

risk function $R = \int L(d, h(x)) d P(x, d)$

d
target
 P

approximates the
prob. distribution

$$R_{\text{emp}} > \frac{1}{e} \sum_{i=1}^e (d_p - h(X_p))^2 \quad L = \text{card}(YR)$$

VC-dim: it labels with probability $1-\delta$

$$R \leq R_{\text{emp}} + \underbrace{E[\epsilon]}_{\substack{\text{dim} \\ \text{of target}}} \left(\frac{1}{\delta}, \text{VC}, \frac{1}{\epsilon} \right)$$

VC-confidence

VC-dim: numero delle componenti di H plausibili del modello

- ϵ is a function that grows with VC (VC-dim), that decreases with (higher) δ and ϵ .
- We know that R_{emp} decreases using complex models (with high VC-dim) (e.g. the

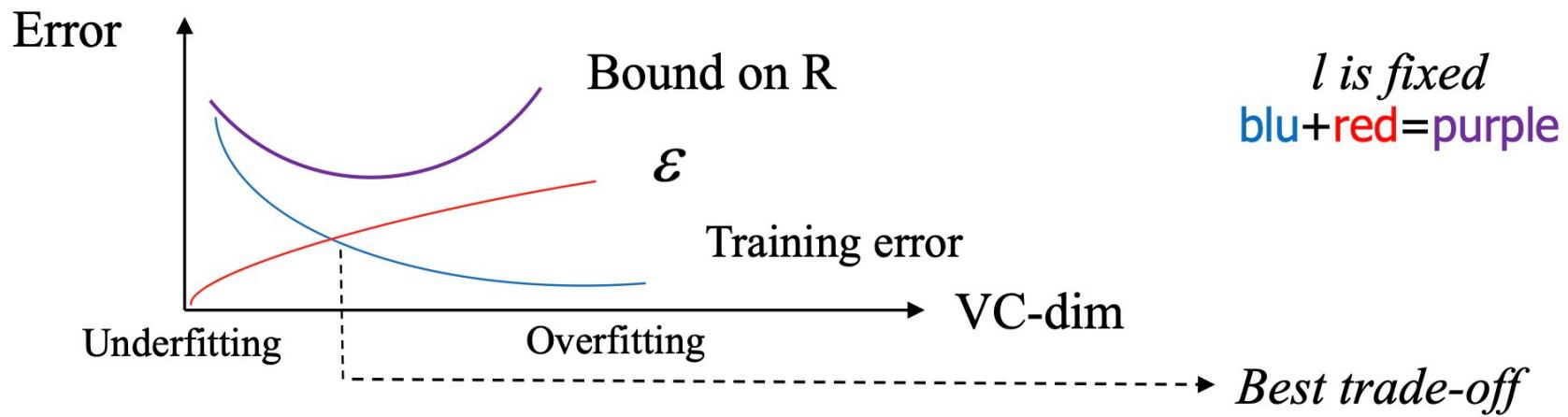
polynomial degree in the example)

- delta is the confidence, it rules the probability that the bound holds (e.g. low delta 0.01 → the bound holds with probability 0.99)

Se \cup_c è grande, ma R_{emp} diminuisce \Rightarrow rischio overfitting

- Higher l (data) → lower VC confidence and bound close to R
- Too simple model (low VC-dim) can be not suff. due to high R_{emp} (underfitting)
- Higher VC-dim (fix l) → lower R_{emp} but VC-conf. and hence R may increase (overfitting)

Structural risk minimization: minimize the bound !



Loss con Tikhonov

$$L(\omega) = \sum_{P=1}^l (y_p - x_p^\top \omega)^2 + \lambda \|\omega\|^2 \rightarrow \sum_i \omega_i^2$$

- 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 value → minimize the loss the focus is on the second term, hence the data error (first term) could grow too much, i.e. moving to underfitting

Perceptron learning Algorithm

- 1 Inizializzo i pesi ω ad un vettore prossimo a zero
- 2 scelgo learning rate η
- 3 fino a condizione di stop non rottinfilto:

Delta rule $w_{new} = w + \frac{1}{2}\eta(d - out) x$ $d \in \{-1, +1\}$, $out = \text{sign}(w^\top x)$

\hookrightarrow verso y_2

Perceptron Convergence Theorem

Th: For linearly separable tasks, the perceptron algorithm converges after a finite number of steps.

Idea: find lower & upper bound to $\|w\|^2$ as a function of q steps. ✓ and q steps ✓
 lower b. upper b.

annun $w(0) = 0$, $\eta = 1$, q annun di pattern sbagliati (tutti hlin negativi)

- Lower bound on $|w(q)|^2$: $w^* \cdot w(q) = w^* \cdot \sum_{j=1}^q x_{ij} \geq q \alpha$ $\alpha = \min_i (w^* x_i)$

now Square

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

- Upper bound

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

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

\leftarrow alla q^{th} iterazione

iterazioni

$$\|\omega(q)\|^2 \leq \sum \|x_{ij}\|^2 \leq q\beta$$

$\beta = \max_i \|x_{il}\|$

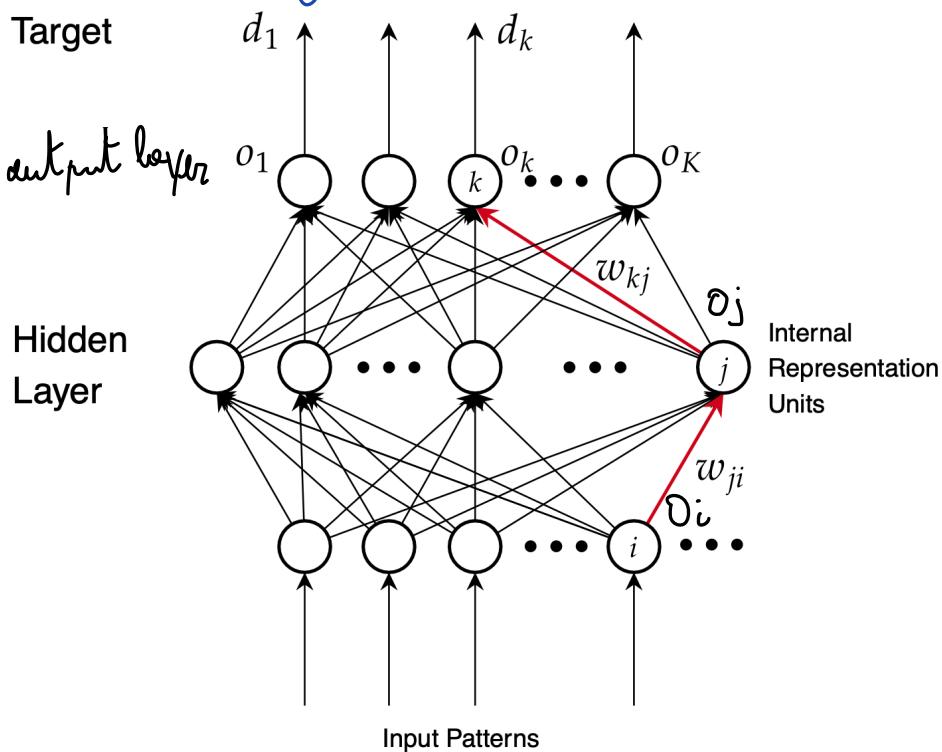
$$q\beta \geq \|\omega(q)\|^2 \geq (q\alpha)^2 / \|\omega^*\|^2$$

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

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

$$q \leq \beta/\alpha^2$$

Backpropagation algorithm



- initialize all the weights w in the network and η

- while $E_{\text{tot}} \geq E$:

$$\forall w \in W : \Delta_w = \frac{\partial E_{\text{tot}}}{\partial w}$$

$$w_i^{\text{new}} = w_i + \eta \Delta w_i + \dots$$

Compute E_{tot} and out

$$\bar{E}_{\text{tot}} = \sum_p \bar{E}_p, \quad E_p = \frac{1}{2} \leq (d_x - o_x)^2$$

$$\Delta w = - \frac{\partial \bar{E}_{\text{tot}}}{\partial w} = - \sum_p \frac{\partial \bar{E}_p}{\partial w} = \sum_p \Delta_p w$$

$$\Delta_p w_{t,u} = - \frac{\partial E_p}{\partial w_{t,u}} = - \underbrace{\frac{\partial E_p}{\partial \text{net}_t}}_{\text{red box}} \cdot \underbrace{\frac{\partial \text{net}_t}{\partial w_{t,u}}}_{\text{green box}} = f_t \cdot o_u$$

$$\begin{cases} \text{net}_t = \sum_j w_{t,j} o_j \\ o_t = f_t(\text{net}_t) \end{cases} \Rightarrow \frac{\partial \text{net}_t}{\partial w_{t,u}} = \frac{\partial \sum_j w_{t,j} o_j}{\partial w_{t,u}} = o_u \quad (o \text{ per } j \neq u)$$

$$f_t = - \frac{\partial E_p}{\partial \text{net}_t} = \underbrace{- \frac{\partial E_p}{\partial o_t}}_{\text{blue box}} \cdot \frac{\partial o_t}{\partial \text{net}_t} = o_t = f_t(\text{net}_t) = - \underbrace{\frac{\partial E_p}{\partial o_t}}_{\text{blue box}} \cdot f'_t(\text{net}_t)$$

distinguishing the error of the output layer or hidden layer

$t = k$ (output level)

$$\frac{-\partial E_p}{\partial o_k} = -\frac{\partial \frac{1}{2} \sum_{i=1}^K (d_i - o_i)^2}{\partial o_k} = (d_k - o_k) \quad \delta_k = (d_k - o_k) f'_k(\text{net}_k)$$

$t = j$ (hidden unit)

$$\frac{-\partial E_p}{\partial o_j} = \sum_{k=1}^K \left(\frac{-\partial E_p}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial o_j} \right) = \sum_{k=1}^K \delta_k \cdot w_{kj}$$
$$\frac{\partial \sum_s w_{ks} \cdot o_s}{\partial o_j} = w_{kj}$$

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

$$w_{t+1}^{\text{new}} = w_{tu} + \eta \cdot \delta_t \cdot o_u$$

Momentum

$$\Delta w_{\text{new}} = -\eta \frac{\partial E(w)}{\partial w} + \alpha \Delta w_{\text{old}}$$

$$w_{\text{new}} = w_{\text{old}} + \Delta w_{\text{new}}$$

Nesterov Momentum

1 $\bar{w} = w + \alpha \Delta w_{\text{old}}$

2 evaluate the new gradient at this interior point (\bar{w})

$$\Delta \bar{w} = -\eta \frac{\partial E(\bar{w})}{\partial w} + \alpha \Delta w_{\text{old}}$$

3 $w_{\text{new}} = \bar{w} + \Delta \bar{w}$

- Shown to improve the rate of convergence for the batch mode (not for the stochastic case!)

Moving Weight Decay and Momentum

$$\Delta w_{tu} = -\gamma \frac{\partial \text{loss}(w)}{\partial w_{tu}} + \lambda \Delta w_{old,tu} = \gamma \delta_t o_u - \lambda w_{tu} + \lambda \Delta w_{old,tu}$$

Poniamo tenere stoccati i *iperperimetri* arrivando

$$\Delta w_{tu} = \eta \delta_t o_u + \alpha \Delta w_{old,tu}$$

$$w_{new,tu} = w_{tu} + \Delta w_{tu} - \lambda w_{tu}$$

The Corrade Correlation algorithm

- Start with N0 network, a network without hidden units: Training of N0. Compute error for N0. If N0 cannot solve the problem: go to N1.
- In the N1 network a hidden unit is added such that the correlation between the output of the unit and the residual error of network N0 is maximized (by training its weights, see next slide).

- After training, the weights of the new unit are frozen (they cannot be retrained in the next steps) and the remaining weights (output layer) are retrained.
- If the obtained network N1 cannot solve the problem, new hidden units are progressively added which are connected with all the inputs and previously installed hidden units.
- The process continues until the residual errors of the output layer satisfy a specified stopping criteria (e.g the errors are below a given threshold).

$$S_k = \left| \sum_p (O_p - \bar{O}) (\bar{G}_{p,k} - \bar{G}) \right|$$

$$O_p = f(\text{net}_p)$$

$$\frac{\partial |S_k|}{\partial w_j} = \text{sign}(S_k) \frac{\partial S_k}{\partial w_j} = \text{sign}(S_k) \sum_p (\bar{G}_{p,k} - \bar{E}) \frac{\partial}{\partial w_j} (O_p - \bar{O}),$$

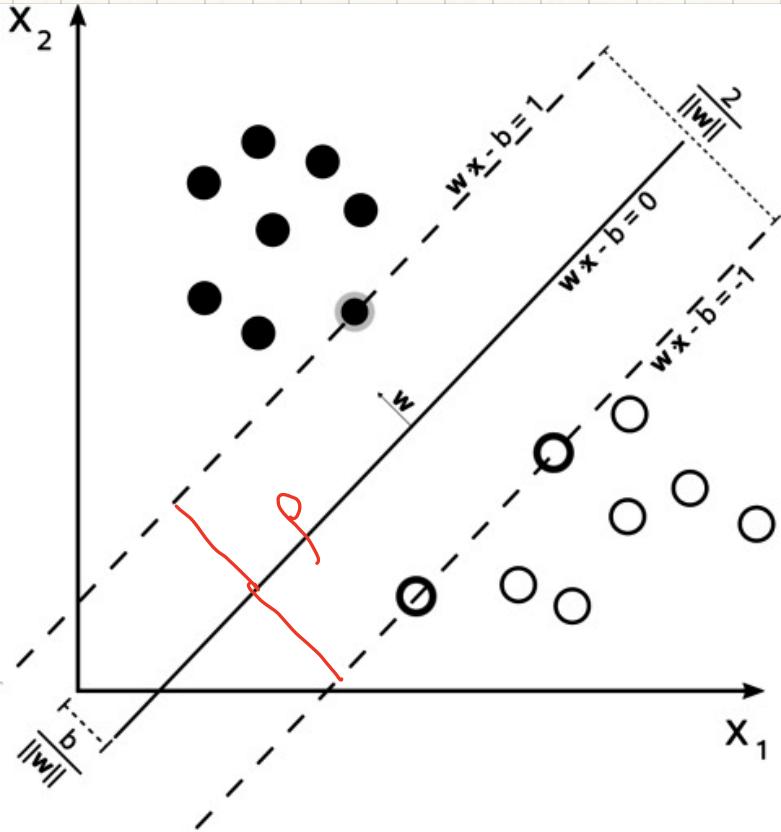
$$\frac{\partial}{\partial w_j} O_p = \frac{\partial}{\partial w_j} f(\text{net}_p) = \frac{\partial f}{\partial \text{net}_p} \cdot \frac{\partial \text{net}_p}{\partial w_j} = f'(\text{net}_p) \cdot x_{pj}$$

$$\min_p (S_k) \leq (E_{P,k} - \bar{E}) f'(w_{t_p}) x_{P_j}$$

- Definition: The VC dimension of a class of functions H is the maximum cardinality of a set of points (configuration) in X that can be shattered by H

$$R \leq R_{emp} + \epsilon(\gamma_e, V_c, \gamma_f)$$

SVM



Dato un punto x nello spazio:

$$x = x_p + \frac{2 \frac{w_0}{\|w_0\|}}{\|w_0\|}$$

prelevare punto in
una aletta

distr. origine - iperplane

$$y(x) = w_0^\top x + b_0$$

$$g\left(x_p + \frac{\omega_0}{\|\omega_0\|}\right) = \omega_0^\top \left(x_p + \frac{\omega_0}{\|\omega_0\|}\right) + b_0 = \omega_0^\top x_p + \omega_0^\top \frac{\omega_0}{\|\omega_0\|} + b_0 =$$

$$\underbrace{\omega_0^\top x_p + b_0}_{g(x_p) = 0} + \frac{1}{2} \|\omega_0\| \Rightarrow g(x) = \frac{1}{2} \|\omega_0\|$$

Ale portys $g(x) = 1$

$$\rho = \frac{1}{\|\omega_0\|} = \frac{\rho}{2} \quad \rho = \frac{1}{\|\omega_0\|}$$

upisanje $\omega_0^\top x + b_0 = 0 \Rightarrow \text{minimum } f = \frac{1}{2} \|\omega_0\|^2 \Rightarrow \text{minimum } \|\omega_0\|$

Hard margin SVM primal form

$$T = \{(x_i, y_i)\}_1^N$$

$$\min_{\omega, b} f(\omega) = \frac{1}{2} \omega^\top \omega \quad \rightarrow \|\omega\|^2$$

$$y_i(\omega^\top x_i + b) \geq 1 \quad i=1, \dots, N$$

Dilemma Leyerengieno

Condizioni necessarie di Karush-Kuhn-Tucker

$$\nabla_x L(x^*, \lambda^*, \mu^*) = 0$$

$$g(x^*) \leq 0, \quad h(x^*) = 0$$

$$\lambda^{*T} g(x^*) = 0$$

$$\lambda^* \geq 0.$$

$$\lambda = \lambda$$

Duale Leyerengieno svolto

$$\min_{\lambda \geq 0} \frac{1}{2} \sum_i \sum_j (\lambda_i \lambda_j y_i y_j x_i^T x_j) - \sum_i \lambda_i$$

$$\sum_i \lambda_i y_i = 0 \quad \lambda_i \geq 0 \quad \forall i$$

Risolvendo il duale Leyerengieno

$$L(\lambda^*, w^*, b^*) = \frac{1}{2} w^T w - \sum_i \lambda_i (1 - y_i (w^T x_i + b^*))$$

$$w^* = \sum_1^n \lambda_i y_i x_i$$

$\lambda > 0$, allora x_i vettore di supporto
 x_i non è vettore di supporto, allora $\lambda = 0$

Ullo fine bisce trovare i λ per risolvere il problema

$$f^* = \lambda - w^T x^{(s)}$$

$$f^* = 1 - \sum_1^n \lambda_i y_i x_i^T x^{(s)}$$

In questo si ha la relazione ottima: $w^{*T} x + f^* = 0 \iff \sum_1^N \lambda^* y_i x_i^T x + f^* = 0$

le misure dipendono interamente dai vettori di supporto

Theorem - Vapnik

Let D be the diameter of the smallest ball around the data points $\mathbf{x}_1, \dots, \mathbf{x}_N$.

For the class of separating hyperplanes described by the equation $\mathbf{w}^T \mathbf{x} + b = 0$ the upper bound to the VC dimension is

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

dimensione del input space

since $\frac{D^2}{\rho^2} = \text{Radius}^2 \|\mathbf{w}\|^2$, VC can be less than $m_0 + 1$ computed for general hyperplanes by constraining the search to the "regularized" one with maximum margin

Soft margin primal problem

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

$$\begin{aligned} \forall i \quad (\mathbf{w}^T \mathbf{x}_i + b) &\geq 1 - \xi_i & \forall i = 1, \dots, N \\ \xi_i &\geq 0 & \forall i = 1, \dots, N \end{aligned}$$

- Low $C \rightarrow$ many TR errors allowed \rightarrow possible misfitting
- High $C \rightarrow$ no TR errors allowed \rightarrow smaller margin \rightarrow possible overfitting

Lagrange dual:

$$\max_{\lambda \geq 0, \mu \geq 0} \min_{w, b, \xi} \frac{1}{2} w^T w + C \sum_i \xi_i + \sum_i \lambda_i (1 - \xi_i - y_i (w^T x_i + b)) - \sum_i \underline{\mu_i \xi_i}$$

Dual Problem

Given the training set $T = \{(\Phi(x_i), y_i)\}_{i=1}^N$ find the optimal values of $\{\lambda_i\}_{i=1}^N$ which maximize the objective function

$$Q(\lambda) = \sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j=1}^N \lambda_i \lambda_j y_i y_j \overbrace{k(x_i, x_j)}^{K_{i,j}}$$

satisfying the constraints

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

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

seleziona C , è la funzione Kernel. Calcola la matrice Kernel K . Risolvere il problema doppio e calcolo λ_i . Calcola b (bias) risolvendo K e λ_i .

Atto un nuovo input x calcolo $\Rightarrow w^T \phi(x) = \sum_{i=1}^N \lambda_i y_i K(x, x_i)$ e calcolo il segno di $w^T \phi(x)$

Polynomial learning machine $k(x, x_i) = (x^T x_i + 1)^P$

Radial Basis Function Net: $k(x, x_i) = e^{-\frac{1}{2\sigma^2} \|x - x_i\|^2}$

Two layer perceptron $k(x, x_i) = \tanh(\beta_0 x^T x_i + \beta_1)$ where $\beta_0 > 0$ and $\beta_1 < 0$

SVM per le regressioni

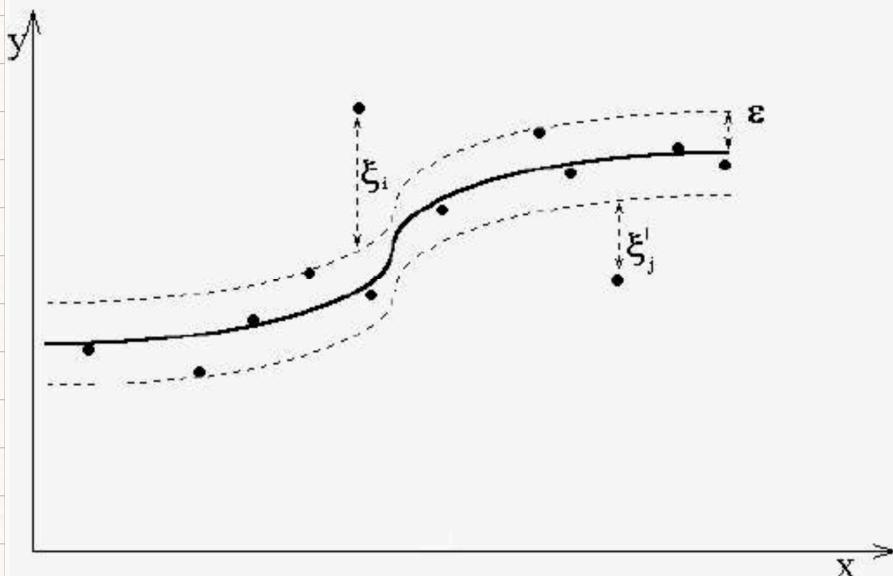
Primal Problem

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

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



$$w = \sum_i^N (\lambda_i - \lambda'_i) \phi(x_i) \text{ formulazione duale}$$

Notiamo in accordo all'input x nei minimi se volere si $f(x)$ con

$$h(x) = \sum_1^N \gamma_i k(x_i, x)$$

Riassumendo

- ➊ Important: we firstly have to select values for the user specified parameters C and ϵ
- ➋ Important: we also have to choose an inner product kernel function k
- ➌ We compute the kernel matrix K
- ➍ We solve the optimization problem in the dual form and get the optimal values of the Lagrangian multipliers (we get the values of $\{\gamma_i\}_{i=1}^N$)
- ➎ We compute the optimal value for the bias
- ➏ We obtain the estimate function as a linear combiner of dot products in a feature space we can ignore $\rightarrow h(x) = \sum_{i=1}^N \gamma_i k(x_i, x)$

Bias variance decomposition

- Expected error of mean of \hat{z} : $\bar{z} = E_p[\hat{z}] = \sum_i^l z_i P(z_i)$
- Variance of \hat{z} : $\text{Var}[\hat{z}] = E[(\hat{z} - \bar{z})^2] = E[\hat{z}^2] - \bar{z}^2 \Rightarrow E[\hat{z}^2] = \text{Var}[\hat{z}] + \bar{z}^2$

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

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

$$E_p[y] = \bar{y} = E_p[\rho(x) + \epsilon] = \rho(x) \Rightarrow E_p[y^2] = E[(y - \rho(x))^2] + \rho(x)^2$$

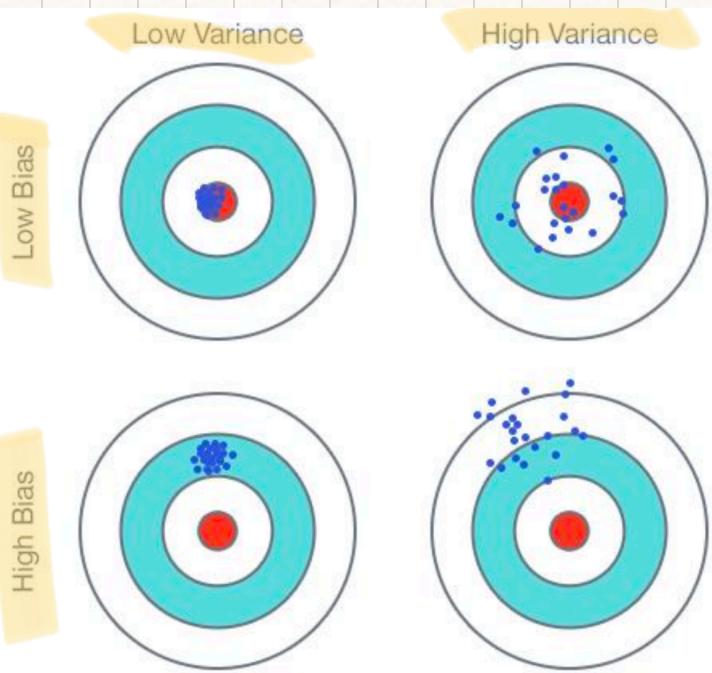
unica realizzazione

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

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

Variance + Bias 2 + $E[\text{noise}^2]$

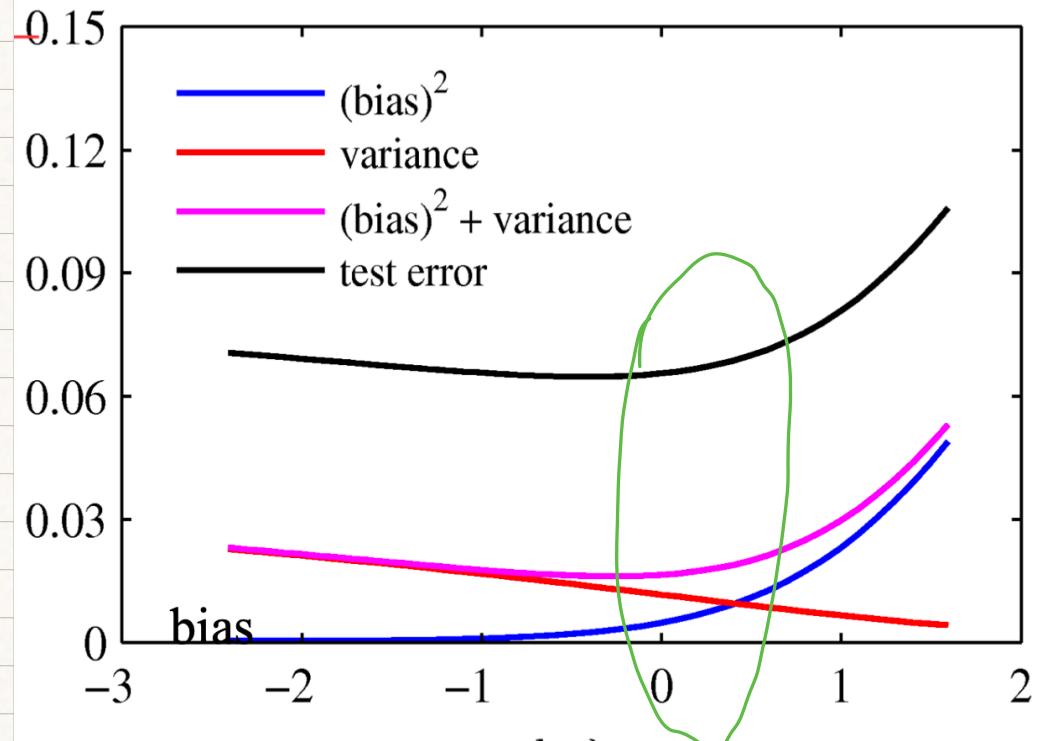
$$\text{Var}[h(x)] + \text{Bias}[h(x)]^2 + \sigma^2$$



Nel caso di regolarizzazione (weight decay)

Il diminuire di λ limitisce le varianze e aumenta il bias.

Il diminuire di λ aumenta le varianze del modello prelutto e limitisce il bias.



Higher complexity

$\ln \lambda$

optimal zone

Underfitting

SOM & k-means

Vector quantization

a submanifold $V \subset \mathbb{R}^m$, induced by reference vectors $w = (w_1, \dots, w_k)$,
 $w_i \in \mathbb{R}^m \quad i=1\dots k$

neighborhood $V \Rightarrow V_i := \{x \in V \text{ s.t. } \|x - w_i\| \leq \|x - w_j\| \forall j\}$
 ↓
 Voronoi polyhedra

each data vector x is described by the corresponding reference vector w_i

Quantization error function

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

(squared error criterion)

probability distribution

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

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

Lloyd and MacQueen's well-known K-MEANS clustering algorithm for any x_i :

$$\frac{\partial E}{\partial w_{i^*}} = \Delta w_{i^*} = \eta \delta_{\text{winner}}(i, i^*) (x_i - w_{i^*}) \cdot 2$$

Learning rate

- The set of reference vectors w that minimizes E is the solution of the VQ problem

SOM Competitive stage

- The winner is $i^*(x) = \arg \min_i \|x - w_i\|$
 w_i weight of unit i
 i index of the unit on the map

The current input vector x is compared with all the unit weights using an Euclidean distance criteria.

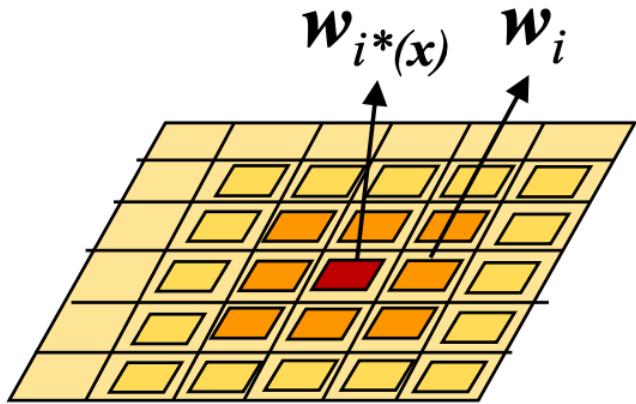
SOM Cooperative stage

$$w_i(t+1) = w_i(t) + \eta(t) h_{i, i^*(x)}(t) [x - w_i(t)]$$

The weight of the winning unit and the weight of the units in its neighborhood are moved closer to the input vector (Hebbian learning). At iteration t:

$$h_{i,i^*}(t) = \exp\left(-\frac{\|q_i - q_{i^*}\|^2}{2\sigma_{nh}^2(t)}\right)$$

coordinates of unit i
 } width



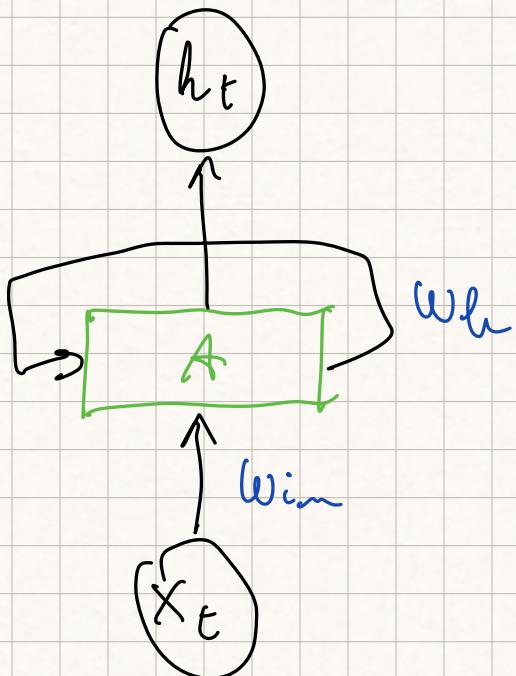
$h_{i,j}(t)$ is a function that decreases monotonically for increasing $\|i-j\|$.

- Learning rate (η) and neighbourhood radius values (σ) are decreased as function of iteration t

RNN

State transition system

$$\begin{cases} x(t) = \varphi(l(t), x(t-1)) \text{ & word RNN} \\ y(t) = g(x(t), l(t)) \end{cases}$$



$$f_t(h_{t-1}, x_t) = w_h h_{t-1} + w_{in} x_t + b_h$$

$$h_t = \tanh(f_t)$$

$$y_t = p(w_{out} h_t + b_{out})$$

