

10/14/2013

Exam 1: 11/12,
11:00-12:10,
10668 ENGINEERING HALL

TODAY: RBF

RECALL: GENERALIZED LINEAR CLASSIFIER

$$y_k(x^n) = \sum_{j=0}^M w_{kj} \phi_j(x^n) \quad \phi_0(x) \equiv 1$$

$\phi_j(x^n)$ = "BASIS FUNCTION"

RADIAL BASIS FUNCTION \equiv A BASIS FUNCTION S.T.

$$\phi_j(x^n) = \phi(\|x^n - \mu_j\|) \quad \text{FOR SOME } \mu_j$$

① ENGINEER CHOOSES ϕ BASED ON INTUITION

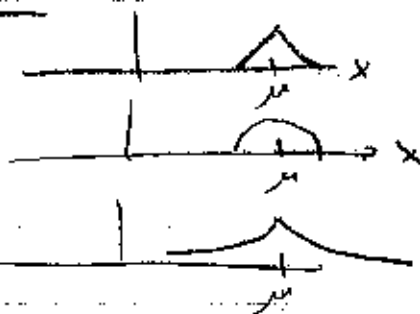
EXAMPLES $\phi(x) = \max(0, 1 - \lambda \|x - \mu\|)$

$$\phi(x) = \max(0, 1 - \frac{\|x - \mu\|^2}{\sigma^2})$$

$$\phi(x) = e^{-\lambda \|x - \mu\|}$$

$$\phi(x) = e^{-\frac{1}{2} \frac{\|x - \mu\|^2}{\sigma^2}}$$

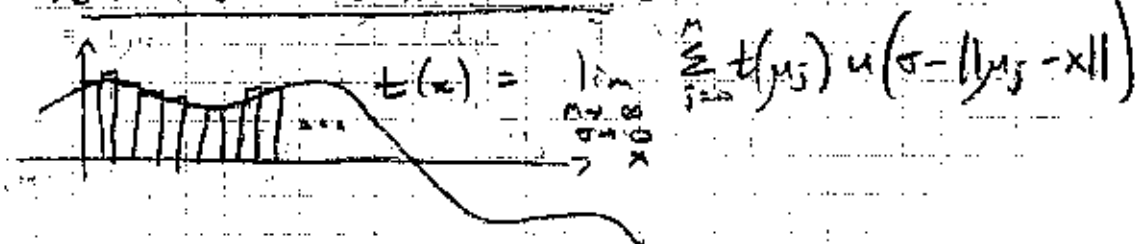
← MOST COMMON



DIGRESSION: ANY TARGET FUNCTION $t(x)$ CAN BE

WRITTEN
$$t(x) = \lim_{\substack{M \rightarrow \infty \\ \sigma \rightarrow 0}} \sum_{j=0}^M w_{kj} \phi_j(x)$$

PROOF: NEWTON'S INTEGRAL THEOREM



② GIVEN Φ , COMPUTER FINDS $\vec{\mu}_j, w_{kj}$

E.G.

$$\vec{\mu}_j, w_{kj} = \operatorname{argmin} \sum_{k=1}^N \sum_{j=1}^C (t_k^j - y_k^j)^2$$

③ FIND w_{kj} GIVEN FIXED $\vec{\mu}_j$: ANALYTIC SOLUTION EXISTS (PSEUDO-INVERSE)

$$w_{kj} = \operatorname{argmin} \frac{1}{2} \sum_{k=1}^N \sum_{j=1}^C \left(t_k^j - \sum_{j=0}^M w_{kj} \Phi_j(x^k) \right)^2$$

$$T [N \times C] \quad T(n, k) = t_k^n = [\bar{t}_1, \dots, \bar{t}_C]$$

$$\Phi [N \times (M+1)] \quad \Phi(n, j+1) = \Phi_j(x^n)$$

$$W [C \times (M+1)] \quad W(k, j+1) = w_{kj} = \begin{bmatrix} w_{k1} \\ \vdots \\ w_{kM} \end{bmatrix}$$

$$E = \frac{1}{2} \sum_{k=1}^N \|\bar{t}_k - \Phi W_k\|^2$$

$$\nabla_{W_k} E = \Phi^T (\Phi W_k - \bar{t}_k) = 0$$

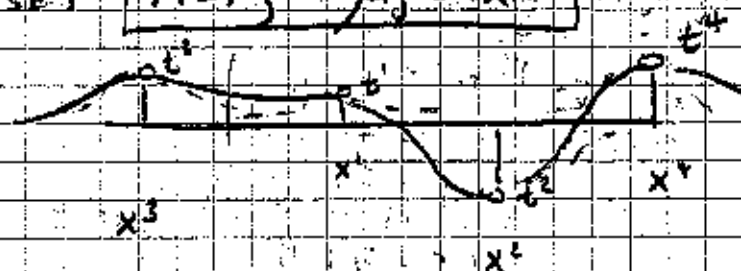
$$\text{AT } \boxed{W_k = (\Phi^T \Phi)^{-1} \Phi^T \bar{t}_k}$$

$$\boxed{W^T = (\Phi^T \Phi)^{-1} \Phi^T T}$$

④ FIND μ_j : THREE METHODS COMMON

① PARZEN WINDOWS / KERNEL ESTIMATOR

SET $M=N, \mu_j = x^j$, $y(x) = \sum_{j=1}^N \phi(\|x - x^j\|)$



(ii) GRADIENT OPTIMIZATION

$$E = \frac{1}{4} \sum_{n=1}^N \sum_{k=1}^C \left(t_k^n - \sum_{j=0}^M w_{kj} \phi(\|x^n - \mu_j\|^2) \right)^2$$

$$\nabla_{\mu_j} E = \sum_{n=1}^N \left[\sum_{k=1}^C \left(\sum_{j=0}^M w_{kj} \phi(\|x^n - \mu_j\|^2) \right) t_k^n - t_k^n \right] \phi'(\|x^n - \mu_j\|^2) (\mu_j - x^n)$$

CALL THIS $\gamma_j(x^n)$

⇒ EM-LIKE GRADIENT DESCENT ALGORITHM

- ①
$$\mu_j = \frac{\sum_{n=1}^N \gamma_j(x^n) x^n}{\sum_{n=1}^N \gamma_j(x^n)}$$

- ② RECOMPUTE w_{kj}

- ③ RECOMPUTE $\gamma_j(x^n)$

- ④ GO TO ① UNTIL CONVERGENCE

(iii) K-MEANS

$$E = \frac{1}{2} \sum_{n=1}^N \min_j \|x^n - \mu_j\|^2$$

$$\nabla_{\mu_j} E = \sum_{n=1}^N \mathbb{1}_{\{j = j_n^*\}} (\mu_j - x^n)$$

- ① COMPUTE $j_n^* = \underset{j}{\operatorname{argmin}} \|x^n - \mu_j\|^2$

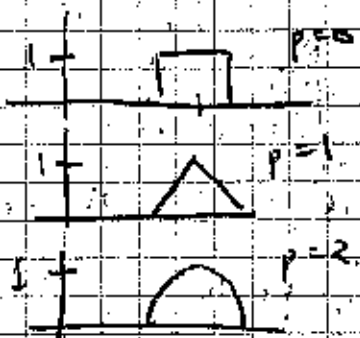
- ② FIND
$$\mu_j = \frac{\sum_{n=1}^N \mathbb{1}_{\{j = j_n^*\}} x^n}{\sum_{n=1}^N \mathbb{1}_{\{j = j_n^*\}}}$$

- ③ REPEAT UNTIL CONVERGENCE

③ How can we set σ , τ ?

A: You can't! (NOT BASED ON TRAINING DATA ALONE)

PROOF: LET $\sigma(x^n, y^n) = \frac{\tau - \|\bar{x}_n - \bar{y}_n\|_p^p}{\tau}$



ETC

LET $M = N$, $\bar{x}_j = \bar{x}^j$, $w_{kj} = t_k^j$

THEN, FOR ANY SMALL ENOUGH σ , FOR ANY p ,

$$\epsilon = \sum_{k=1}^n \sum_{j=1}^n \left(t_{kj} - \sum_{j=1}^n [t_{kj} \cdot \sigma] \right)^2 = 0!$$

KNOWLEDGE BEYOND THE DATASET: REGULARIZATION

$$\epsilon = \frac{1}{2} \sum_{k=1}^n (t_k^n - y(x^n))^2 + \frac{\nu}{2} \int |p_y(x)|^2 dx$$

↑ PENALTY FOR NON-SMOOTHNESS OF $y(x)$
 ↑ ESTIMATE OF NON-SMOOTHNESS

$$\frac{\partial \epsilon}{\partial y(x)}$$

$$= \begin{cases} \nu \hat{p} p_y(x) & x \notin x^n \\ \nu \hat{p} p_y(x) + (y(x) - t_k^n) s(x - x^n) & x \in x^n \end{cases}$$

How much does changing $y(x)$ at this particular x affect the average non-smoothness?

SOLUTION: SETTING $\frac{\partial E}{\partial y(x)} = 0$ AT EACH x

GIVES $y(x) = \sum_{n=1}^N w_n G(x - x^n)$

FOR $G(x)$ S.T. $\hat{P}P G(x - x^n) = \delta(x - x^n)$

EXAMPLE

$$E = \sum_{n=1}^N (y(x^n) - t^n)^2 + \sum_{l=0}^{\infty} \frac{1}{l! 2^l} \int |\nabla^l y(x)|^2 dx$$

NORMALIZE
LIKE THE
TAYLOR SERIES
FOR e^x

2TH POWER OF
CURVATURE

$\frac{\partial E}{\partial y(x)} = 0$ AT EACH x

$$\Rightarrow y(x) = \sum_{n=1}^N w_n e^{-\frac{1}{2} \frac{\|x - x^n\|^2}{\sigma^2}}$$

ONE MORE PERSPECTIVE: REPRESENTER THEOREM

MERCER'S THEOREM

SUPPOSE $\Phi(x - x^n)$ IS A POSITIVE DEFINITE KERNEL,
MEANING THAT FOR ANY $f(x)$,

$$\iint \Phi(x - x^n) f(x) f(x^n) dx dx^n > 0$$

THEN

$$\Phi(x - x^n) = \sum_{j=0}^{\infty} \lambda_j \psi_j(x) \psi_j(x^n)$$

$[\lambda_j, \psi_j]$ = EIGENVALUES, EIGENFUNCTIONS OF Φ

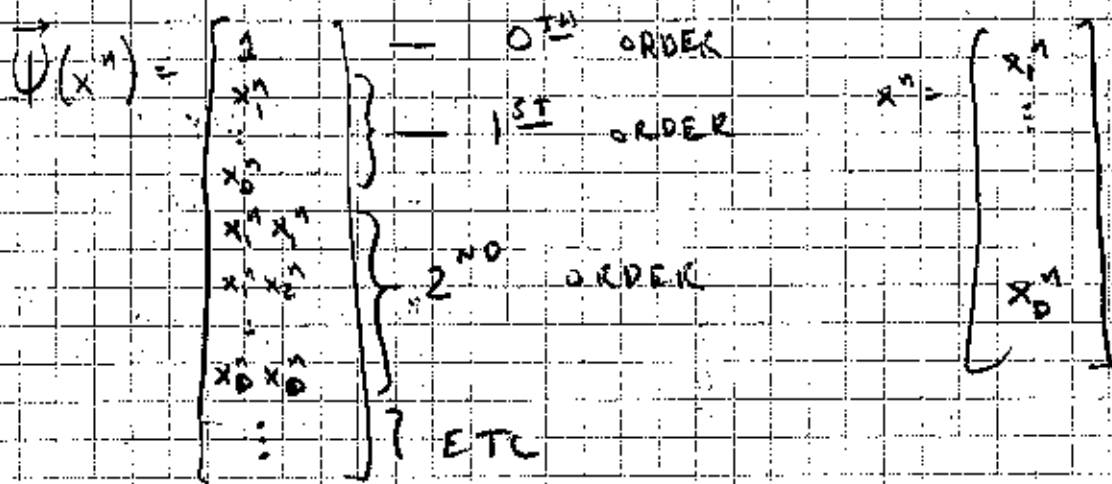
AND THEREFORE L_2 NORMS EXIST:

$$\begin{aligned} \|\vec{\Psi}(x^n) - \vec{\Psi}(x^m)\|^2 &= \|\vec{\Psi}(x^n)\|^2 + \|\vec{\Psi}(x^m)\|^2 - 2\vec{\Psi}(x^n)^T \vec{\Psi}(x^m) \\ &= \sum_{j=0}^{\infty} \lambda_j (\psi_j^2(x^n) + \psi_j^2(x^m) - 2\psi_j(x^n)\psi_j(x^m)) \end{aligned}$$

EXAMPLE

$$e^{-\|x - x^n\|^2} = 1 - (x - x^n)^T (x - x^n) + \frac{1}{2} \left[(x - x^n)^T (x - x^n) \right]^2 - \dots$$

$$= \sum_{j=0}^{\infty} \psi_j(x) \psi_j(x^n) \lambda_j$$



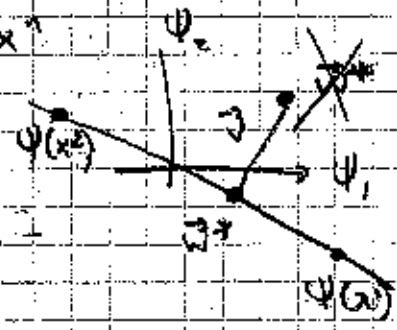
REPRESENTATION THEOREM LET $y(x) = \sum_{j=0}^{\infty} w_j \psi_j(x) = \vec{w}^T \vec{\psi}(x)$

SUPPOSE $\vec{w}^* = \text{argmin} \sum_{n=1}^N \left(t_n - \vec{w}^T \vec{\psi}(x^n) \right)^2 + \frac{\nu}{2} \|\vec{w}\|^2$

BUT IT'S ALWAYS TRUE THAT

$$\vec{w}^* = \underbrace{\sum_{n=1}^N \alpha_n \psi(x^n)}_{\text{IN THE SPAN OF } x^n} + \underbrace{\vec{v}}_{\text{ORTHOGONAL TO THE SPAN OF } x^n} \quad \text{FOR}$$

$$\|\vec{w}^*\|^2 = \left\| \sum_{n=1}^N \alpha_n \psi(x^n) \right\|^2 + \|\vec{v}\|^2$$



$$\sum_{n=1}^N \left(t_n - \vec{w}^T \vec{\psi} - \sum_{m=1}^N \alpha_m \vec{w}^T \psi(x^m) \right)^2$$

↑ N UNKNOWN
CHOOSE TO MAKE THIS ZERO!

THUS $\vec{v}^{opt} = 0$, $\vec{w}^* = \sum_{n=1}^N \alpha_n \psi(x^n)$

THEN $y(x) = \vec{w}^T \psi(x) = \sum_{n=1}^N \alpha_n \psi(x)^T \psi(x^n)$

$y(x) = \sum_{n=1}^N \alpha_n \phi(x - x^n)$ MINIMIZES ϵ