(w)}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_{p=1}^l (y_p - x_p^T w)^2 = \sum_{p=1}^l 2(y_p - x_p^T w) \cdot \frac{\partial (y_p - x_p^T w)}{\partial w_j} = \\ &= \sum_{p=1}^l 2(y_p - x_p^T w) \left(\frac{\partial y_p}{\partial w_j} - \frac{\partial (x_p)_1 w_1}{\partial w_j} - \frac{\partial (x_p)_2 w_2}{\partial w_j} - \dots - \frac{\partial (x_p)_j w_j}{\partial w_j} - \dots \right) \\ &= -2 \sum_{p=1}^l (y_p - x_p^T w) (x_p)_j = -\sum_{p=1}^l 2 y_p x_{pj} \end{aligned}$$

Lossy ALGs For Classification: $E(w) = \sum_{p=1}^l (y_p - x_p^\top w)^2$

$$\frac{\partial E(w)}{\partial w_j} = 2 \sum_{p=1}^l (y_p - x_p^\top w) \cdot \underbrace{\frac{\partial (y_p - x_p^\top w)}{\partial w_j}}_{=}, \quad -(x_p)_j$$

$$= 2 \sum_{p=1}^l (y_p - x_p^\top w) \cdot \left(\frac{\partial y_p}{\partial w_j} - \cancel{\frac{\partial (x_p)_1}{\partial w_j} w_1} - \cancel{\frac{\partial (x_p)_2}{\partial w_j} w_2} - \dots - \cancel{\frac{\partial (x_p)_j}{\partial w_j} w_j} \right)$$

$$= -2 \sum_{p=1}^l (y_p - x_p^\top w) (x_p)_j = \boxed{- \sum_{p=1}^l 2 y_p x_{p,j}}$$

\Rightarrow USE NOW **GRADIENT DESCEND**

- la derivazione di prima conduce un ALGORITMO ITERATIVO

$$\frac{\partial E(w)}{\partial w_j} = -2 \sum_{p=1}^l y_p x_{p,j} = -2 \sum_{p=1}^l (y_p - x_p^\top w) (x_p)_j$$

- GRADIENT = ASCENT DIRECTION: we can move verso il minimo w^n in GRADIENT DESCEND = $-\frac{\partial E(w)}{\partial w} = \Delta_w$

\Rightarrow LOCAL SEARCH:

- inizializza il vettore dei pesi
- itero modificandolo fino a dessegnare per minimizzazione

- GRADIENT vector (for each component)

$$\Delta_w = -\frac{\partial E(w)}{\partial w} = \begin{bmatrix} -\frac{\partial E(w)}{\partial w_1} \\ -\frac{\partial E(w)}{\partial w_2} \\ \dots \\ -\frac{\partial E(w)}{\partial w_j} \end{bmatrix} = \begin{bmatrix} \Delta w_1 \\ \Delta w_2 \\ \dots \\ \Delta w_m \end{bmatrix}$$

la error function
 in n multidimensional space w ,
 is called
 (we not visualize)

- AS iterative approach: we LEARNING rule on DFTM
 changing w proportionally to the (opposite) of GRADIENT
- $$w_{new} = w + \text{eta} \cdot \Delta_w$$
- learning rate

A simple algorithm:

- 1) Start with weight vector w_{initial} (small), fix η ($0 < \eta < 1$).
gradient direction in which the function decreases + rapidly decreases
 - 2) Compute $\Delta w = \text{"gradient of } E(w) = -\frac{\partial E(w)}{\partial w}$ (or for each w_j)
learning rate
 - 3) Compute $w_{\text{new}} = w + \eta * \Delta w$ (or for each w_j)
affine transformation of previous
direction opposite of gradient
- Repeat (2) until convergence or $E(w)$ is "sufficiently small" (①)
stop: le errore è basso || raggiunto le massime iterazioni.
- $\Delta w/l$: least mean squares (dividing by l , that will be the standard case)
epoch error and Δw_j
 $\Rightarrow \text{error} = \text{sum of the error for all the patterns}$
 - Batch versions (Δw after each "epoch" of l training patterns)
 - η : step size = learning rate: speed/stability trade-off: can be (gradually) decreased to zero (guarantee convergence, avoiding oscillation around the min.): many variants will be introduced later

Micheli SE FSA TROPPO GRANDE \Rightarrow si oscilla intorno al minimo globale
TROPPO PICCOLO \Rightarrow converge lentamente o bloccato in un minimo locale 46

• BATCH: gradient is the sum over all the l patterns

$$\frac{\partial E_w}{\partial w_j} = -2 \sum_{p=1}^l (y_p - x_p^T w) x_{pj}$$

ultimate weight after the sum

• ONLINE/STOCHASTIC

- update weights with the error computed for each pattern

$$\frac{\partial E_p(w)}{\partial w_j} = -2(y_p - x_p^T w) x_{pj} = -\Delta_p w_j$$

• GRADIENT DESCEND CORRECTED TO

ERROR CORRECTION RULE

DELTAW RULE

- change w_j proportionally to the error $(y_{\text{target}} - \text{output})$
- $(\text{target } y - \text{output}) = 0 \Rightarrow \text{NO ERROR} \Rightarrow \text{NO CHANGES}$

- if input is positive ($\text{input}_j > 0$) and error $> 0 \Rightarrow$ output is too low (Δ positive)
 \Rightarrow increase w_j
 \Rightarrow \uparrow output
 \Rightarrow less error

- if input is < 0 and error $< 0 \Rightarrow$ output is too high (Δ negative)
 \Rightarrow decrease w_j
 \Rightarrow reduce output
 \Rightarrow less error

$$\Delta w_j = 2 \sum_{p=1}^l (x_p)_j (y_p - x_p^T w)$$

$$w_{\text{new}} = w + \eta * \Delta w$$

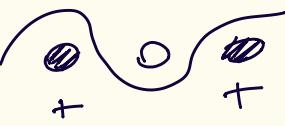
ACCURACY = media of correct classified $(l - \text{num_err})$

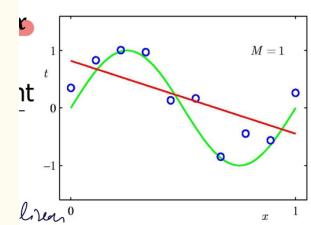
$$L(h(x_p), d_p) = \begin{cases} 0 & \text{if } h(x_p) = d_p \\ 1 & \text{else} \end{cases}$$

$$\text{mean-error} = \frac{1}{l} \left[\sum_{p=1}^l L(h(x_p), d_p) \right]$$

INDUCTIVE BIAS

- Language bias: the H is a set of linear functions (too rigid and restrictive)
- SEARCH BIAS: ordered the search algorithm least square minimization goal

es.  \Rightarrow need a non-linear model



\Rightarrow COME ESTENDERE UN MODELLO LINEARE

• in $h_w(x) = w_1 x + w_0$ or $h_w(x) = w^T x$

(**linea**) is the **way** in which the **coefficient w occurs** in the regression **equation** (not red line)

\Rightarrow we can use **TRANSFORMED INPUTS**, x, x^2, x^3, x^4 and other **oblique change** CON RELAZIONI NON LINEARI tra input e output

with LMS:

$$h_w(x) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_M x^M$$

\Rightarrow POLYNOMIAL REGRESSION

\Rightarrow call this **LINEAR BASIS EXPANSION**

$$h_w(x) = \sum_{k=0}^K w_k \phi_k(x)$$

- aggiungere tra input vector con l'aggiunta di variabili che sono la trasformazione dell'utile x in accordo delle funzione phi ($\phi_k: \mathbb{R}^m \rightarrow \mathbb{R}$)

es. polynomial representation: $x: \phi(x) = x_j^2$

or
 $\phi(x) = x_j x_i$

non-lineare trasformazioni of simple inputs: $\phi(x) = \log(x_j)$ $\phi(x) = \sqrt{x_j}$

• Riscriui il modello lineare come somma di trasformazioni del vettore originale x tramite la funzione phi con ogni componente perata w_k

• NUMERO DI PARAMETRI $K > m$ (before was n)

• il modello è lineare nei parametri (anche in phi, non in x)

• phi are fixed before see TR

• curse of dim

(P.S.): Posso misurare relazioni più complicate (non lineari) \Leftrightarrow with large basis of function like overfitting \Rightarrow need control complexity

=> CONTROL COMPLEXITY

Lo polinomio con grado ALTO

1) Ridge / Tikhonov REGULARIZATION:

smoothed model

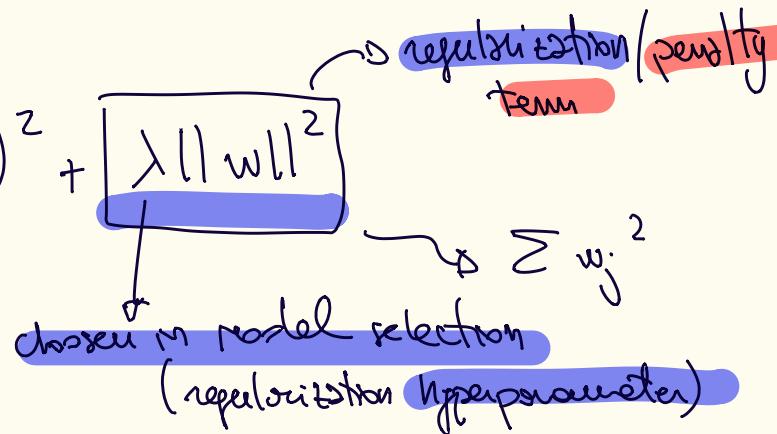
- aggiunge dei vincoli alla somma dei valori $|w_j|$ penalizzando modelli con alti valori di $|w|$ => less complex model

every data term (MSE)

$$\text{loss}(w) = \sum_{p=1}^l (y_p - x_p^T w)^2$$

↓
data term

control the size of the weight



→ qui la objective function si chiama loss (usata per stimare la funzione lost per ellinare il modello)

- perché $E(\text{errore})$ è utile

- per CALCOLARE il modello e usare nel DATA TERM Dopo $L(w)$

GRADIENT APPROACH:

$$\Rightarrow \text{new rule } w_{\text{new}} = w + \text{eta} \cdot \Delta w \quad [-2\lambda w]$$

\Rightarrow WEIGHT DECAY: add $2\lambda w$ al gradiente

$$\text{loss}(w) = \sum_{p=1}^l (y_p - x_p^T w)^2 + \lambda ||w||^2$$

$$\Rightarrow \text{gradient} \Rightarrow w_{\text{new}} = w + \text{eta} \cdot \Delta w - 2\lambda w$$

$\lambda = 0$ Good fitting
 \Rightarrow Min Regress (overfitting)
 \Rightarrow increase VC-dim

Tikhonov regularization: trade-off



$$\text{Loss}(w) = \sum_{p=1}^l (y_p - x_p^T w)^2 + \lambda ||w||^2$$

- Note the balancing (trade-off) between the two terms:

- Small **lambda** value → minimizing the loss the focus is on obtaining a small error data term (first term, minimize just the training error) with a too complex model (high norm of the weights), the risk is of overfitting.
- High **lambda** → minimizing the loss the focus is on the second term, hence the data error (first term) could grow too much, i.e. moving to underfitting
- The trade-off is ruled by the value of **lambda**
- We will see soon some examples

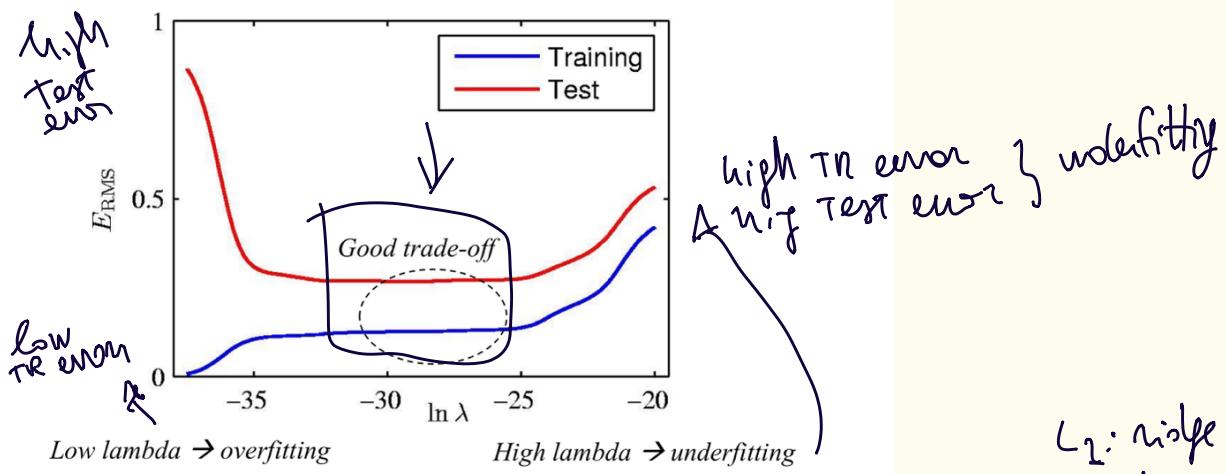
& quando le Gradienza non è un angolari bias (control hypothesis space) BUT SEARCH bias perché e nelle loss

Tikhonov and SLT

- The **penalty term** penalizes high value of the weights and tends to drive all the weights to smaller values
 - E.g. Some weights values can go even to zero
- It **implements a control** of the model complexity
- This leads to a model with **less (or proper) VC-Dim** (with a trade-off obtained through **just a (1) parameter** that you can control: the λ)

→ mi serve lambda e non devo decidere il grado del polinomio in valutazione perché è più generale e posso applicarlo con solo col polinomio

⇒ con 1 parametro controllo la complessità



L_2 : ridge regression
 L_1 : Lasso

The L_2 norm penalizes the square value of the weight and tends to drive all the weights to smaller values.

On the other hand, the L_1 norm penalizes the absolute value of the weights and tends to drive some weights to exactly zero (while allowing some weights to be large) → toward feature selection!

- Unfortunately $\|\cdot\|_1$ (absolute value) introduce a non differentiable loss → needs other approaches

(KNN): K-nearest neighbor

DATA (before EACER)

- STORE TR data
- wait until a data point is presented
- construct ~~ad-hoc~~ hypothesis to classify that data point

ALGORITHM:

- 1) Store TR data $\langle x_p, y_p \rangle \quad p = 1, \dots, l$
- 2) given an input x with dimension m
- 3) Find the nearest training example x_i
 \Rightarrow find i s.t. have minimum $d(x, x_p)$
 $\Rightarrow i(x) = \arg \min_p d(x, x_p)$
 \Rightarrow euclidean distance

$$d(x, x_p) = \sqrt{\sum_{t=1}^m (x_t - x_{pt})^2} = \|x - x_p\|$$
- 4) then output x

with $K=1$ NN

Can some limitation of linear model

- 1-nn return + for x_q
- 5-nn return - for x_q

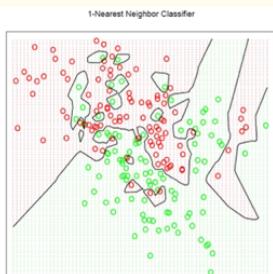


Figure 2.3:
The same classification example in two dimensions as in Figure 2.1.
GREEN = 0, RED = 1

- Very flexible !
- No misclassifications in TR data: 0 training error: what for test data ?
- Decision boundaries is *not* linear: it is quite *irregular*
- May be unnecessary noisy (e.g. for scenario 1)

K-Nearest Neighbors

- A natural way to classify a new point is to have a look at its neighbors,
- and take an average:

$$\text{avg}_k(\mathbf{x}) = \frac{1}{k} \sum_{x_i \in N_k(\mathbf{x})} y_i$$

- where $N_k(\mathbf{x})$ is a neighborhood of \mathbf{x} that contains exactly k neighbors (closest patterns according to distance d):
- k-nearest neighborhood: **K-nn**.

- If there is a clear dominance of one of the classes in the neighborhood of an observation \mathbf{x} , then it is likely that the observation itself would belong to that class, too. Thus the classification rule is the **majority voting** among the members of $N_k(\mathbf{x})$. As before,

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \text{avg}_k(\mathbf{x}) > 0.5 \\ 0 & \text{otherwise} \end{cases} \quad \text{for targets } y \text{ in } \{0,1\}$$

- For regression task: use directly the avg: mean over K-nn

look wise of dim.

NN

ARTIFICIAL NEURON: PROCESS UNIT

- node, neuron, unit

- Input: external source or other units

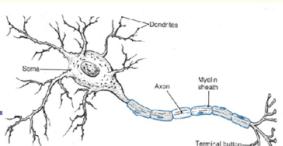
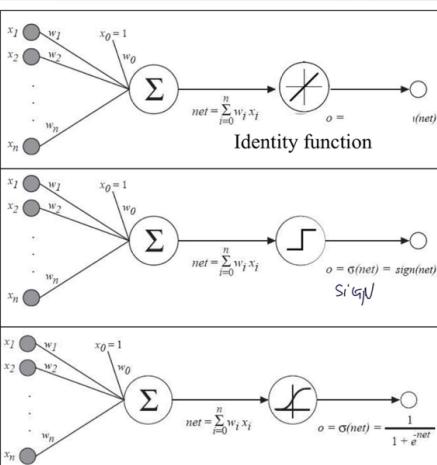
- Input connections: weights w : free-parameters

$$\left\{ \begin{array}{l} \text{net}_i(\mathbf{x}) = \sum_j w_{ij} x_j \quad \text{weighted sum} \\ o_i(\mathbf{x}) = f(\text{net}_i(\mathbf{x})) \end{array} \right. \quad \begin{array}{l} \text{is called} \\ \text{not input} \\ \text{to unit } i \end{array}$$

(w_{ij}) is the weight of the unit i coming from j unit to i

(f) is the activation function: linear, LTU, ...

Neuron: Three activation functions



$$h(\mathbf{x}) = \sum_i w_i x_i$$

Linear activation function

\Rightarrow classification

Perceptron (see LTU) \circ

Threshold activation function

Logistic function: sigmoid
(see later)

$$\text{avg}_k(\mathbf{x}) = \frac{1}{K} \sum_{x_i \in N_k(\mathbf{x})} y_i$$

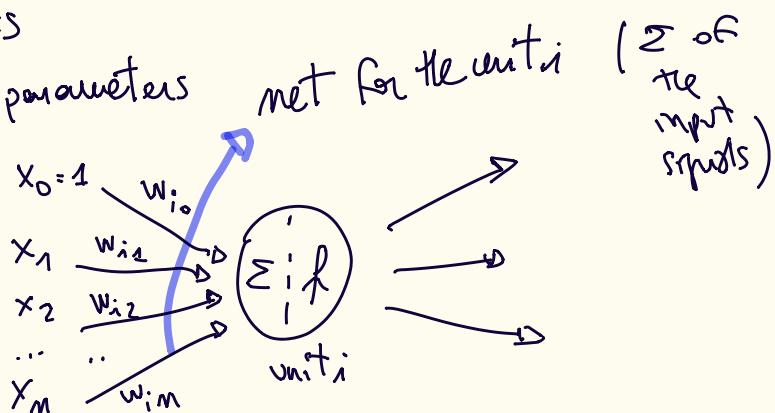
↑
neighborhood

KNN FOR MULTI-CLASS (not important)

$$h(\mathbf{x}) = \arg \max_v \sum_{x_i \in N_k(\mathbf{x})} 1_{v, y_i}$$

$$1_{v, y_i} = \begin{cases} 1 & \text{if } v = y_i \\ 0 & \text{else} \end{cases}$$

- count cbns of neighborhood



PERCEPTRON: • biologically inspired model

• with a simple computational unit

• can be composed and connected to build networks:
from LTU to MLP

- Neuron i get active when \sum connectors w_{ij} coming from unit j connected to it are active, plus a bias, is > 0
- BINARY CLASSIFICATION task

1 perceptron cannot solve the Xor \Rightarrow MLP for: with 1 hidden layer

- Useful concept of internal (re)representation of input variables via internal (hidden) units
- Developing high level (hidden) features is a key factor in NN
- The representation in the hidden layer makes easier the task to the output layer (last unit).
 - E.g. See the figure on the right: point $(1,0)$ (in the sud-est corner) became $h_1=0 h_2=1$ (north-west corner) \rightarrow now it is a linear separable problem

A look ahead:

- Such composition of sub-/intermediate operations to perform a complex task can be extended through many layers of abstraction
- In NN such internal representation/intermediate features can be learned
- Learning internal distributed representation \rightarrow representation learning and deep learning concepts (link to next parts of the ML course)

23

problems
becomes
linearly separable
with 1 hidden layer

Re-represent the input space in a way that makes easier the task

• 1 hidden layer to implement all possible function (4 levels just not)

LEARNING FOR 1 UNIT MODEL:
due to single perceptron linear sep. problem \rightarrow solved

1) Adaline: adaptive linear neuron (what I know): correspond to LTS and LMS for classification

- gradient descent solution
- an approach that we will generalize to MLP
- repetition
- classification

2) PERCEPTRON:

- non-linear model: w includes THRESHOLD ACTIVATION FUNCTION
- ONLY FOR CLASSIFICATION
- CONVERGENCE THEOREM

- perceptron learning algorithm: minimize # misclassified patterns:
find w s.t. $\text{sign}(w^T x) = d$ (or $w^T x > 0$)
- online algorithm: a step made for each input pattern

1. Initialize the weights (either to zero or to a small random value)
2. pick a learning rate η (this is a number between 0 and 1) $\Rightarrow \eta \in [0, 1]$
3. until stopping condition is satisfied (e.g. weights don't change):

For each training pattern (x, d) ($d = +1$ or -1): x : input pattern vector d : desired output

• Compute output activation $out = \text{sign}(w^T x)$ $[+1, -1]$ (omitting T)

• If $out = d$, don't change weights, i.e. "minimize (only) misclassifications"

• If $out \neq d$, update the weights:

$$w_{\text{new}} = w + \eta dx$$

use $\text{sign}(w^T x) < 0 \Rightarrow w_{\text{new}} = w + \eta x$
 $\Rightarrow w_{\text{new}} = w + \eta x$
add $+ \eta x$ if $w^T x \leq 0$ and $d = +1$
- ηx if $w^T x > 0$ and $d = -1$

or in a different form:

Hebbian learning on the rule
 $w_{\text{new}} = w + \eta d x$
 $w_{\text{new}} = w + \eta (d - out) x =$
 $= w + \int x$
 $\int = \text{error signal}$
change w proportionally to the error

- a PERCEPTRON represent a LINEAR DECISION BOUNDARY (LDB) only
- PERCEPTRON can solve LINEAR SEPARATION PROBLEM
- a PERCEPTRON always LEARN a SOLUTION \Rightarrow PERCEPTRON CONVERGENCE THEOREM

PERCEPTRON CONVERGENCE THEOREM:

- IL PERCEPTRON GARANTISCE LA CONVERGENZA (classifica correttamente tutti gli input patterns) IN UN NUMERO FINITO DI STEPS IF THE PROBLEM IS LINEAR SEPARABLE (posso ottenere diverse soluzioni (first solution is not unique))

PERCEPTRON LEARNING ALGO VS LMS

$$w_{\text{new}} = w + \eta (d\text{-out}) x$$

- Apparently similar but note that:
 - LMS rule was derived (by the gradient) without threshold activation functions: minimization of the error of the linear unit (using directly wx)
 - The perceptron uses $out = \text{sign}(w^T x)$
- Hence, for training:
 - $\delta = (d - w^T x)$ for the LMS approach versus
 - $\delta = (d - \text{sign}(w^T x))$ for the perceptron learning alg.
- Of course the model trained with LMS can still be used for classification applying the threshold function $h(x) = \text{sign}(w^T x)$ (LTU)

Differences Summary	
Perceptron Learn. Alg.	LMS Alg.
<ul style="list-style-type: none"> • Min. misclassifications ($out = \text{sign}(w^T x)$) • it always converges for linear separable problems (in a finite number of steps) to a perfect classifier • (else no convergence) • difficult to extend to networks of units (NN) 	<ul style="list-style-type: none"> • Min. $E(w)$ with $out = w^T x$ • asymptotic convergence (*), also for non linear separable problems • not always zero classification errors, linear separable problems • it can be extended to networks of units (NN), using the gradient based approach

(*) example of asymptotic convergence (training with MSE loss) 41

Proof: perceptron convergence theorem

• focus on: POSITIVE PATTERN TASK (classification)

- assume: $\bullet (x_i, d_i)$ in the TR with $d_i = +1 \parallel -1$ (right pattern)

OSS: LINEAR SEPARABLE means $\Rightarrow \exists w^*$ solution s.t.

(desired output + pattern i)

$$d_i \cdot (w^* \cdot x_i) \geq \alpha = \min_i d_i (w^* \cdot x_i) > 0$$

is always positive

$$\Rightarrow \text{hence } w^* (d_i \cdot x_i) \geq \alpha$$

- defining $x_i' = (d_i \cdot x_i)$ then w^* is a solution iff w^* is a solution of $(x_i', +1)$

if pattern is positive
hyper plane is positive
the product is positive

(iP) w^* solve ^{the task} $\Rightarrow d_i (w^* \cdot x_i) \geq \alpha \Rightarrow (w^* \cdot d_i \cdot x_i) \geq \alpha$

$$\Rightarrow (w^* \cdot x_i') \geq \alpha \Rightarrow w^* \text{ is a solution of } (x_i', +1)$$

only if w^* is a solution of $(x_i', +1)$ is a unit $\Rightarrow (w^* \cdot d_i \cdot x_i \geq \alpha) \Rightarrow d_i (w^* \cdot x_i) \geq \alpha$

transformed problem

$\Rightarrow w^*$ solve for x_i

assuming: $w(0) = 0$ (at step 0) \Rightarrow WAT step 0 (updated weight is zero)

$$\cdot \eta = 1$$

$$\cdot \beta = \max |x_i|^2 \rightarrow \text{euclidean norm of the vector patterns in the TR}$$

- after q errors (all false negative)

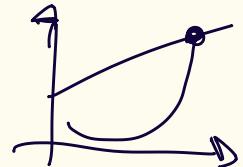
$$w(q) = \sum_{j=1 \rightarrow q} x_{ij} \quad \begin{array}{l} \text{ij: patterns \& misclassified patterns} \\ \text{sum up } x \text{ thru the updating of weights} \end{array}$$

$$\text{because } w(j) = w(j-1) + x_{ij}$$

(Rule: $w_{\text{new}} = w + \eta dx$) but all d are +1, only positive patterns, and $\eta = 1$)
 $\hookrightarrow d=1 \text{ (all pos.)}; \eta = 1$

BASIC IDEA: can find lower and upper bound to $|w|^2$ as a function of q^2 steps
 (lower \hookrightarrow q steps (upper))

\Rightarrow can find q (# steps) s.t. algo converges \rightarrow



PT 1: lower bound on $|w(q)|^2 = w \cdot w(q) = w \sum_{j=1 \rightarrow q} x_{ij} \geq q\alpha$

Cauchy Schwartz: $(w \cdot v)^2 \leq |w|^2 |v|^2$ where $|w| = \|w\|_2$

Recall that
 $\alpha = \min_i (w^* \cdot x_i)$

$$|w^*|^2 |w(q)|^2 \geq (w \cdot w(q))^2 \geq (q\alpha)^2$$

$$\Rightarrow |w(q)| \geq \frac{(q\alpha)^2}{|w^*|^2}$$

PT 2: upper bound
 on $|w(q)|^2$

$$|a+b|^2 = \sum_i (a_i + b_i)^2 = \sum_i a_i^2 + \sum_i 2a_i b_i + \sum_i b_i^2 = |a|^2 + 2ab + |b|^2$$

$$|w(q)|^2 = |w(q-1) + x_{iq}|^2 = |w(q-1)|^2 + 2w(q-1) \cdot x_{iq} + |x_{iq}|^2$$

for q -th error this is < 0

$$\hookrightarrow |w(q)|^2 \leq |w(q-1)|^2 + |x_{iq}|^2$$

and by iteration

$$(w(0) = 0) \Rightarrow |w(q)|^2 \leq \sum_{j=1 \rightarrow q} |x_{ij}|^2 \leq \alpha \beta$$

PT 3: upper & lower

$$\alpha \beta \geq |w(q)|^2 \geq \frac{(\varphi \alpha)^2}{|w^*|^2}$$

$$\varphi \beta \geq \varphi^2 \alpha^2 \quad \beta \geq \alpha^2 \Rightarrow$$

$$\beta = \max_i |x_i|^2$$

$$\boxed{\varphi \leq \frac{\beta}{\alpha^2}}$$

the upper bound

convergence theorem: th. P. converge (class corrct. \times & th' pattern) in
a finite steps if prob. lin. separable

- given (x_i, d_i) in the TR with $d_i = 1 \text{ or } -1$

linear separable $\Rightarrow \exists w^*$ solution s.t. $d_i(w^* x_i) \geq \alpha = \min_i d_i(w^* x_i) > 0$

$$\text{hence } w^*(d_i x_i) \geq \alpha$$

defining $x_i' = (d_i x_i)$ then w^* is a solution iff
 w^* is a solution of $(x_i', 1)$

\Rightarrow if w^* solve the task $\Rightarrow d_i(w^* x_i) \geq \alpha \Rightarrow (w^* d_i x_i) \geq \alpha \Rightarrow$
 $\Rightarrow (w^* x_i') \geq \alpha \Rightarrow w^*$ is a solution of $(x_i', +1)$

only if w^* is a solution of $(x_i', +1) \Rightarrow (w^* d_i x_i) \geq \alpha \Rightarrow d_i(w^* x_i) \geq \alpha$
 $\Rightarrow w^*$ solve for x_i

- assume: $w(0) = 0$

$$\eta = 1$$

$$\beta = \max_i \|x_i\|^2$$

$$\boxed{\text{AFTER } q \text{ errors}} \quad w(q) = \sum_{j=1 \dots q} x_{i,j}$$

↓
index
of
miscls
pat.

Cut rule: $w_{\text{new}} = w_{\text{old}} + \eta \frac{d}{\|x\|} x$

$$\Rightarrow w(j) = w(j-1) + x_{i,j}$$

• can find lower & upper bound to $|w|^2$ as a function of q^2 steps
end q steps (upper bound)

\Rightarrow can find q (# steps) s.t. alg converges

• lower bound $|w(q)|^2$: $w^* w(q) = \underbrace{w^*}_{w^* w^*} \sum_{j=1 \dots q} x_{i,j} \geq q\alpha$

$$|w^*|^2 |w(q)|^2 \geq (w^* w^*)^2 \geq (q\alpha)^2$$

convg swrtz

$$|w(q)|^2 \geq \frac{(q\alpha)^2}{|w^*|^2}$$

lower bound

$$\alpha = \min_i (w^* x_i)$$

UPPER BOUND

$$|w(q)|^2 = |w(q-1) + x_{iq}|^2 = |w(q-1)|^2 + \underbrace{2w(q-1)x_{iq}}_{\text{for the } q^{\text{th}} \text{ error is } < 0} + |x_{iq}|^2$$

$$|w(q)|^2 \leq |w(q-1)|^2 + |x_{iq}|^2$$

by iteration $w(0)=0$

$$\beta = \max_i |x_i|^2$$

$$|w(q)|^2 \leq \sum_{j=1-p}^p |x_{ij}|^2 \leq q\beta$$

UPPER BOUND

UPPER B

$$q\beta \geq |w(q)|^2 \geq \frac{(q\alpha)^2}{|w^*|^2} \text{ lower}$$

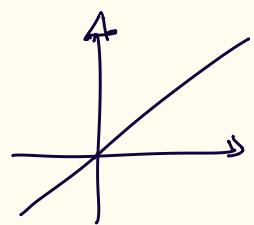
$$\alpha' = \frac{\alpha^2}{|w^*|^2}$$

$$q\beta \geq q^2 \alpha'$$

$$\beta \geq q\alpha'$$

$$q \geq \frac{\beta}{\alpha'}$$

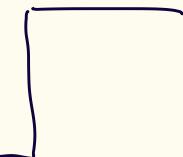
ACTIVATION FUNCTIONS for neurons



1) Linear:
 $\begin{cases} \text{net}_i = \sum_j w_{ij} x_j \\ o_i = f(\text{net}_i(x)) \end{cases}$ $f(x) = \sum w_i x_i$
 $\Rightarrow \text{derivative is } 1$

2) Threshold (perception/CTU):

$$\begin{cases} \text{net}_i = \sum w_{ij} x_j \\ o_i = \text{SIGN}(\text{net}_i(x)) \end{cases} \Rightarrow \text{DERIVATIVE } \nabla \text{ DEFINED} \quad (\text{why isn't used})$$

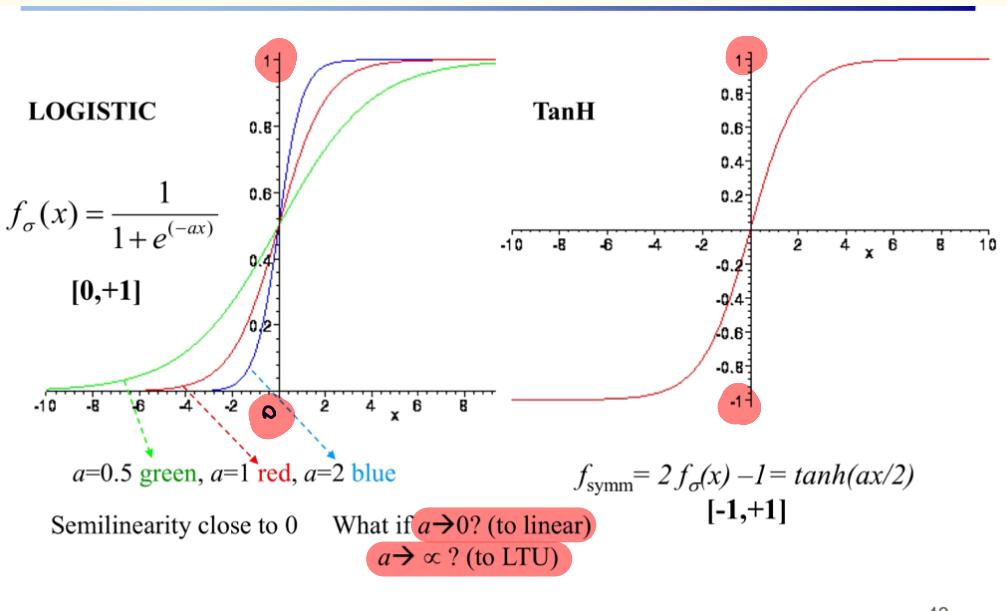
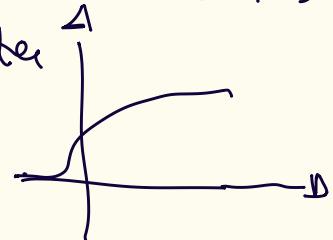


3) Logistic / Sigmoid: NON-LINEAR SQUASHING FUNCTION: assumes values only in $[0, 1]$

2) Logistic

$$f_\sigma(x) = \frac{1}{1 + e^{-ax}} \quad a - \text{slope parameter}$$

• PROPERTY: SMOOTHED DIFFERENTIABLE THRESHOLD FUNCTION



Derivative Logistic

$$a=1$$

$$\frac{\partial f_\sigma(x)}{\partial x} = f_\sigma(x)(1 - f_\sigma(x))$$

Derivative tanh

$$\frac{\partial f_{\text{TANH}}(x)}{\partial x} = 1 - f_{\text{TANH}}(x)^2$$

- Rectifier \rightarrow **ReLU** (Rectified Linear Unit): see Deep Learning section

It becomes a default choice for Deep models, so it (and its variants) deserves more attention discussing Deep Learning

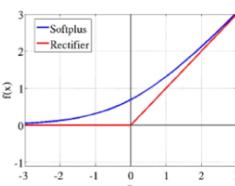
$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

$$f(x) = \max(0, x)$$



- Softplus** (smooth approximation)

$$f(x) = \ln(1 + e^x)$$



- An updated list (Note: many perform comparably!):

https://en.wikipedia.org/wiki/Activation_function

LMS with Sigmoid for

- Sigmoid $\hat{=}$ smoothed differentiable threshold function
- \Rightarrow can derive LMS algorithm \Rightarrow computing GRADIENT OF MEAN SQUARE LOSS FUNCTION
 - from $\sigma(x) = x^T w \Rightarrow \sigma(x) = f_\sigma(x^T w)$
 - \Rightarrow FIND w to minimize the DUAL SUM OF SQUARES

$$E(w) = \sum_{p=1}^l (d_p - o_n(x_p))^2 = \sum_{p=1}^l (d_p - f_\sigma(\text{net}(x_p)))^2$$

COMPUTE THE GRADIENT

- start derivative of E_p w.r.t w_j :

$$\frac{\partial E_p}{\partial w_j} = \frac{\partial E_p}{\partial \text{net}} \cdot \frac{\partial \text{net}}{\partial w_j} = \frac{\partial \text{net}(x_p)}{\partial w_j} = \sigma'(x_{p,j} w_i) \frac{\partial x_{p,j} w_i}{\partial w_j} = (x_p)_j$$

$$\begin{aligned} \frac{\partial E_p}{\partial \text{net}} &= \left(\frac{\partial E_p}{\partial \text{out}} \right) \cdot \left(\frac{\partial \text{out}}{\partial \text{net}} \right) \Rightarrow \frac{\partial f(\text{net})}{\partial \text{net}} = f'_{\text{net}} \\ \frac{\partial (d_p - o_{n,p})^2}{\partial \text{out}_p} &= 2(d_p - o_{n,p}) \cdot 2 \frac{(d_p - \text{out}_p)}{\partial \text{out}_p} = -2(d_p - o_{n,p}) \end{aligned}$$

$$= -2(d_p - o_{n,p}) \cdot f'_{\text{net}}(x_p)_j$$

$$E(w) = \sum_{p=1}^l (d_p - \sigma(x_p))^2$$

$$\begin{cases} \text{net}_i = \sum_j w_{ij} x_j \\ o_i = f_\sigma(\text{net}_i(x)) \end{cases}$$

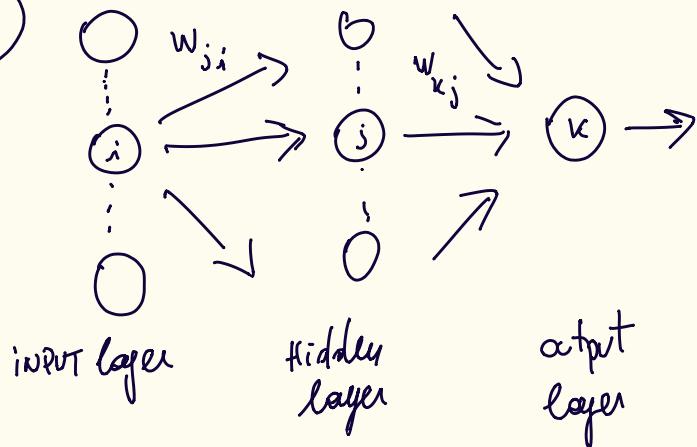
$$\frac{\partial E_p}{\partial w_j} = \frac{\partial E_p}{\partial \text{net}} \cdot \left(\frac{\frac{\partial \text{net}(x_p)}{\partial w_j}}{\frac{\partial \text{net}}{\partial w_j}} \right) \rightarrow \frac{2 \sum_i x_{pi} w_i}{2 w_j} = x_{pj}$$

$$\frac{2(d_p - o_p)^2}{2 o_p} \cdot \left(\frac{\frac{\partial o_p}{\partial \text{net}}}{\frac{\partial \text{net}}{\partial w_j}} \right) = f'_o(\text{net})$$

$$2(d_p - o_p) \cdot \frac{\partial(d_p - o_p)}{\partial o_p} = -2(d_p - o_p)$$

$$\Rightarrow -2(d_p - o_p) \cdot f'_o(\text{ret})$$

(MLP)



$$h(x) = f_k \left(\sum_j w_{kj} f_j \left(\sum_i w_{ji} x_i \right) \right)$$

feed
forward

non-linear function
(sigmoid)

input

- UNITÀ CONNESSE DA, DA DEI LINK PRETTI
- ORGANIZZATE IN LAYER
- INPUT layer poiché nello hidden unit
- feedforward architecture
- TWO LAYER FEEDFORWARD NN

$$h(x) = f_u \left(\sum_j w_{uj} f_j \left(\sum_i w_{ji} x_i \right) \right)$$

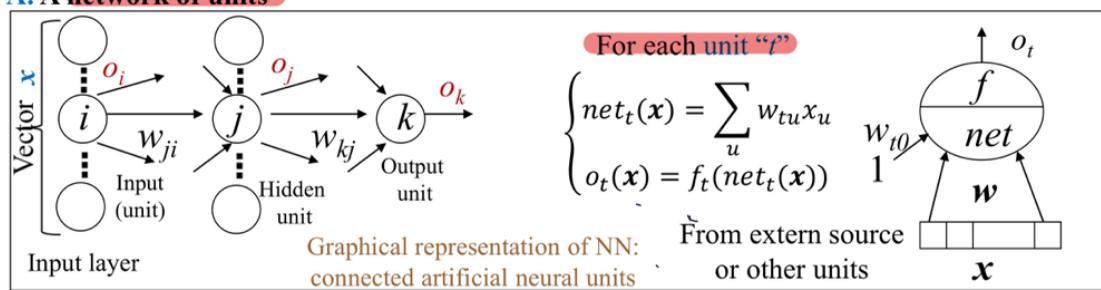
$$h(x) = f_u \left(\sum_j w_{uj} f_j \left(\sum_i w_{ji} x_i \right) \right)$$

The NN and their Units

Standard feedforward NN (with one hidden layer)



A network of units



Notation:

- The index t denotes a generic unit, it can be either k or j (*)
- The index u denotes a generic input component, it can be either i or j (*)
- In the unit, x is a generic input from an external source (input vector) or from other units (*) according to the position of the unit in the network
- If we load the pattern x in the input layer, we can use the notation with o for both the inputs and the hidden units outputs. Hence, inside the network, the input to each unit t from any source u (through the connection w_{tu}) is typically denoted as o_u (we will use this style in the back-prop derivation).

57

MLP need
A non-linear
activation function
in the hidden
layer
 \Rightarrow otherwise
becomes
a linear
model

$$\begin{cases} \text{net}_t(x) = \sum_u w_{tu} x_u \\ \text{OUT}_t(x) = f_t(\text{net}_t(x)) \end{cases}$$

$$\begin{cases} \text{net}_t(x) = \sum_u w_{tu} x_u \\ \text{OUT}_t(x) = f_t(\text{net}_t(x)) \end{cases}$$

- change w can have f functions \Rightarrow different ff space
- output: discrete || continuous values \Rightarrow classification / response task (just change ACTIVATION FUNCTIONS OF THE OUTPUT UNITS)
- change w \Rightarrow sigmoid + linear activation function

$$h(x) = f_k \left(\sum_j w_{kj} f_j \left(\sum_i w_{ji} x_i \right) \right)$$

NN and LBE

• NN is a function non-linear on the parameters W

• W param. is inside a non-linear activation function

$$h_k(x) = f_k \left(\sum_j w_{kj} f_j \left(\sum_i w_{ji} x_i \right) \right)$$

STANDARD: fixed LBE: $\sum_j w_j \phi_j(x)$ • linear w.r.t. parameters • phi preselected

NN: fleible/adaptive

LBE:

- consider this as phi function

$$h(x) = f_k \left(\sum_j w_{kj} f_j \left(\sum_i w_{ji} x_i \right) \right)$$

$$\phi_j(x, w) = f_j \left(\sum_i w_{ji} x_i \right) \quad \text{phi depends on } w$$

is a LBE where the phi IS NOT Fixed, chosen arbitrarily by you

we have weights inside, changed by TR algorithm

\Rightarrow change automatically according to your task

• each basis function (hidden unit) compute a new non-linear derived features, ADAPTIVELY to the TR DATA

• error function less convex \Rightarrow non-linear depending w.r.t. to w \Rightarrow non-linear optimization problem

2. Universal approximation (important!)

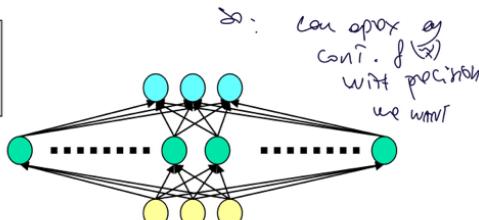


• NN are universal approximators
any degree of accuracy you like

- Many early results (Cybenko 1989, Hornik et al. 1993, etc...)
- Shortly:
 - A single hidden-layer network (with logistic activation functions) can approximate (arbitrarily well) every continuous (on hyper cubes) function (provided enough units in the hidden layer)
[e.g. image it generalizes approximation by finite Fourier series]
 - A MLP network can approximate (arbitrarily well) every input-output mapping (provided enough units in the hidden layers)

Given ε , $\exists h(x)$ s.t. $|f(x) - h(x)| < \varepsilon$
 $\forall x$ in the hypercube

$$h(x) = \sum_j w_{kj} f_j \left(\sum_i w_{ji} x_i \right)$$



- given ε , $\exists h(x)$
such that

$$|f(x) - h(x)| < \varepsilon$$

$\forall x$ in hypercube

Existence theorem! This not provide the algorithm nor the number of units

NN - learning algorithm

: adapt free parameters w to obtain best approximation of the target function

= minimize error via LMS

$$E(w) = R(w_p) = \frac{1}{l} \sum_{p=1}^l (L(h(x_p)), d_p)$$

BACKPROP.

$$L(h(x_p), d_p) = \text{MSE} \\ = (d_p - h(x_p))^2$$

FIND w to min.
residual sum of squares

• STIMARE IL CONTRIBUTO DEGLI HIDDEN UNITS PER

L'ERRORE ALL'OUTPUT LEVEL

=> Computing generalized Δ rule

$$E_{\text{tot}} = \sum_p E_p = \frac{1}{2} \sum_{k=1}^K (d_k - o_k)^2 \quad \text{find } w \text{ to minimize } E_{\text{tot}}$$

Algorithm 1: Back-propagation (the same presented in general for the LMS)

```
1. Initialize all the weights  $w$  in the network and  $\eta$ ;
   Compute units outputs and  $E_{\text{tot}}$ ;
  while  $E_{\text{tot}} > \epsilon$  (desired value or other criteria) do
    |  $\forall w \in w: \Delta w = -\frac{\partial E_{\text{tot}}}{\partial w}$  (Step 1);
    |  $w^{new} \leftarrow w + \eta \Delta w + \dots$  (Step 2); }
    | Compute units outputs and  $E_{\text{tot}}$ ;
  end
  1) compute partial derivatives, gradient
  2) update the weights
```

$$\Delta_w = -\frac{\partial E_{\text{tot}}}{\partial w} = -\sum_p \frac{\partial E_p}{\partial w} = \Delta_p w$$

$$\Delta_p w_{tm} = -\frac{\partial E_p}{\partial w_{tm}} = \boxed{-\frac{\partial E_p}{\partial \text{net}_t}} \cdot \boxed{\frac{\partial \text{net}_t}{\partial w_{tm}}} = \boxed{\delta_t} \boxed{o_m}$$

$$\delta_t = -\frac{\partial E_p}{\partial \text{net}_t} = \boxed{-\frac{\partial E_p}{\partial o_t}} \cdot \frac{\partial o_t}{\partial \text{net}_t} = \boxed{\frac{\partial E_p}{\partial o_t}} \cdot f'_t(\text{net}_t) \rightarrow \frac{\partial f_t(\text{net}_t)}{\partial \text{net}_t} = f'_t(\text{net}_t)$$

$$\therefore \frac{\partial E_p}{\partial o_t} = \frac{\partial \frac{1}{2} \sum_{k=1}^K (d_k - o_k)^2}{\partial o_t} = d$$

GRADIENT COMPUTATION

$$\Delta_w = - \frac{\partial E_w}{\partial w} = - \sum_{p=1}^e \frac{\partial E_p}{\partial w} = \sum_{p=1}^e \Delta_p w$$

$$net_T = \sum_{i=1}^m w_{ti} o_i$$

$$o_t = f_t(\text{net}_t)$$

- compute the component $\Delta_p w_{T,U}$ of the vector w the $\Delta_p w_{T,U}$

$$\Delta_p w_{T,U} = - \frac{\partial E_p}{\partial w_{T,U}} = \frac{\partial E_p}{\partial \text{net}_t} \cdot \frac{\partial \text{net}_t}{\partial w_{T,U}} = \delta_t \cdot o_U$$

$$o_U = \frac{\partial \text{net}_t}{\partial w_{T,U}} = \frac{\partial \sum_j w_{tj} o_j}{\partial w_{T,U}} = * = o_U$$

* cuz $w_{tn} \cdot o_n$ is constant wrt $w_{T,U}$
but for $\delta_t = 0$

$$\Rightarrow \Delta_p w_{T,U} = \delta_t \cdot o_U$$

$$- \frac{\partial E_p}{\partial \text{net}_t} = - \frac{\partial E_p}{\partial o_t} \cdot \frac{\partial o_t}{\partial \text{net}_t} = \frac{\partial f_t(\text{net}_t)}{\partial \text{net}_t} = f_t'(\text{net}_t)$$

2 CASES

① OUTPUT UNIT : assume $t=k$ (consider an output unit of the network)

$$- \frac{\partial E_p}{\partial o_k} = - \frac{\partial \frac{1}{2} \cdot \sum_{R=1}^K (d_R - o_R)^2}{\partial o_k} = 2 \cdot \frac{1}{2} (d_k - o_k) = (d_k - o_k)$$

$\delta_k = (d_k - o_k) \cdot f_k'(\text{net}_k)$

output unit

② HIDDEN LAYER $t=j$ > chain rule

$$- \frac{\partial E_p}{\partial o_j} = \sum_{k=1}^K - \frac{\partial E_p}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial o_j} = \sum_{k=1}^K \delta_k \cdot \frac{\partial \sum_{j=1}^K w_{kj} \cdot o_j}{\partial o_j} = \sum_{k=1}^K \delta_k \cdot w_{kj}$$

SO SUBSTITUTING IN

$$\delta_j = \left(\sum_{k=1}^K \delta_k \cdot w_{kj} \right) \cdot f_j'(\text{net}_j)$$



Summary

We derived, for two layers,

$$\Delta_p w_{tu} = \delta_t \cdot o_u$$

where δ_t is the error signal available to the unit t , and o_u is the input to t from a generic unit u , i.e. the input through the connection weighted by w_{tu} . In particular:

- ✓ $\Delta_p w_{kj} = \delta_k \cdot o_j$ if $t = k$ (output unit) and the input is from unit j
- ✓ $\Delta_p w_{ji} = \delta_j \cdot o_i$ if $t = j$ (hidden unit), and input i

The other two equations specify the error signals:

$$\delta_k = (d_k - o_k) \cdot f'_k(\text{net}_k) \quad \text{if } t = k \text{ (output unit)}$$

$$\delta_j = \left(\sum_{k=1}^K \delta_k w_{kj} \right) \cdot f'_j(\text{net}_j) \quad \text{if } t = j \text{ (hidden unit)}$$

backprop.

$$\Delta_w = -\frac{\partial E_{\text{tot}}}{\partial w} = -\sum_{p=1}^e \frac{\partial E_p}{\partial w} = \sum_{p=1}^e \Delta_p w$$

$$\cdot \Delta_p w_{T,U} = \frac{-\partial E_p}{w_{T,U}} = -\frac{\partial E_p}{\partial \text{net}_T} \cdot \left(\frac{\partial \text{net}_T}{\partial w_{T,U}} \right) = \delta_T o_M$$

$$= \frac{\partial \sum_j w_{+j} o_j}{\partial w_{T,U}} \quad \text{when } j \neq U \\ = o_M$$

$$\therefore \frac{-\partial E_p}{\partial \text{net}_T} = \frac{-\partial E_p}{\partial o_T} \cdot \left(\frac{\partial o_T}{\partial \text{net}_T} \right) = \delta_T$$

$$\frac{\partial f'_T(\text{net}_T)}{\partial \text{net}_T} = f'_T \text{ net}_T$$

2 cases

$t = \text{output}$
 $t = \text{hidden unit}$

$$\boxed{T=K \text{ output}} \Rightarrow \boxed{\delta_k} = \frac{\partial E_p}{\partial o_K} \cdot \frac{\partial o_K}{\partial \text{net}_K} = \boxed{(\delta_K - o_K) \cdot f'_K \text{ net}_K}$$

$$\frac{-\partial \frac{1}{2} \sum_{R=1}^K (d_R - o_R)^2}{\partial o_K} = (\delta_K - o_K)$$

$$\therefore \delta_j = \boxed{\frac{-\partial E_p}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j}} = \boxed{\sum_{k=1}^K \frac{\partial E_p}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial o_j} = \sum_{k=1}^K \delta_k w_{kj} \circ f'_j(\text{net}_j)}$$

$$\frac{\partial f'_j(\text{net}_j)}{\partial \text{net}_j} = f'_j'(\text{net}_j)$$

CHAIN RULE

$$\sum_{k=1}^K \delta_k \cdot w_{kj} = \sum_{k=1}^K \delta_k \cdot \frac{\partial \text{net}_k}{\partial o_j} = \sum_{k=1}^K \delta_k \cdot \frac{\partial \text{net}_k}{\partial \text{net}_T} \cdot \frac{\partial \text{net}_T}{\partial o_j}$$

$$\therefore \frac{\partial \sum_s w_{ks} \cdot o_s}{\partial o_j} = (w_{kj})$$

$$\Rightarrow \boxed{\delta_j = \sum_{k=1}^K \delta_k w_{kj}}$$

Hyperparameters: variazioni da scorrere x initializzazone e fine il TRAIN

LEARNING CURVE ANALYSIS

(1) EARLY STOPPING: use a validation set to decide when to stop
(look over a set of epochs)

- since training means

make grows # params

\Rightarrow go back limits complexity

if validation error ↑
 \Rightarrow stop

\Rightarrow go back to best solution
before see increase
in val set

(2) REGULARIZATION

training \Rightarrow weights increase: optimize the loss considering the weight values

\Rightarrow TIKONOV: add penalty term to the error function

$$\text{loss}(w) = \sum_p (d_p - o(x_p))^2 + \gamma \|w\|^2$$

effect: weight decay: $w_{\text{new}} = w + \text{eta} \cdot \Delta_w - \lambda w$

(if applied to linear model: ridge regression)

• Loss: \rightarrow the objective function: MSE + PENALTY TERM

• Error: Risk for the data term (MSE, dependency on data): $\sum_p (d_p - o(x_p))^2$
TO EVALUATE MODEL ERROR

Regularization NOT CONTROL THE STABILITY

BUT CONTROL THE COMPLEXITY

○ —

CMTS: Control the complexity, MODEL SELECTION issues
too few hidden units \rightarrow underfitting
too many \rightarrow overfitting

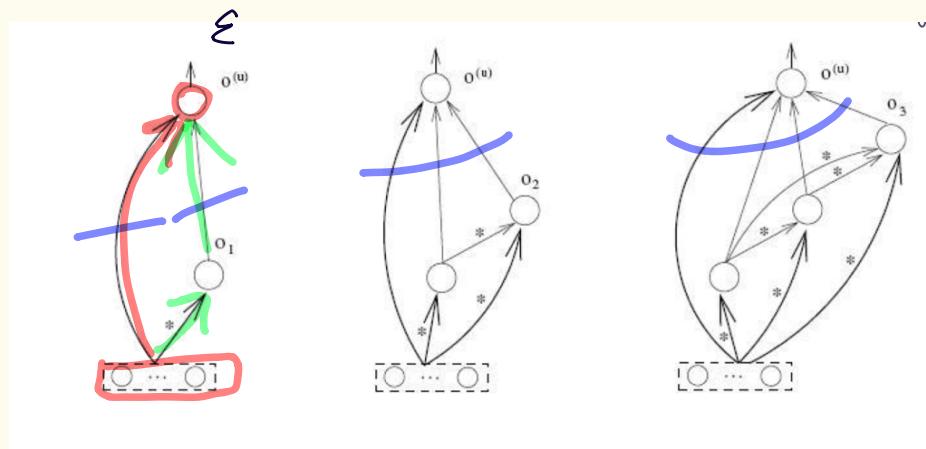
\Rightarrow 2 METHODS: CONSTRUCTIVE APPROACH
PRUNING METHODS

CONSTRUCTIVE APPROACH: incremental, build network starting with a minimal configuration, add new units and connections during training

example: CASCADE CORRECTION learning algorithm

- START WITH A MINIMAL NETWORK
- ADD UNITS UNTIL THE ERROR IS BELOW
- CC ^{OP} LEARN ABOUT NETWORK WEIGHTS AND NETWORK TOPOLOGY (THREE UNITS)
- H SPACE OF FLEXIBLE SIZE SINCE THREE UNITS DECIDED BY THE ALGORITHM
- TRAIN A SINGLE UNIT FOR EACH STEP
(then decide to add or not the units)

Initial step on input ϵ allows 1 output unit: see errors down \Rightarrow stop?



altrimenti aggiunge un
nodo nascosto e lo guarda,
e lo allena per essere
consistente con l'errore che
ha avuto dall'output layer

- si ALLENA LA RETE
(non ha nascosto nato
ma gli esisti)

- add one fuzzy weight
and complete the
output of this hidden
unit to the residual
error

Learning for CC

Dip. Informatica
University of Pisa

- Specifically, the algorithm works interleaving the minimization of the total error function (LMS), e.g. by a simple backpropagation training of the output layer, and the maximization of the (non-normalized) correlation, i.e. the covariance, of the new inserted hidden (candidate) unit with the residual error.

$$S = \sum_k \left| \sum_p \frac{(\bar{o}_p - \text{mean}_p(\bar{o}_p))}{\bar{o}} \frac{(E_{p,k} - \text{mean}_p(E_{p,k}))}{\bar{E}} \right|$$

E: residual error
with $\bar{o}_{p,k}$ at the output layer
 p : pattern, k : output unit

- **Exercise:** derivation (sketch on blackboard)
hints: $df/dx = \text{sign}(f) df'/dx$; assume no effect of variations of \bar{o}_p over $\bar{o} = \text{mean}_p(\bar{o}_p)$

- Result:

$$\frac{\partial S}{\partial w_j} = \sum_k \text{sgn } S_k \sum_p (E_{p,k} - \text{mean}_p(E_{p,k})) f'(net_{p,h}) I_{p,j} \quad \downarrow \text{Input}$$

- Weight update: why $+ dS/dw_j$?
 $\Rightarrow +$ we are maximizing gradient ascent instead of LMS desc.

A. Micheli

$$\begin{aligned} \frac{\partial S}{\partial w_j} &= \frac{\partial \left(\sum_{k=1}^K \left| \sum_{p=1}^P (\bar{o}_p - \text{avg}_p(\bar{o})) \cdot (E_{p,k} - \text{avg}_p(E_{p,k})) \right| \right)}{\partial w_j} \\ &= \sum_{k=1}^K \left[\text{sgn}(S_k) \cdot \sum_{p=1}^P \frac{\partial (\bar{o}_p - \text{avg}_p(\bar{o})) \cdot (E_{p,k} - \text{avg}_p(E_{p,k}))}{\partial w_j} \right] \quad \text{COMBIN BULC} \\ &= \sum_{k=1}^K \left[\text{sgn}(S_k) \cdot \sum_{p=1}^P \left(E_{p,k} - \text{avg}_p(E_{p,k}) \right) \cdot \frac{\partial (\bar{o}_p - \text{avg}_p(\bar{o}))}{\partial net_p} \cdot \frac{\partial net_p}{\partial w_j} \right] \\ &= \sum_{k=1}^K \left[\text{sgn}(S_k) \cdot \sum_{p=1}^P \left(\bar{E}_{p,k} - \text{avg}_p(\bar{E}_k) \right) \cdot 1 \cdot \frac{\partial net_p}{\partial w_j} \cdot \frac{\partial \bar{o}_p}{\partial net_p} \right] \\ &= \sum_{k=1}^K \left(\text{sgn}(S_k) \cdot \sum_{p=1}^P \left(\bar{E}_{p,k} - \text{avg}_p(\bar{E}_k) \right) \cdot f'(net_p) (I_{p,j}) \right) \end{aligned}$$

SLT

- can approximate estimation of NL set error (for model selection) and TS error (model assessment)
=> SRM: structural risk minimization and VC-minimization

VC-DIMENSION: measure of flexibility (complexity of or CASS of)

- $N_{\text{data}} = 2^N$ ^{possibili} _{discreti} HYPOTESIS (also for hypothesis space H)
- H fragments X iff H può rappresentare tutte le possibili distinzione su X (o error)
- VC-DIMENSION DI UNA CLASSE DI FUNZIONI H È LA CARDINALITÀ MASSIMA DI UN INSIEME DI PUNTI IN X (INSTANCE SPACE) CHE PUÒ ESSERE FRAMMENTATO DA H

$\xrightarrow{\quad}$ instance space
 $X=3$

$$VC(H) = p$$

- H frammenta sempre un insieme (una configurazione) di p punti
- H non frammenta nemmeno insieme (configurazione) di $p+1$ punti

(quindi basta trovare una configurazione di punti che non' viene frammentata)

$$VC(H) \geq 3 \quad VC(H) < 4 \quad \Rightarrow \quad VC(H) = 3$$

-
- In general: The VC dimension of a class of linear (separator/decision) hyperplanes (LTU) in a n -dimensional space is $n+1$ on 2 dim 3
- Note: Upper-bound, it can be less constraining the model, we'll see how... if
- our aim is to reduce VC-dim dim n =>

o V = # punti
che possono essere
classeificate
senza errori

es [Determin] must be linear in the free parameters of a linear model

STRUCTURAL RISK MINIMIZATION

- STR USA la VC-dimension come CONTROLLING PARAMETER per minimizzare il GENERALIZATION BOUND SU R
=> want minimize the risk
- assuming FINITE VC-DIMENSION posso definire modelli ANNIATI delle spartizioni delle ipotesi in base alle VC-dimension

$$H_1 \subseteq H_2 \subseteq \dots \subseteq H_m$$

$$VC(H_1) \leq VC(H_2) \leq \dots \leq VC(H_m)$$

Some example?
NN \leq \leq NN 100 units or epochs
• NN with increasing number of Hidden units (roughly), but also the #epochs (we know that also the size of weights care and a NN is a variable-size H)
• Polynomial of increasing degree M nikonov
• Increasing values for c , where $\|w\| < c$ for regularization (or controlled by λ)
• Increasing number of nodes in a decision tree (roughly)
• Note: important hyper-parameters → related to VC-dim
 $\|w\| \leq \frac{C}{\lambda}$
 $\|w\| \leq \dots \leq \dots$ with 17 steps

Growing VC-dim: empirical (TR) error decreases, VC-confidence increases

Structural Risk Minimization: find a trade-off on the bound

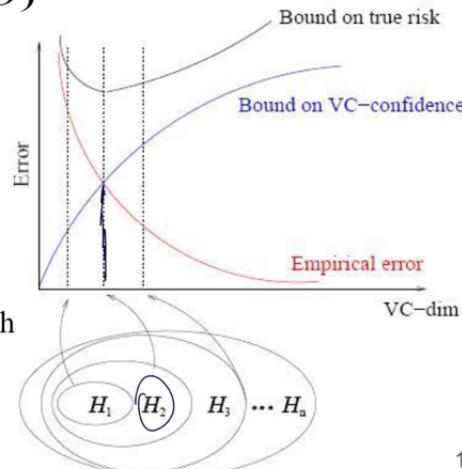
(as seen in the introduction using ℓ_1):

$$R \leq R_{emp} + \varepsilon(1/l, VC, 1/\delta)$$

- $\mathcal{H}_1 \subseteq \mathcal{H}_2 \subseteq \dots \subseteq \mathcal{H}_n$

- $VC(\mathcal{H}_1) \leq \dots \leq VC(\mathcal{H}_n)$

e.g. NN with increasing number of
Hidden units



18

SVM

- N: # examples (instead of l)
- m: dimension of input vector (instead of n)
- b: bias (instead of w.)

SVM

- linear machine (LTV, perceptron)
- maximization on separation margin
- structural risk minimization

- given training set $T = \{(x_i, d_i)\}_{i=1}^N$

- find hyperplane of equation $w^T x + b = 0$ to separate the examples

$$w^T x_i + b \geq 0 \quad \text{for } d_i = +1$$

$$w^T x_i + b < 0 \quad \text{for } d_i = -1$$

- $f(x) = w^T x + b$ è chiamata discriminative function
- $h(x) = \text{sign}(f(x))$ è l'ipotesi

SEPARATION MARGIN = ρ = il doppio della distanza tra l'iperpiano e i punti più vicini
 • passo avere due iperpiani = non con iperelli diverse \Rightarrow solo tra i due interstizi non si può trovare soluzioni separate

• voglio l'OTTIMO iperpiano \rightarrow quello che massimizza il margine ρ

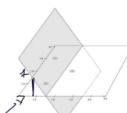
$$\begin{aligned} \text{OPTIMO} \rightarrow & \rightarrow w_0^T x + b_0 = 0 \quad \text{done } w_0^T \circ \text{ is optimal} \\ \rightarrow (\text{trovare}) \rho = \frac{2}{\|w\|} & \quad \text{MASSIMIZZARE } \rho \iff \text{MASSIMIZZARE } \|w\| \end{aligned}$$

FORMA CANONICA DELL'IPERPIANO $\rightarrow d_i(w^T x_i + b) \geq 1 \quad \forall i = 1..N$

$$\begin{cases} w^T x_i + b \geq 0 & \text{if } d_i = +1 \\ w^T x_i + b \leq -1 & \text{if } d_i = -1 \end{cases}$$

Re-scaling w and b (see scaling-freedom property) so that the closest points to the separating hyperplane satisfy $|g(x_i)| = |w^T x_i + b| = 1$, we can write:

$$\begin{cases} w^T x_i + b \geq 1 & \text{if } d_i = +1 \\ w^T x_i + b \leq -1 & \text{if } d_i = -1 \end{cases}$$



In a compact form:

$$d_i(w^T x_i + b) \geq 1 \quad \forall i = 1, \dots, N$$

Support Vector

A support vector $x^{(s)}$ satisfies the previous equation exactly:

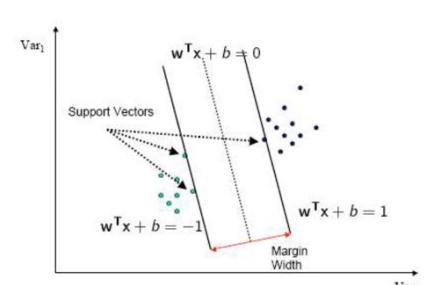
$$d^{(s)}(w^T x^{(s)} + b) = 1$$

\rightarrow POINTS AT THE edge of the MARGIN

support vector $x^{(s)}$ soddisfa l'eq.

$$d^{(s)}(w^T x^{(s)} + b) = 1$$

if satisfies this, all points are correctly classified



OBBLIGO DI TROVARE LA LARGHEZZA DEL MARGIN (Size)

Distance to the hyperplane

Let's denote $g(x)$ as the discriminant function: $g(x) = \mathbf{w}^T \mathbf{x} + b$.

Recall that \mathbf{w}_0 is a vector orthogonal to the hyperplane

Let's denote the distance between x and the optimal hyperplane with r

(according to the orthogonal vector to the hyperplane)
(exactly \mathbf{w}_0)

AIM:
measure
the margin

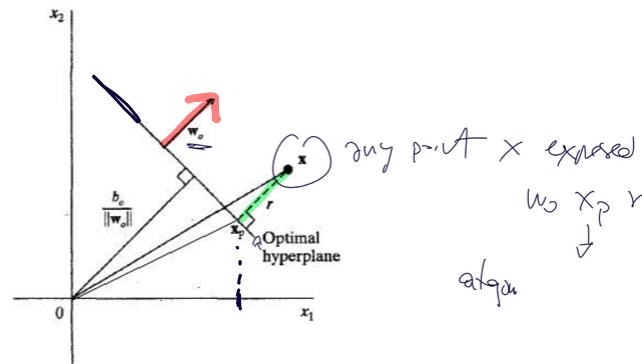


Figure: Distance between a point and the optimum hyperplane.

$$\text{point on the hyperplane} \rightarrow \mathbf{x} = \mathbf{x}_p + r \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|} \rightarrow \text{real movement by } r \text{ and in the direction of } \mathbf{w}_0 \text{ the direction of the normal vector}$$

Introduction to Support Vector Machines

12 / 64

$$x = x_p + r \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|}$$

$$x = x_p + r \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|}$$

$$g(x) = \mathbf{w}_0^T \mathbf{x} + b_0$$

$$g(x) = g\left(x_p + r \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|}\right) = \underbrace{\mathbf{w}_0^T x_p + b_0}_{g(x_p)} + \mathbf{w}_0^T r \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|} = g(x_p) + \underbrace{\mathbf{w}_0^T r \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|}}_{= r \|\mathbf{w}_0\|} =$$

$$= g(x_p) + r \mathbf{w}_0^T \frac{\mathbf{w}_0}{\|\mathbf{w}_0\|} = r \frac{\|\mathbf{w}_0\|^2}{\|\mathbf{w}_0\|} = r \|\mathbf{w}_0\|$$

x_p è un punto
sull'iperpiano
quindi è bello

$$g(x) = r \|\mathbf{w}_0\|$$

$$r = \frac{g(x)}{\|\mathbf{w}_0\|}$$

• considero la distanza tra l'iperpiano
e il support vector positivo $x^{(s)}$

$$\text{i.e. } r \text{ per } x^{(s)} = \frac{g(x^{(s)})}{\|\mathbf{w}_0\|} = \frac{1}{\|\mathbf{w}_0\|} = \frac{P}{2}$$

→ MARGIN
this distance is half of the margin

$$P = \frac{2}{\|\mathbf{w}_0\|}$$

odendo trovò w_0 e b che mi dà la soluzione con il margine massimo
for class. with w errors

HARD MARGIN SVM

given training examples $T = \{(x_i, d_i)\}_{i=1}^N$

find optimal w_0 e b che minimizzano

$$\Psi = \frac{1}{2} w^T w \quad (\min \|w\|)$$

satisfying zero-errors condition (constraint)

$$d_i(w^T x_i + b) \geq 1 \quad \forall i = 1, \dots, N$$

- objective function $\Psi(w)$ is quadratic and convex in w

- constraints are linear in w

- solving this problem scales with size of input space n (later change on N)

WHY INTRODUCING SVM?

GOOD IN TERM OF GENERALIZATION

Why does this improve the generalization?

We fixed the training error on linearly separable problems. Then minimizing the norm of w is equivalent to minimizing the VC dimension and thus to minimizing the capacity term (VC-confidence) ϵ in

$$R[h] \leq R_{emp}[h] + \epsilon(VC, N, \delta)$$

"O VC-confidence thru objective f(x) of primal form"

Relevant!!!

Theorem - Vapnik

Let D be the diameter of the smallest ball around the data points x_1, \dots, x_N . For the class of separating hyperplanes described by the equation $w^T x + b = 0$ the upper bound to the VC dimension is

$$VC \leq \min\left(\lceil \frac{D^2}{\rho^2} \rceil, m_0\right) + 1$$

since $\frac{D^2}{\rho^2} = \text{Radius}^2 \cdot \|w\|^2$, VC can be less than $m_0 + 1$ computed for general hyperplanes by constraining the search to the "regularized" one with maximum margin: regularization \Rightarrow reduce VC \Leftrightarrow control size of weights $\|w\|^2$

Apprendimento Automatico 1 - 2014 () Introduction to Support Vector Machines 24 / 64

\Rightarrow can find an optimal solution in term of generaliz. with no error and no hyperparameters

• Fix the Regul (even posed to sub)
(no error cut is liner. sep.)
 \Rightarrow just minimize $\|w\|$
is equivalent to
minimize the VC-dimension
 \Rightarrow minimize the VC-confidence
imperatively directly through min.

SVM minimize VC-confidence
trough maximization of the margin
thru VAPNIK THEOREM

$$VC \leq \min\left(\lceil \frac{D^2}{\rho^2} \rceil, m_0\right) + 1$$

↳ diameter of the ball
+ dimension of input space

Forma canonica $d_i (w^T x_i + b) \geq 1$

support vector $x^{(s)}$

$$d^{(s)} (w^T x^{(s)} + b) = 1$$

distanza dell'iperpiano da un punto x

$$x = x_p + r \frac{w_0}{\|w_0\|}$$

LUNGHEZZA iperpiano ottimo

$$g(x) = w^T x_0 + b_0$$

$$g(x_p + r \frac{w_0}{\|w_0\|}) = \underbrace{w^T x_p + b_0}_{g(x)} + w^T r \frac{w_0}{\|w_0\|} = g(x) + r w^T \frac{w_0}{\|w_0\|} = r \frac{\|w_0\|^2}{\|w_0\|} = r \|w_0\|$$

$$r = \frac{g(x)}{\|w_0\|} \quad \text{per } x^{(s)}$$

$$\Rightarrow r = \frac{g(x^{(s)})}{\|w_0\|} = \frac{1}{\|w_0\|} = \frac{P}{2}$$

$$g(x) = r \|w_0\| \quad P = \frac{2}{\|w_0\|}$$

nuovo w e b ottimi per MASSIMIZZARE IL MARGINE: HARD MARGIN

- given T samples $T = \{x_i, d_i\}_{i=1}^N$ find w & b s.t. minimize
(COST FUNCTION) $\psi(w) = \frac{1}{2} w^T w$ satisfying two new conditions:

$$d_i (w^T x_i + b) \geq 1$$

THEOREM: maximum margin solution

correspond the MINIMUM possible VC-dimension \Rightarrow best possible generalization capacity

$$R \leq R_{\text{emp}} + \epsilon(N, \delta)$$

$$VC \leq \min \left(\left\lceil \frac{D^2}{\rho^2} \right\rceil, M_0 \right) + 1$$

OPTIMAL VALUES of w and b

$$w_0 = \sum_{i=1}^N \alpha_{0,i} d_i x_i$$

\Rightarrow so the optimal hyperplane

$$w_0^T x + b_0 = 0 \Rightarrow \sum_{i=1}^N \alpha_{0,i} d_i x_i^T x + b_0 = 0$$

KKT - KARLSON CONDITION

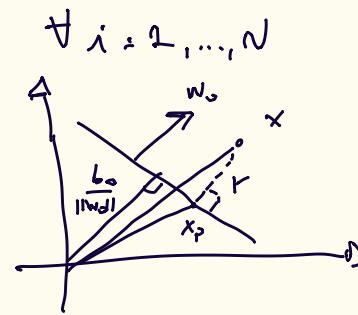
$$\alpha_i (d_i (w^T x_i + b) - 1) = 0$$

If $d_i > 0 \Rightarrow d_i (w^T x_i + b) = 1$ and x_i is a support vector

If x_i is NOT a support vector $\Rightarrow \alpha_i = 0$

\Rightarrow can restrict the optimization to N_s : $w_0 = \sum_{i=1}^{N_s} \alpha_{0,i} d_i x_i$

HYPERSPACE DEPENDS ONLY FROM SUPPORT VECTOR



- Solving dual problem obtaining α
- Compute $w_0 = \sum_{i=1}^N \alpha_{0,i} d_i x_i$
 - Compute $b_0 = 1 - w_0^\top x^{(s)}$ corresponding to positive support vector
 - so $w_0 = \sum_{i=1}^N \alpha_{0,i} d_i x_i$
- decision surface is given by $w_0^\top x + b_0 \Rightarrow \sum_{i=1}^N \alpha_{0,i} d_i x_i^\top x + b_0 \approx 0$
- \Rightarrow classify point x with the sign of $f(x)$

Why does this provide an elegant approach?

For linearly separable data there are many solutions. Vapnik propose an "optimal separating hyperplane" maximizing the margin providing

- an unique solution (not for iterative Perceptron Learn. Alg.) with zero errors (not for LMS) for the binary classifier
- an automated approach to SRM that minimizes VC-confidence (by maximizing the margin) as part of the training process (directly into the optimization process), without hyper-parameters in the linear-separable case
- the use of solver in the class of constrained Quadratic Programming (instead of gradient descent) with a nice dual form (showing support vectors and the dot product among patterns)
- a solution focused on "selected" training data, the support vectors
 - (nice) not depending directly on the point far away from the decision boundary
 - (less nice) hoping that the points at the boundaries are not the noisiest ones.

But what for noisy or non linearly separable data?

\Rightarrow SOFT MARGIN: admit errors (not every separable, solve with error)

• a data point inside the margin (also classification errors)

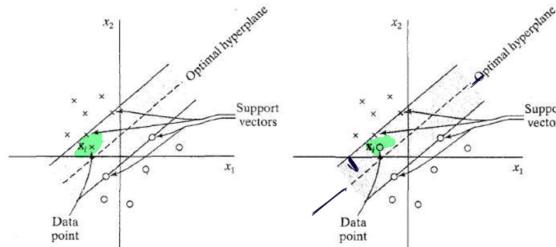


Figure: Patterns violating the margin (with and without classification error).

At least one point violate the exact-fit condition:

$$d_i(w^\top x_i + b) \geq 1$$

Soft Margin.

Note: admitting points inside the margin allows us to have a larger margin.

Apprendimento Automatico 1 - 2014 () Introduction to Support Vector Machines

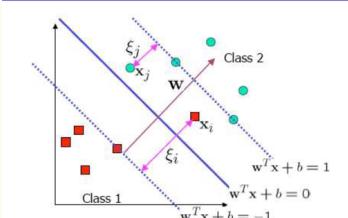
26 / 64

\Rightarrow admitting errors:
 - admit points in the margin
 allow us to have a large margin
 \Rightarrow better generalization

introduce slack-variables

$$\xi_i \geq 0 \quad \forall i=1, \dots, N \quad d_i(w^\top x_i + b) \geq 1 - \xi_i \quad \forall i=1, \dots, N$$

A SUPPORT VECTOR x_i SATISFY: $d_i(w^\top x_i + b) \geq 1 - \xi_i \quad \forall i=1, \dots, N$



VAPNIK THEOREM NOT HOLD

$$\text{POLYNOMIAL FORM } \Psi(w, \varepsilon) = \frac{1}{2} w^T w + C \sum_{i=1}^N \varepsilon_i$$

$$\text{constraint } d_i(w^T x + b) \geq 1 - \varepsilon_i \quad \forall i=1, \dots, N$$

C is REGULARIZATION hyper-parameter

$$\varepsilon_i \geq 0 \quad \forall i=1, \dots, N$$

low $C \Rightarrow$ TR error admitted \Rightarrow possible underfitting

high $C \Rightarrow$ no TR allowed \Rightarrow possible overfitting

we lost SVM

MARGIN WITH C

\Rightarrow HIGH C

A trade off between
Regularization for
and VC-confidence

KAN-TACKER CONDITION

$$\alpha_i (d_i(w^T x_i + b) + \varepsilon_i) = 0 \quad \forall i=1, N$$

$$\sum_i \varepsilon_i = 0 \quad \forall i=1, \dots, N$$

1 constraint

Lagrange multiplier

$0 < \alpha_i < C \Rightarrow \varepsilon_i = 0$ (on the margin)

$\alpha_i = C \Rightarrow \varepsilon_i \geq 0$ (inside the margin)

- Solve the dual problem with respect to $\{\alpha_i\}_{i=1}^N$
- Find w_0 and b_0 from $\{\alpha_{o,i}\}_{i=1}^N$

$$w_0 = \sum_{i=1}^N \alpha_{o,i} d_i x_i$$

$$b_0 = d_j - \sum_{i=1}^N \alpha_{o,i} d_i x_i^T x_j \text{ for a pattern } j \text{ such that } 0 < \alpha_j < C$$

(or an average of all the solutions for numerical stability)

The non zero Lagrangian Multipliers correspond to support vectors.

How do we use it?

As before:

$$g(x) = \sum_{i=1}^N \alpha_{o,i} d_i x_i^T x + b_0$$

$$h(x) = \text{sign}(g(x))$$

PER i PROBLEMI NON LINEARI \Rightarrow need to map in the feature space

\Rightarrow concept of Kernel

non linear function map $\Phi: \mathbb{R}^{m_0} \rightarrow \mathbb{R}^{m_1} \quad x \mapsto \Phi(x)$

Tr set becomes $T = \{(\Phi(x), d_i)\}_{i=1}^N \quad \forall i=1, N$

\Rightarrow hyperplane becomes $w^T \Phi(x) + b = 0$

choice Φ is CHALLENGING in SVM

WITH LBF : non-linear mapping of input patterns to high-dimensional
FEATURE SPACE

(COVER'S THEOREM: the patterns are linearly separable
with high probability in the feature space under such condition (increase
the dimension from original space to high dim. space where we
project the data))

⇒ However, using high-dim. feature space (large LBF) can lead to
OVERFITTING = ~~KERNEL APPROX.~~
(bad complexity & classifier)

→ MANAGE THE FEATURE SPACE IN THE CONTEXT OF REGULARIZED MODEL

↓
complexity of LBF is related to the # free-parameters (# dimension of the new phi expansion) and input dimensionality

THE MARGIN RULE THE COMPLEXITY AND NOT THE DIMENSION OF INPUT

Solving the problem in the Feature-Space

Linear hyperplane in feature space

Non linear function map

$$\begin{aligned}\Phi : \mathbb{R}^{m_0} &\rightarrow \mathbb{R}^{m_1} \\ x &\mapsto \Phi(x)\end{aligned}$$

- The problem is formulated as before, but the training set is now $T = \{(\Phi(x_i), d_i)\}_{i=1}^N$
- The hyperplane is now $w^T \Phi(x) + b = 0$

Incorporating the bias in the weight vector:

- $w(0) = b$
- $\Phi_0(x) = 1$

$$\Phi(x) = (\phi_0(x) = 1, \phi_1(x), \dots, \phi_{m_1}(x))^T$$

Decision surface equation: $w^T \Phi(x) = 0$ i.e. a linear expansion of the basis functions (where (i) the dimension of the enlarged space is allowed to get very large and (ii) the complexity depends from the margin, not directly on the space dimension)

↑
VC-DIMENSION ↑ INCREASE
ACCORDING TO THE
Feature space

The weight vector is a linear combination of the feature vectors:

$$w = \sum_{i=1}^N \alpha_i d_i \Phi(x_i)$$

Thus the hyperplane equation is:

$$\sum_{i=1}^N \alpha_i d_i \underbrace{\Phi^T(x_i) \Phi(x)}_{\text{dot product}} = 0$$

Note the dot product is now in the feature space

Evaluating $\Phi(x)$ could be intractable

Kernel trick

Fortunately under certain conditions we do NOT need to evaluate directly $\Phi(x)$.

We do NOT even need to know the feature space itself!

This is possible using a function to compute directly the dot products in the feature space

Kernel: $k : \mathbb{R}^{m_0} \times \mathbb{R}^{m_0} \rightarrow \mathbb{R}$

k is called inner product kernel function

$$: k(x_i, x) = \Phi^T(x_i) \Phi(x)$$

Property: a kernel function is a symmetric function $k(x_i, x) = k(x, x_i)$

The kernel function provides us a powerful abstraction tool.

inner product
K(DIRECTIONS)
compute the product
 $\Phi^T(x_i) \Phi(x)$
between

SOFT MARGIN AND KERNEL

Primal Problem

Given the training set $T = \{(\Phi(\mathbf{x}_i), d_i)\}_{i=1}^N$ find the optimal value of \mathbf{w} which minimizes the objective function

$$\Psi(\mathbf{w}, \xi) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i$$

under the constraints

$$d_i(\mathbf{w}^T \Phi(\mathbf{x}_i)) \geq 1 - \xi_i \quad \forall i = 1, \dots, N$$

$$\xi_i \geq 0 \quad \forall i = 1, \dots, N$$

Please note that \mathbf{w} is now in the feature space and thus its dimensionality can lead to an intractable problem.

\mathbf{w} depends on dimensional Φ
 Φ can very high dimensional space

Dual Problem

Given the training set $T = \{(\Phi(\mathbf{x}_i), d_i)\}_{i=1}^N$ find the optimal values of $\{\alpha_i\}_{i=1}^N$ which maximize the objective function

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j d_i d_j \overbrace{k(\mathbf{x}_i, \mathbf{x}_j)}^{K_{i,j}}$$

satisfying the constraints

$$\sum_{i=1}^N \alpha_i d_i = 0$$

$$0 \leq \alpha_i \leq C \quad \forall i = 1, \dots, N$$

Not
SUSY

Remember C is a user specified non-negative parameter (regularization parameter)
The kernel includes the mapping to the feature space. The dependences on N becomes now an advantage (due to the possible high dim of the feature space)

$$\times \frac{1}{2\sigma^2} - \gamma$$

Architecture view of a SVM

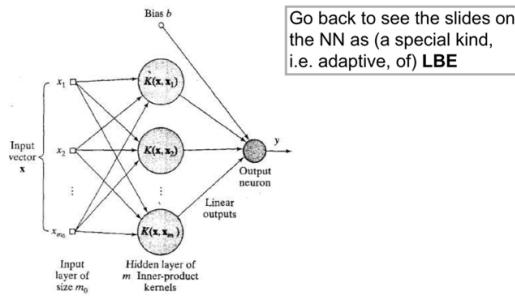


Figure: "Architecture" of a SVM.

- Number of hidden units = Number of support vectors
(... but not exactly the same unit, e.g. hidden units in NN have adaptive free parameters!)

This is just a "picture" view of the SVM. Think by yourself on the comparison in terms of complexity control, efficiency, and other characteristics.

Some examples of kernels

- Polynomial Learning Machine: $k(\mathbf{x}, \mathbf{x}_i) = (\mathbf{x}^T \mathbf{x}_i + 1)^p$ where p is a user specified parameter
- Radial Basis Function Net: $k(\mathbf{x}, \mathbf{x}_i) = e^{-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}_i\|^2}$ where σ^2 is a user specified parameter
A.k.a. Gaussian kernel
[narrow peaked Kernels (small sigma) imply that reply for \mathbf{x}_i is only d_i (what does it resemble?)]
- Two-layer perceptron $k(\mathbf{x}, \mathbf{x}_i) = \tanh(\beta_0 \mathbf{x}^T \mathbf{x}_i + \beta_1)$ where $\beta_0 (> 0)$ and $\beta_1 (< 0)$ are user specified parameters

Important

For polynomial learn. machines and RBN you always get an inner product kernel, while for a two-layer perceptron the Mercer's th. only holds for some choices of β_0 and β_1

Note

Also note that using a kernel function of RBF type always leads to a feature space with an infinite number of dimensions

SVM: non-linear regression

Regression: function that describes the dependency of a scalar from a vector

$$d = f(\mathbf{x}) + \nu$$

↑ ↓
 Unknown noise
 function

- Training set $\{(\mathbf{x}_i, d_i)\}_{i=1}^N$

- estimate the target d using a LINEAR EXPANSION of non-linear functions

$$\{\phi_j(\mathbf{x})\}_{j=0}^{m-1}$$

$$g = f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) \quad \text{where } \mathbf{w} = (w(0) = b, w(1), \dots, w(m-1))$$

transformation ↪ $\phi(\mathbf{x}) = (\phi_0(\mathbf{x}) = 1, \phi_1(\mathbf{x}) = x_1, \dots, \phi_{m-1}(\mathbf{x}))^T$

We have a function that describes the dependency of a scalar from a vector:

$$d = f(\mathbf{x}) + v$$

where:

- the function f is unknown
- the noisy term v is statistically independent of the input vector \mathbf{x}

We only have a training set $T = \{(\mathbf{x}_i, d_i)\}_{i=1}^N$

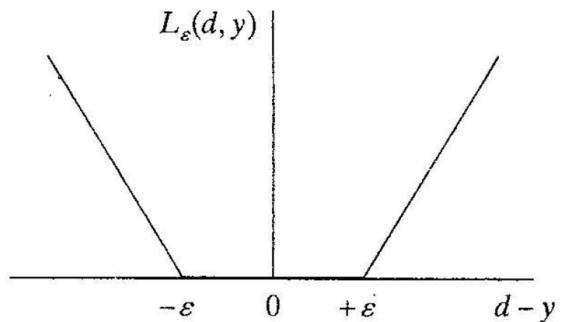
We estimate d using a linear expansion of non-linear functions $\{\phi_j(\mathbf{x})\}_{j=1}^{m_1}$

$$y = h(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x})$$

Where: $\mathbf{w} = (w(0) = b, w(1), \dots, w(m_1))^T$

$$\Phi(\mathbf{x}) = (\phi_0(\mathbf{x}) = 1, \phi_1(\mathbf{x}), \dots, \phi_{m_1}(\mathbf{x}))^T$$

$$L_\epsilon(d, y) = \begin{cases} |d - y| - \epsilon & \text{if } |d - y| \geq \epsilon \\ 0 & \text{otherwise} \end{cases}$$



ϵ -tube

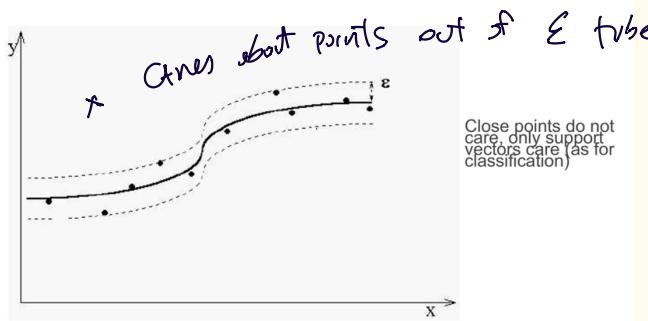


Figure: ϵ -tube around d .

Note: again only the support vectors guide the solution, as the cost function ignores any training data close to the model prediction, i.e. within a threshold ϵ (for class. were for points that lie beyond the margin)

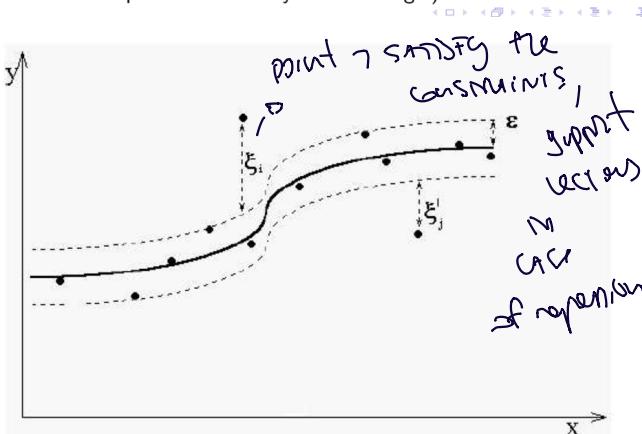


Figure: ϵ -tube and slack variables.

Formulating the optimization problem

Introduce the non-negative slack variables ξ'_i and ξ_i for all $i=1,\dots,N$

$$-\xi'_i - \epsilon \leq d_i - \mathbf{w}^T \Phi(\mathbf{x}_i) \leq \epsilon + \xi_i \quad \forall i = 1, \dots, N$$

This leads to the following constraints:

$$\left. \begin{array}{l} d_i - \mathbf{w}^T \Phi(\mathbf{x}_i) \leq \epsilon + \xi_i \\ \mathbf{w}^T \Phi(\mathbf{x}_i) - d_i \leq \epsilon + \xi'_i \\ \xi_i \geq 0 \\ \xi'_i \geq 0 \end{array} \right\} \forall i = 1, \dots, N$$

Primal Problem

Primal Problem

Given the training set $T = \{(\mathbf{x}_i, d_i)\}_{i=1}^N$ find the optimal values of \mathbf{w} such that the following objective function is minimized

$$\Psi(\mathbf{w}, \xi, \xi') = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N (\xi_i + \xi'_i)$$

under the constraints

$$\left. \begin{array}{l} d_i - \mathbf{w}^T \Phi(\mathbf{x}_i) \leq \epsilon + \xi_i \\ \mathbf{w}^T \Phi(\mathbf{x}_i) - d_i \leq \epsilon + \xi'_i \\ \xi_i \geq 0 \\ \xi'_i \geq 0 \end{array} \right\} \forall i = 1, \dots, N$$

Computing the weight vector

Solving the dual problem we obtain the optimal values of the Lagrangian multipliers $\{\alpha_i\}_{i=1}^N$ and $\{\alpha'_i\}_{i=1}^N$

Then we can compute the optimal value of vector \mathbf{w} :

$$\mathbf{w} = \sum_{i=1}^N (\alpha_i - \alpha'_i) \Phi(\mathbf{x}_i)$$

↑
 $\mathbf{w} = \sum_{i=1}^N \gamma_i \Phi(\mathbf{x}_i)$

Computing the estimate

The estimate function is defined as $h(\mathbf{x}) = y = \mathbf{w}^T \Phi(\mathbf{x})$

Using the linear expansion for \mathbf{w} we get:

$$h(\mathbf{x}) = \sum_{i=1}^N \gamma_i \Phi^T(\mathbf{x}_i) \Phi(\mathbf{x}) = \sum_{i=1}^N \gamma_i k(\mathbf{x}_i, \mathbf{x})$$

Support vector

In this case support vectors correspond to non-zero values of γ_i

BIAS-VARIANCE: investiga gli errori e in punto guardando le realizzazioni del TUE set
• la decomposizione dell'errore ($y - \hat{h}(x)$) in un punto x ; 3 componenti account

1) BIAS: quantifica la discrepanza tra TRUE FUNCTION e $\hat{h}(x)$
(se H è piccolo, il BIAS è ALTO) (MEAN $\hat{h}(x)$ over possible set)

2) VARIANCE: quantifica la VARIABILITÀ della risposta del modello \hat{h}
per diverse realizzazioni del training set (dove all'elenco disponibili)
 \Rightarrow flexible model if change the set \Rightarrow TRAIN ≠ TEST

3) NOISE: le label includono un RANDOM ERROR (if given x there are more possible d, targets)

- assume Regression Code

\rightarrow BIAV-VAR ANALYSIS

• given a data point x WHAT IS THE EXPECTED PREDICTION ERROR? error at that point
• $E_p[y - \hat{h}(x)]^2$ \Rightarrow DECOMPOSE IT IN THE 3 COMPONENTS
→ $E_p[y - \hat{h}(x)]^2 = E_p[\hat{h}(x) - E_p[\hat{h}(x)]]^2 + E_p[(y - \hat{h}(x))^2]$ averaged over possible realizations of the set

Avg when you chose the m set

• let be Z una variabile casuale con valori $z_i = 1, \dots, l$ with prob. distn. $P(Z)$

• la MEDIA (\circ expected value): $\bar{Z} = E_p[Z] = \sum_{i=1}^l z_i P(z_i)$ center

• la VARIANZA: $\text{VAR}[Z] = E[Z - \bar{Z}]^2 = E[Z^2] - \bar{Z}^2$ $E[Z^2] = \bar{Z}^2 + \text{VAR}[Z]$

$$E_p[y - \hat{h}(x)]^2 = E_p[\hat{h}(x)^2 - 2y\hat{h}(x) + y^2] = \\ = E_p[\hat{h}(x)^2] + E_p[y^2] - 2E_p[y]\bar{E_p}[\hat{h}(x)] =$$

• let $\bar{h}(x) = E_p[\hat{h}(x)]$

↓
leme $E_p[\hat{h}(x)^2] = E_p[(\hat{h}(x) - \bar{h}(x))^2 + \bar{h}(x)^2]$
 $\hat{h}(x) \equiv h(x)$

$$\Rightarrow E_p[y^2] = E_p[(y - f(x))^2] + f(x)^2$$

$$= E_p[\hat{h}(x) - \bar{h}(x)]^2 + \bar{h}(x)^2 - 2E_p[y]E_p[\hat{h}(x)] +$$

$$+ E_p[(y - f(x))^2] + f(x)^2 =$$

$$= E_p[(\hat{h}(x) - \bar{h}(x))^2] +$$

$$+ (\bar{h}(x) - f(x))^2 +$$

$$+ E_p[(y - f(x))^2]$$

VARIANCE

BIAS²

NOISE²

$$\text{cuz } E_p[g] = f(x) \\ = E_p(f(x) + \epsilon)$$

$$= \text{VAR}[h(x)] + \text{BIAS}[h(x)]^2 + \mathbb{E}_p[\epsilon^2]$$

$$= \text{VAR}[h(x)] + \text{BIAS}[h(x)]^2 + \sigma^2$$

$$\text{EXPECTED PREDICTION ERROR} = \text{VARIANCE} + \text{BIAS}^2 + \text{NOISE}^2$$

ENSEMBLE LEARNING: bagging and boosting

Idea: use $\frac{1}{n}$ of TR set or random initialization of network

\Rightarrow models with high variance

But the output mean is good!
(not the mean error)

\Rightarrow ensemble learning

Recall Regularization:

$$\text{Loss}(\mathbf{w}) = \sum_p (d_p - o(\mathbf{x}_p))^2 + \lambda \|\mathbf{w}\|^2$$

$\lambda \uparrow \Rightarrow \text{BIAS} \uparrow, \text{VARIANCE} \downarrow$

Left (data) term (TR data) Penalty term

Varying λ we can obtain

- less regularized (complex) solutions (low λ)
- more regularized (less complex) solutions (high λ)

- The following example (from Bishop) (should be) using RBF network with LMS

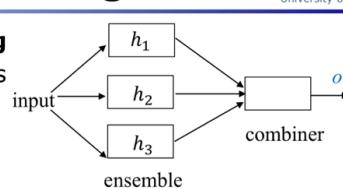
(By the way, note: that the intercept is typically left out of the penalty term)

20

ENSEMBLE LEARNING

Dip. Informatica
University of Pisa

- A first look to ensemble learning
- Take advantage of multiple models



Simplest case ("voting schema"):

- Regression
 - simple average committee $o(\mathbf{x}) = \frac{1}{K} \sum_{i=1}^K h_i(\mathbf{x})$
 - $\text{loss}(E_i(h_i)) \leq E_i(\text{loss}(h_i))$ i.e.
 - (mean₀) square error of a committee \leq mean of the (mean₀) square error
e.g. $(t-o)^2, t=5, 2 \text{ models} \rightarrow (5-3)^2 = 4 < ((5-2)^2 + (5-4)^2)/2 = 5$
or even: $(5-5)^2 = 0$ while $((5-6)^2 + (5-4)^2)/2 = 1$
- Classification
 - Take a vote over many classifier (given they provide "diversified" responses): why? (or classify after the mean of the continuous outputs)

Alternative ("stacking schema"): the combiner is again a ML model

△ Mirhali

27

Send all the outputs to a combiner
- take the mean
- produce output
• the loss of the committee
(mean of h_i)
is always \leq loss of each
AUG_i:
 $\text{Loss}(\text{mean}(h_i)) \leq$ reg / clas
~~loss(h_i)~~
 $\text{AUG}_i(\text{loss}(h_i))$

SOM - K-MEANS

- UNSUPERVISED NN : w Teacher
- TR = $\langle d \rangle$
- clustering \rightarrow VECTOR QUANTIZATION for partitioning clustering
- dimensional reduction / visualization (PCA)
- SIMILARITY MEASURE: dipende dalle metriche (es. euclidean distance)
- VECTOR QUANTIZATION: Tecnica per l'ENCODING di DATA MANIFOLD (Projected in 2D) (riduzione le dimensioni)
 - un submanifold $V \subseteq \mathbb{R}^N$ using un finite set $w = (w_1, \dots, w_k)$ of reference vectors $w_i \in \mathbb{R}^m$ (use just finite of this prototype per rappresentare i dati)
 - un data vector $x \in V$ il BEST MATCHING o' winning reference vector $w_i * (x)$ ($*$: è il vettore per cui $x \in V$ e w_i è il reference vector (cluster center) che rappresenta x)
 - di w per cui il DISTORTION error $d(x, w_i * (x))$ è minima (il vettore più vicino)
 - QUESTA PROCEDURA DIVIDE il MANIFOLD V in sub-regioni $V_i \Rightarrow V_i = \{x \in V \text{ s.t. } \|x - w_i\| \leq \|x - w_j\| \forall j\}$
- QUANTIZATION: from continuous to discrete values
- clustering has interest/useful group of data (to always reduce the dimensionality)
- Goal: PARTIZIONE OTTIMA DI UNA DISTRIBUZIONE SCONOSCIUTA del x-input space in regions (cluster) approssimate da un cluster centroid (o prototype)
 - Assegna ogni punto al suo cluster center
$$x \rightarrow c(x) \text{ or } w_i * (x) \quad (\text{prototype})$$
- considera SQUARED EUCLIDEAN DISTANCE

$$d(x_i, c(x_i)) = \|x_i - c(x_i)\|^2$$

w: prototype vectors
↑
prototype center

$$E = \int f(d(x, w_{i^*(x)})) \cdot p(x) dx = \int \|x - w_{i^*(x)}\|^2 p(x) dx$$



Quantization error function

$$F = \sum_i^l \sum_j^K \|x_i - w_j\|^2 \delta_{winner}(i, j)$$

$$E = \int f(d(x, w_{i^*(x)})) p(x) dx = \int \|x - w_{i^*(x)}\|^2 p(x) dx$$

from continuous to discrete
↓ (squared error criterion)

$$E = \sum_i^l \sum_j^K \|x_i - w_j\|^2 \delta_{winner}(i, j) \quad \text{Discrete version}$$

$\delta_{winner}(i, j)$ characteristic function of the receptive field of w_j
 1 if j is the winner for x_i , 0 otherwise

$p(x)$ is the probability distribution of x

$$\begin{aligned} E &= \int f(d(x, w_{i^*(x)})) \cdot p(x) dx = \int \|x - w_{i^*(x)}\|^2 p(x) dx = \\ &= \sum_i^l \sum_j^K \|x_i - w_j\|^2 \cdot \delta_{winner}(i, j) \end{aligned}$$

- need s_j that is reception field of $w_j \Rightarrow 1$ if j is winner
for x_i
0 else
- (cut to one
so b rel cluster)