

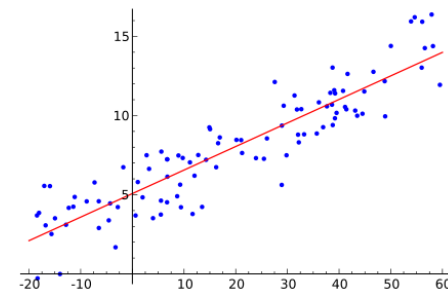
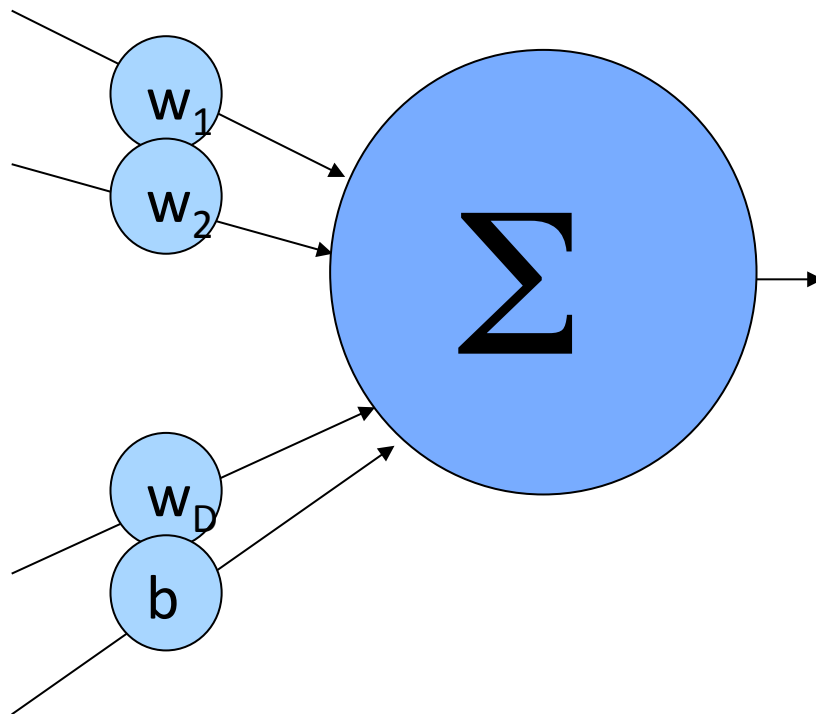
UCB - CS189
Introduction to Machine Learning
Fall 2015

Lecture 8: Kernel machines

Isabelle Guyon
ChaLearn

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

Last time



Linear regression

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

Today

Kernel machines

PARAMETRIC (Perceptrons)

$$f(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x})$$

$$\mathbf{w} = \sum_k \alpha_k \Phi(\mathbf{x}^k)$$

(Large margin) Perceptron

$$\Delta \mathbf{w} \sim y_k \Phi(\mathbf{x}^k) \quad \text{if } y_k f(\mathbf{x}^k) < 1$$

$$\sim \mathbf{1}(1 - z_k) y_k \Phi(\mathbf{x}^k) \quad z_k = y_k f(\mathbf{x}^k)$$

(Rosenblatt 1958)

Logistic regression

$$\Delta \mathbf{w} \sim S(-z_k) y_k \Phi(\mathbf{x}^k)$$

(Cox 1958)

LMS regression or classification

$$\Delta \mathbf{w} \sim (y_k - f(\mathbf{x}^k)) \Phi(\mathbf{x}^k) \sim (1 - z_k) y_k \Phi(\mathbf{x}^k)$$

(Widrow-Hoff, 1960)

NON PARAMETRIC (Kernel machines)

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

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

Potential Function algorithm

$$\Delta \alpha_k \sim y_k \quad \text{if } y_k f(\mathbf{x}^k) < 1$$

$$\sim \mathbf{1}(1 - z_k) y_k$$

(Aizerman et al 1964)

Dual logistic regression

$$\Delta \alpha_k \sim S(-z_k) y_k$$

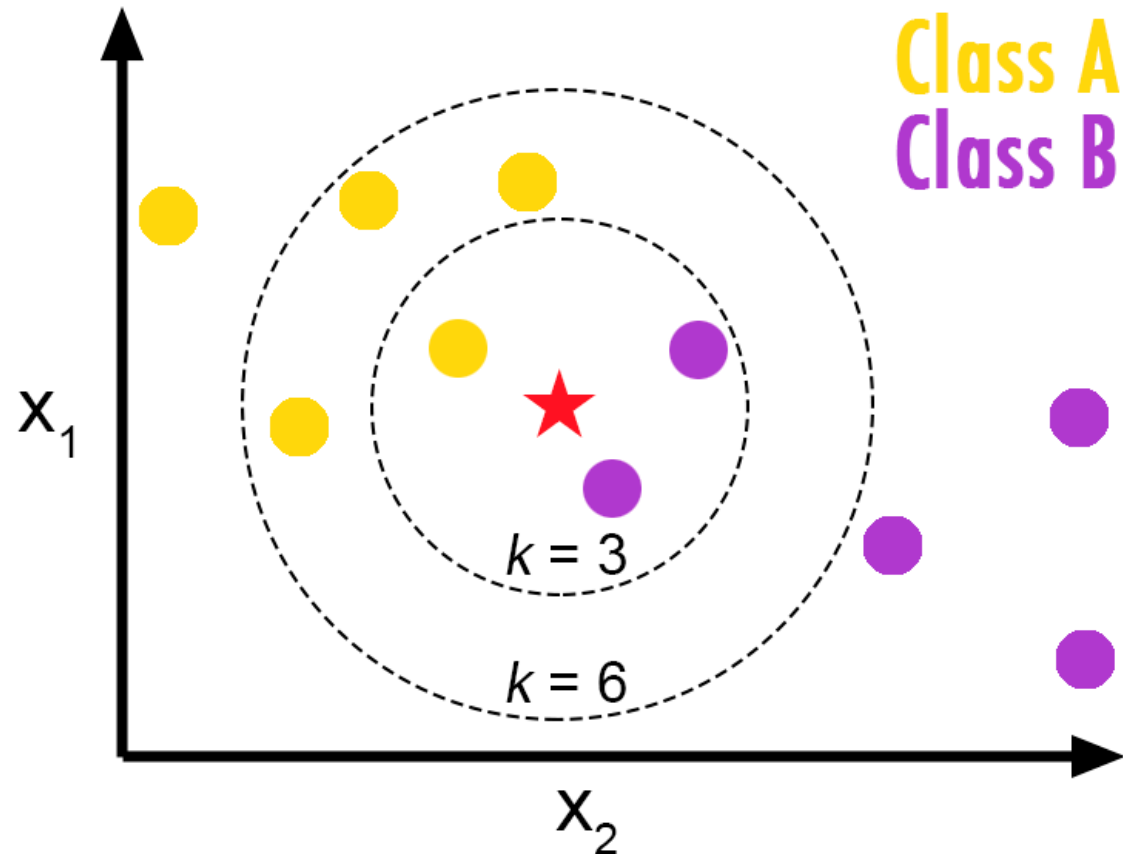
Dual LMS

$$\Delta \alpha_k \sim (y_k - f(\mathbf{x}^k)) \sim (1 - z_k) y_k$$

Linear regression

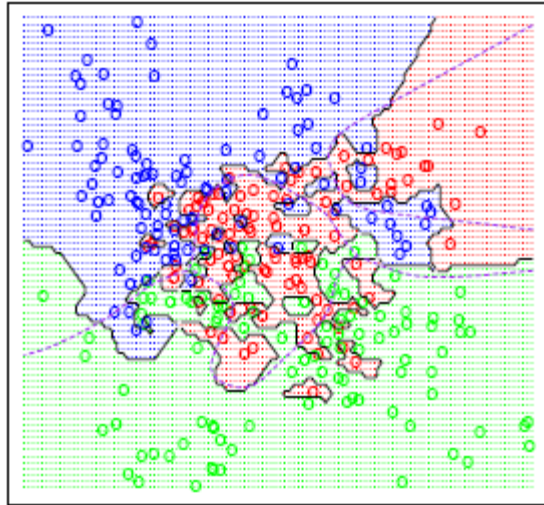
- **Problem setting:**
 - N training input/output pairs $\{(\mathbf{x}^k, y^k)\}$; y is continuous.
 - Linear model $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$ (w_0 corresponds to $x_0 \equiv 1$, $\dim(\mathbf{x})=d$)
 - Solve a system of N equations $y^k = w_0 + w_1 x_0 + \dots$ of $(d+1)$ unknown.
- **Matrix notation:**
 - $X \mathbf{w}^T = \mathbf{y}$
 - $X = [x^k_i]$ of $\dim(N, d)$; $\mathbf{y} = [y^k]$ of $\dim(N, 1)$; $\mathbf{w} = [w_i]$ of $\dim(1, d)$.
- **Solution:**
 - Case 1: $N > d$, over-determined, no exact solution. $\mathbf{w}^T = X^+ \mathbf{y}$ is the best solution in the least-square sense. $\text{RSS} = \sum_k (f(\mathbf{x}) - y^k)^2$.
 - Case 2: $N < d$, under-determined, $\mathbf{w}^T = X^+ \mathbf{y}$ is the solution of $\min \|\mathbf{w}\|$.
- **Pseudo-inverse:**
 - $X^+ = \lim_{\lambda \rightarrow 0^+} (X^T X + \lambda I)^{-1} X^T = \lim_{\lambda \rightarrow 0^+} X^T (X X^T + \lambda I)^{-1}$
 - Inverse the smallest matrix and adjust $\lambda > 0$ by cross-validation (see SVD trick).
- **Kernel trick:**
 - in case 2, $N \ll d$, replace $X X^T$ by a (N, N) kernel matrix $K = k(\mathbf{x}^k, \mathbf{x}^h)$
 - $\boldsymbol{\alpha} = (K + \lambda I)^{-1} \mathbf{y}$, gives you a non linear regression function $f(\mathbf{x}) = \sum_k \alpha_k k(\mathbf{x}, \mathbf{x}_k)$.

Ancestor of kernel methods: KNN

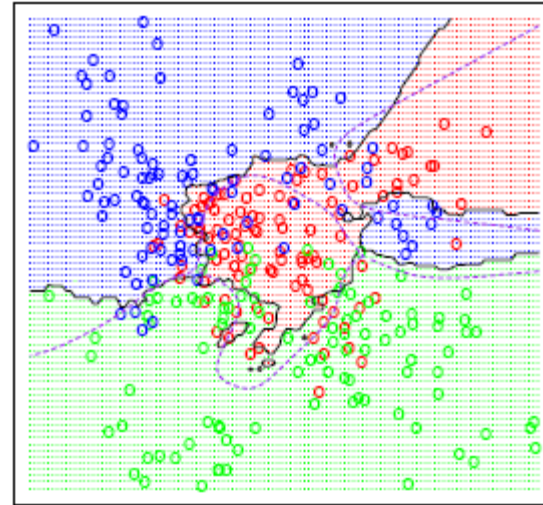


Assign to ★ the class label of the majority of the k closest examples.

Optimum k in k-nearest neighbors



k=1

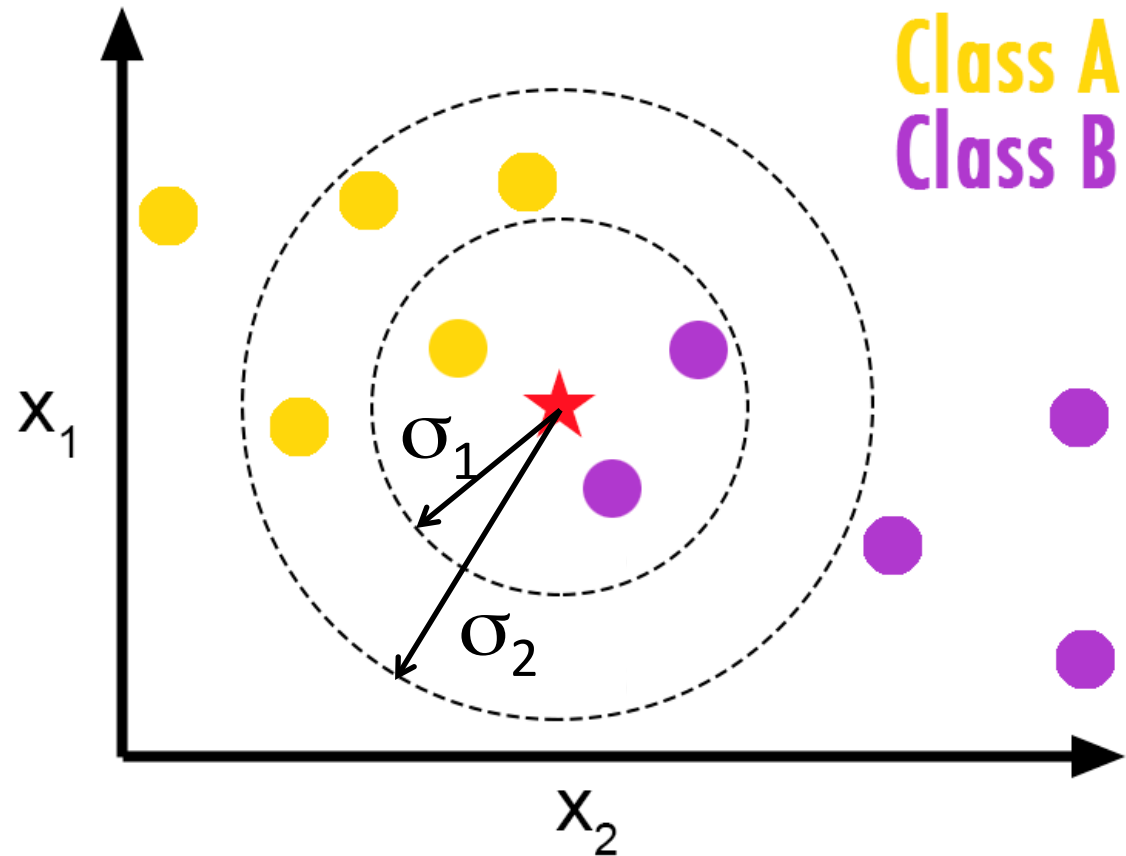


k=15

Fit vs. robustness tradeoff in kNN:

- kNN does not have any parameters (**non parametric method**).
- The “complexity” (of the decision function) *decreases* with k.
- k is a *hyper-parameter* adjustable by *cross-validation*.
- kNN can also be used for regression.

Parzen windows



Assign to ★ the class label of the majority of the examples enclosed in a sphere of radius σ_r .

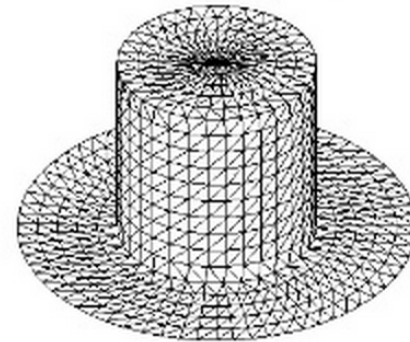
Parzen windows is a kernel method

- Assign to \mathbf{x} the class label of the majority of the examples enclosed in a sphere of radius σ .

- $f(\mathbf{x}) = \sum_{k=1:N} y_k k(\mathbf{x}, \mathbf{x}_k)$

- Top hat kernel:

$$k(\mathbf{x}, \mathbf{x}_k) = \mathbf{1}(\|\mathbf{x} - \mathbf{x}_k\|^2 < \sigma^2)$$



Parzen windows is a kernel method

- Assign to \mathbf{x} the class label of the majority of the examples enclosed in a sphere of radius σ .

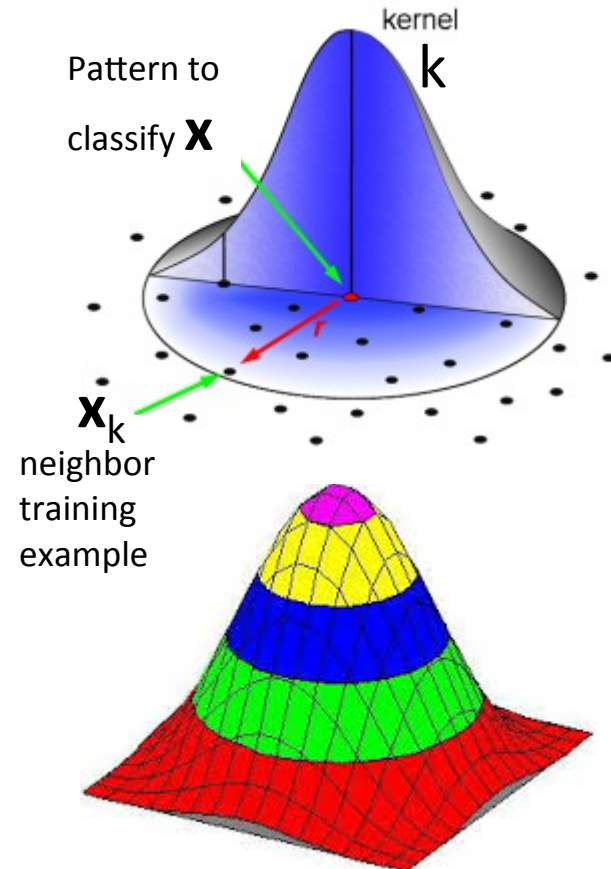
- $$f(\mathbf{x}) = \sum_{k=1:N} y_k k(\mathbf{x}, \mathbf{x}_k)$$

- Top hat kernel:

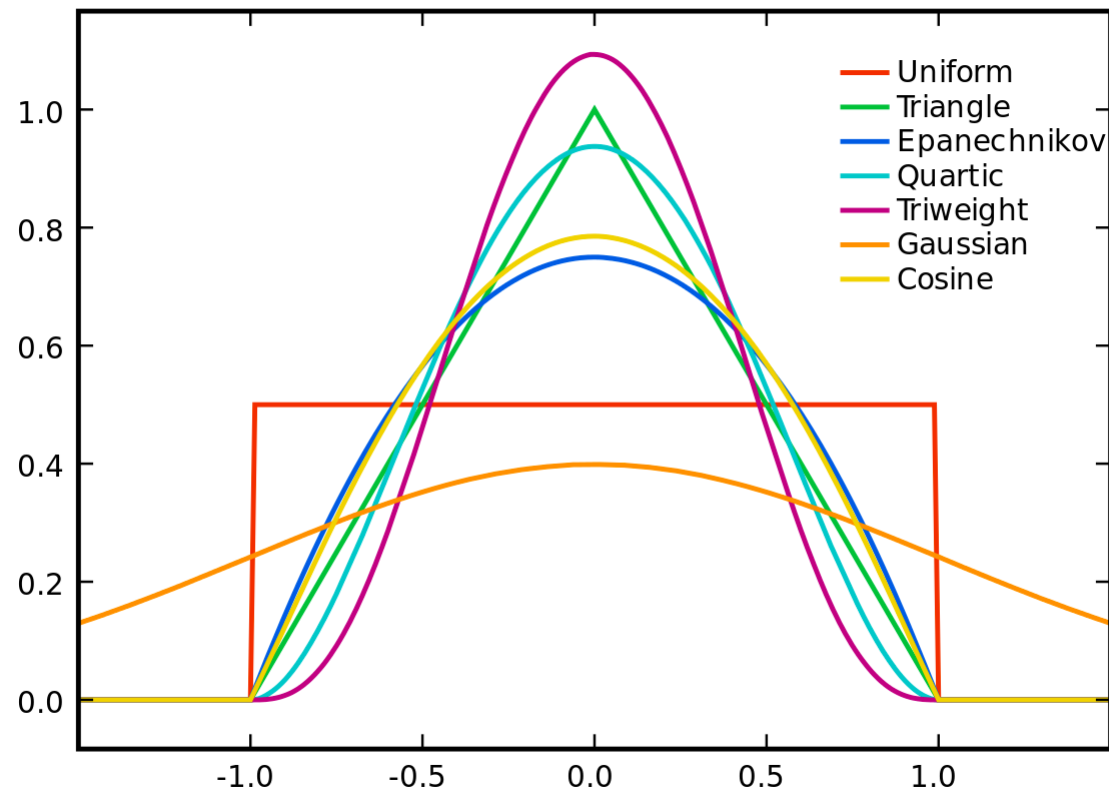
$$k(\mathbf{x}, \mathbf{x}_k) = \mathbf{1}(\|\mathbf{x} - \mathbf{x}_k\|^2 < \sigma^2)$$

- Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}_k) = \exp -(\|\mathbf{x} - \mathbf{x}_k\|^2 / 2\sigma^2)$$



Some radial kernels



Source: [https://en.wikipedia.org/wiki/Kernel_\(statistics\)](https://en.wikipedia.org/wiki/Kernel_(statistics))

See also: <http://crsouza.com/2010/03/kernel-functions-for-machine-learning-applications>

<http://www.cs.berkeley.edu/~bartlett/courses/281b-sp08/8.pdf>

Non radial kernels

$$f(\mathbf{x}) = \sum_{k=1:N} y_k k(\mathbf{x}, \mathbf{x}_k)$$

- Linear kernel:

$$k(\mathbf{x}, \mathbf{x}_k) = \mathbf{x} \cdot \mathbf{x}_k$$

The Parzen windows $f(\mathbf{x})$ for the linear kernel is just Hebb's rule!

- Polynomial kernel:

$$k(\mathbf{x}, \mathbf{x}_k) = (1 + \mathbf{x} \cdot \mathbf{x}_k)^q$$

What makes $k(\mathbf{x}, \mathbf{x}')$ a “good” kernel?

- Symmetric $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x})$
- Kernel matrix $K = [k(\mathbf{x}_k, \mathbf{x}_h)]_{k=1:N, h=1:N}$ invertible, possibly after “regularization” $(K + \lambda I)$, $\lambda > 0$. Satisfied if all eigenvalues ≥ 0 (PSD matrix).

- True in particular if:
 - $k(\mathbf{x}, \mathbf{x}')$ is a PSD kernel i.e. satisfies Mercer’s condition.

$$\sum_{i=1}^n \sum_{j=1}^n K(x_i, x_j) c_i c_j \geq 0$$

for all finite sequences of points x_1, \dots, x_n of $[a, b]$ and all choices of real numbers c_1, \dots, c_n

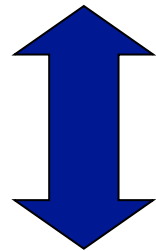
There exists an eigen decomposition $k(\mathbf{x}, \mathbf{x}') = \sum_{i=1:\infty} \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{x}')$

- K is a Gram matrix (outer product $K = XX^T$ or $K = \Phi\Phi^T$)
- K is a covariance matrix.
- The sigmoid kernel $\tanh(\mathbf{x} \cdot \mathbf{x}')$ is NOT PSD.

Remember the “kernel trick”

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

NON
PARAMETRIC



Dual forms

- $f(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x})$
- $\mathbf{w} = \sum_k \alpha_k \Phi(\mathbf{x}^k)$

PARAMETRIC

Reproducing kernels

- A Hilbert space H is a space of functions $\{f(\mathbf{x})\}$.
- A “reproducing kernel” Hilbert space (RKHS) is endowed with a dot product:

$$\langle f, g \rangle_H = \int f(\mathbf{x}) g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

- H has a (unique) positive definite “reproducing kernel”:

$$k(\mathbf{s}, \mathbf{t}) = \langle k(\cdot, \mathbf{s}), k(\cdot, \mathbf{t}) \rangle_H$$

- Reproducing kernel property:

$$f(\mathbf{x}) = \langle k(\cdot, \mathbf{x}), f \rangle_H$$

Representer theorem

Kimeldorf, George S.; Wahba, Grace (1970)

Schölkopf, Bernhard; Herbrich, Ralf; Smola, Alex J. (2001)

- $f(\mathbf{x})$ is a function in RKHS H , w. $\langle f, g \rangle_H = \int f(\mathbf{x}) g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$, and kernel k : $f(\mathbf{x}) = \langle k(\cdot, \mathbf{x}), f \rangle_H$.
- Given a risk functional $R[f]$ and a regularizer $\Omega[f]$ and training examples $\{\mathbf{x}^k, \mathbf{y}^k\}$ on which we evaluate $R[f] + \lambda\Omega[f]$, $\lambda > 0$.
- $f^*(\mathbf{x}) = \operatorname{argmin}_f R_{\text{train}}[f] + \lambda\Omega[f]$ admits a representation

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

- This means that we only have to solve for α_k and that the solution lies in the span of $k(\cdot, \mathbf{x}^k)$ that can be thought of as special kind of features. The kernel determines the type of function used and the type of regularization (smoothness).

See for details: https://en.wikipedia.org/wiki/Representer_theorem

Mercer kernels and kernel trick

$$k(\mathbf{x}, \mathbf{x}') = \sum_{i=1:\infty} \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{x}')$$

$\langle f, g \rangle = \int f(\mathbf{s}) g(\mathbf{s}) p(\mathbf{s}) d\mathbf{s}$; $\psi_i(\mathbf{x})$ are eigen functions, $\langle \psi_i, \psi_j \rangle = \delta_{ij}$
 $\lambda_i \psi_i(\mathbf{x}) = \langle k(\cdot, \mathbf{x}), \psi_i \rangle$ with eigen values $\lambda_i > 0$

$$\langle f, g \rangle_H = \sum_{i=1:\infty} (1/\lambda_i) \langle f, \psi_i \rangle \langle g, \psi_i \rangle$$

Use $\phi_i(\mathbf{x}) = \sqrt{\lambda_i} \psi_i(\mathbf{x})$

$$\begin{aligned} f(\mathbf{x}) &= \sum_k \alpha_k k(\mathbf{x}^k, \mathbf{x}) = \sum_k \alpha_k \sum_i \phi_i(\mathbf{x}^k) \phi_i(\mathbf{x}) \\ &= \sum_i \left(\sum_k \alpha_k \phi_i(\mathbf{x}^k) \right) \phi_i(\mathbf{x}) \\ &= \sum_i w_i \phi_i(\mathbf{x}) \end{aligned}$$

See for details: https://en.wikipedia.org/wiki/Representer_theorem

Kernel algebra

Cristianini and Shaw-Taylor, 2001

Proposition 3.22 (Closure properties) *Let κ_1 and κ_2 be kernels over $X \times X$, $X \subseteq \mathbb{R}^n$, $a \in \mathbb{R}^+$, $f(\cdot)$ a real-valued function on X , $\phi: X \rightarrow \mathbb{R}^N$ with κ_3 a kernel over $\mathbb{R}^N \times \mathbb{R}^N$, and \mathbf{B} a symmetric positive semi-definite $n \times n$ matrix. Then the following functions are kernels:*

- (i) $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_1(\mathbf{x}, \mathbf{z}) + \kappa_2(\mathbf{x}, \mathbf{z})$,
- (ii) $\kappa(\mathbf{x}, \mathbf{z}) = a\kappa_1(\mathbf{x}, \mathbf{z})$,
- (iii) $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_1(\mathbf{x}, \mathbf{z})\kappa_2(\mathbf{x}, \mathbf{z})$,
- (iv) $\kappa(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})f(\mathbf{z})$,
- (v) $\kappa(\mathbf{x}, \mathbf{z}) = \kappa_3(\phi(\mathbf{x}), \phi(\mathbf{z}))$,
- (vi) $\kappa(\mathbf{x}, \mathbf{z}) = \mathbf{x}'\mathbf{B}\mathbf{z}$.

Proposition 3.24 *Let $\kappa_1(\mathbf{x}, \mathbf{z})$ be a kernel over $X \times X$, where $\mathbf{x}, \mathbf{z} \in X$, and $p(x)$ is a polynomial with positive coefficients. Then the following functions are also kernels:*

- (i) $\kappa(\mathbf{x}, \mathbf{z}) = p(\kappa_1(\mathbf{x}, \mathbf{z}))$,
- (ii) $\kappa(\mathbf{x}, \mathbf{z}) = \exp(\kappa_1(\mathbf{x}, \mathbf{z}))$,
- (iii) $\kappa(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|^2 / (2\sigma^2))$.

Favorite kernel

$$k(\mathbf{x}, \mathbf{x}') = (a + \mathbf{x} \cdot \mathbf{x}')^b \exp(-d \|\mathbf{x} - \mathbf{x}_k\|^c)$$

Of 4 hyper-parameters a , b , c , d .

Simplify: take $a = 1$ and $c = 2$.

Then: $b = 0 \rightarrow$ Gaussian kernel

$d = 0 \rightarrow$ Polynomial kernel

Adjust the hyper-parameters by cross-validation.

Kernel machines

PARAMETRIC (Perceptrons)

$$f(\mathbf{x}) = \mathbf{w} \bullet \Phi(\mathbf{x})$$

$$\mathbf{w} = \sum_k \alpha_k \Phi(\mathbf{x}^k)$$

(Large margin) Perceptron

$$\begin{aligned} \Delta \mathbf{w} &\sim y_k \Phi(\mathbf{x}^k) \quad \text{if } y_k f(\mathbf{x}^k) < 1 \\ &\sim \mathbf{1}(1 - z_k) y_k \Phi(\mathbf{x}^k) \quad z_k = y_k f(\mathbf{x}^k) \end{aligned}$$

(Rosenblatt 1958)

Logistic regression

$$\Delta \mathbf{w} \sim S(-z_k) y_k \Phi(\mathbf{x}^k)$$

(Cox 1958)

LMS regression or classification

$$\Delta \mathbf{w} \sim (y_k - f(\mathbf{x}^k)) \Phi(\mathbf{x}^k) \sim (1 - z_k) y_k \Phi(\mathbf{x}^k)$$

(Widrow-Hoff, 1960)

NON PARAMETRIC (Kernel machines)

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

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

Potential Function algorithm

$$\begin{aligned} \Delta \alpha_k &\sim y_k \quad \text{if } y_k f(\mathbf{x}^k) < 1 \\ &\sim \mathbf{1}(1 - z_k) y_k \end{aligned}$$

(Aizerman et al 1964)

Dual logistic regression

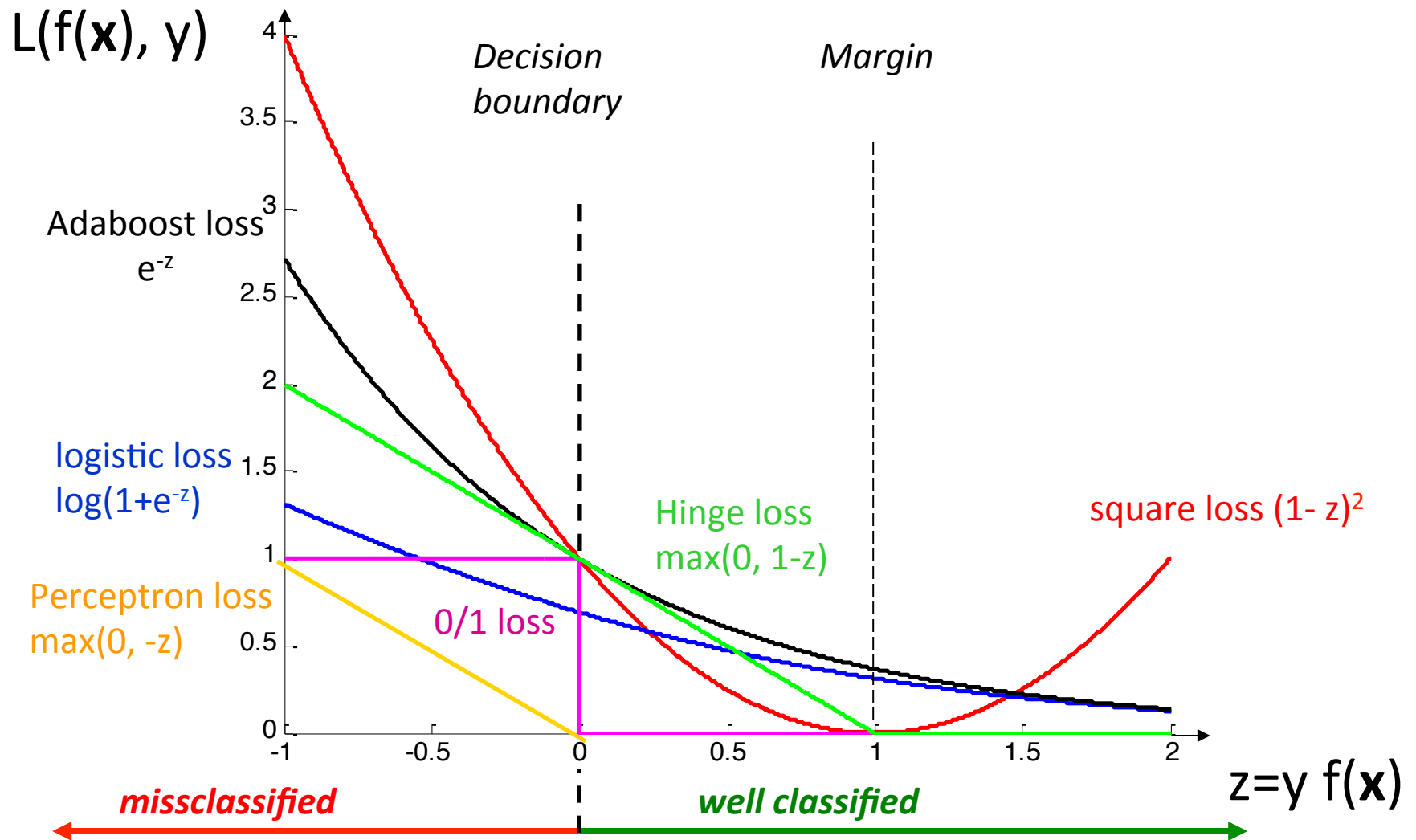
$$\Delta \alpha_k \sim S(-z_k) y_k$$

Dual LMS

$$\Delta \alpha_k \sim (y_k - f(\mathbf{x}^k)) \sim (1 - z_k) y_k$$

Loss Functions

The risk is the average of the loss.



Regularized risk minimization

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

- Minimize:

$$R_{\text{reg}}[\mathbf{w}] = \sum_k L(f(\mathbf{x}^k), y_k) + \lambda \Omega(\mathbf{w})$$

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|^2$$

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|$$

Stochastic gradient

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

1) Compute $\partial L / \partial \mathbf{w}_i = \partial L / \partial f \cdot \partial f / \partial \mathbf{w}_i = \partial L / \partial f \Phi_i(\mathbf{x})$

- Classification: $\partial L / \partial \mathbf{w}_i = \partial L / \partial z \cdot \partial z / \partial \mathbf{w}_i$ with $z = y f(\mathbf{x})$. $\partial z / \partial \mathbf{w}_i = y \Phi_i(\mathbf{x}) \rightarrow \partial L / \partial \mathbf{w}_i = \partial L / \partial z \cdot y \Phi_i(\mathbf{x})$

2) Compute $\partial \Omega / \partial \mathbf{w}_i$

- $\Omega(\mathbf{w}) = \|\mathbf{w}\|^2 \rightarrow \partial \Omega / \partial \mathbf{w}_i = 2\mathbf{w}_i$

3) Compute the negative gradient of $L(f(\mathbf{x}), y) + \lambda \Omega(\mathbf{w})$

$$\Delta \mathbf{w}_i = -\eta \partial L / \partial f \Phi_i(\mathbf{x}) - \gamma \mathbf{w}_i$$

$$\Delta \mathbf{w} = -\eta \partial L / \partial f \Phi(\mathbf{x}) - \gamma \mathbf{w}$$

Kernel version

Φ -space version: $\Delta \mathbf{w} = -\eta \partial L / \partial f \Phi(\mathbf{x})$

- Learning (\mathbf{x}^k, y^k) : $\Delta \mathbf{w} = -\eta \partial L / \partial f|_{(\mathbf{x}^k, y^k)} \Phi(\mathbf{x}^k)$
- $\mathbf{w} = \sum_k \alpha_k \Phi(\mathbf{x}^k)$, $\Delta \mathbf{w} = \Delta \alpha_k \Phi(\mathbf{x}^k)$

$$\Delta \alpha_k = -\eta \partial L / \partial f|_{(\mathbf{x}^k, y^k)}$$

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

Kernel version **with shrinkage**

Φ -space version: $\Delta \mathbf{w} = -\eta \partial L / \partial f \Phi(\mathbf{x}) - \gamma \mathbf{w}$

- Learning (\mathbf{x}^k, y^k) :

$$\Delta \mathbf{w} = -\eta \partial L / \partial f|_{(\mathbf{x}^k, y^k)} \Phi(\mathbf{x}^k) - \gamma \mathbf{w}$$

- $\mathbf{w} = \sum_k \alpha_k \Phi(\mathbf{x}^k)$, $\Delta \mathbf{w} = \Delta \alpha_k \Phi(\mathbf{x}^k)$

$$\begin{aligned} \Delta \alpha_k &= -\eta \partial L / \partial f|_{(\mathbf{x}^k, y^k)} - \gamma \alpha_k && \text{for example } k \\ \Delta \alpha_h &= -\gamma \alpha_h && \text{for other examples} \end{aligned}$$

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

Example: Least Mean Square or LMS

Widrow-Hoff, 1960

- $L_{\text{square_loss}} = (f(\mathbf{x}) - y)^2$
- $L/\partial f = 2 (f(\mathbf{x}) - y)$
- $\partial f/\partial w_i = \Phi_i(\mathbf{x})$ for $f(\mathbf{x}) = \sum_j w_j \Phi_j(\mathbf{x})$
- $\partial L/\partial w_i = 2 (f(\mathbf{x}) - y) \Phi_i(\mathbf{x})$
- $\Delta w_i = -\eta (f(\mathbf{x}) - y) \Phi_i(\mathbf{x}) - \gamma \mathbf{w}$

Loss variation

$$\Delta \mathbf{w} = \eta (y - f(\mathbf{x})) \Phi(\mathbf{x}) - \gamma \mathbf{w}$$

LMS learning rule

$$\Delta \alpha_k = \eta (y^k - f(\mathbf{x}^k)) - \gamma \alpha_k$$

for example k

$$\Delta \alpha_h = -\gamma \alpha_h$$

for other examples

Comparison Hebb's rule

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

- Hebb's rule:

$$\Delta w_i = \eta y \Phi_i(\mathbf{x})$$

$$\Delta \alpha_k = \eta y^k$$

- LMS rule:

$$\Delta w_i = \eta (y - f(\mathbf{x})) \Phi_i(\mathbf{x})$$

$$\Delta \alpha_k = \eta (y^k - f(\mathbf{x}^k))$$

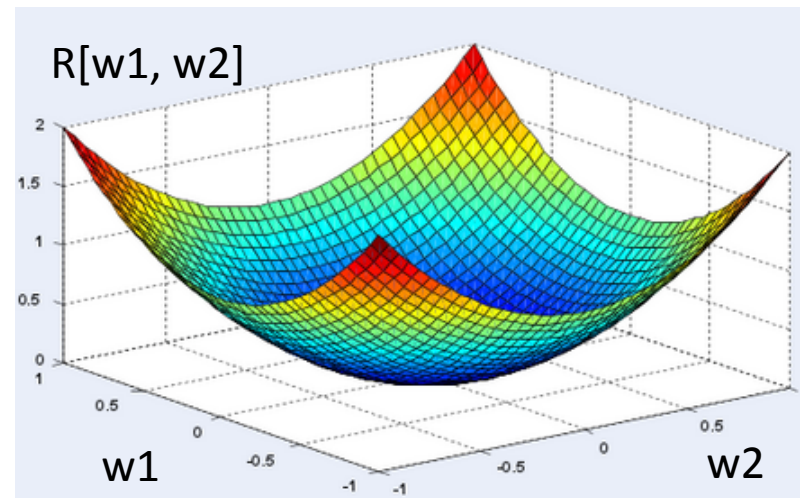
1

Exploiting convexity?

- Stochastic gradient does not exploit convexity, is this the best approach?
 - Yes for BIG data.
 - No if either N or d is small.
- For regression, regularized inverse:

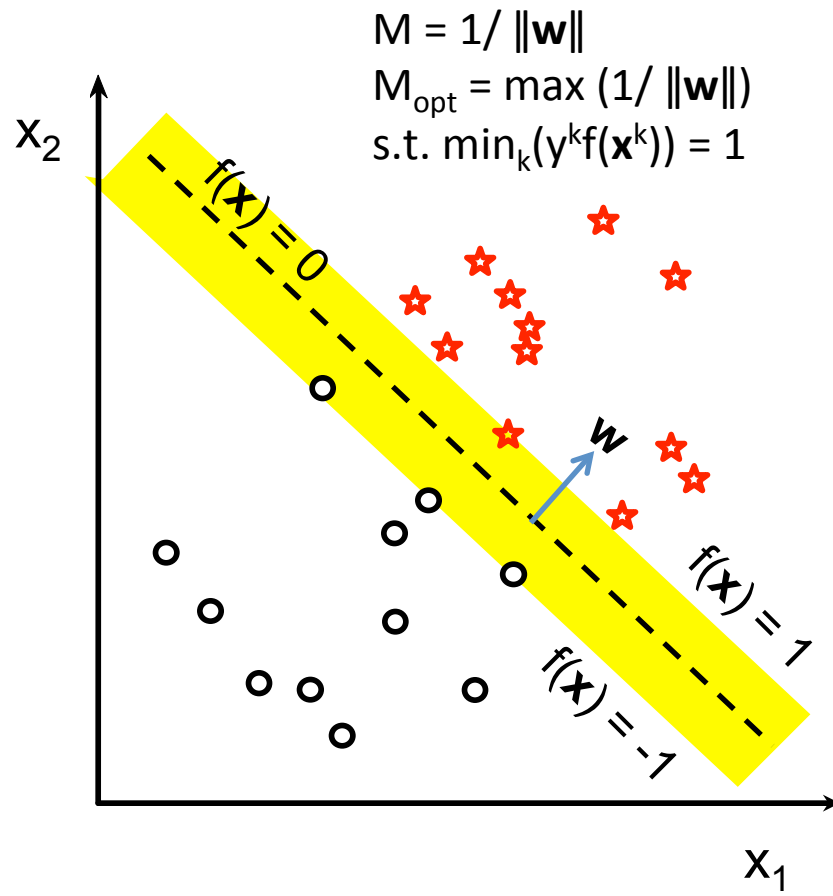
$$\mathbf{w}^T = (X^T X + \lambda I)^{-1} X^T \mathbf{y} \text{ or } X^T (X X^T + \lambda I)^{-1} \mathbf{y}$$

$$\boldsymbol{\alpha} = (K + \lambda I)^{-1} \mathbf{y}$$

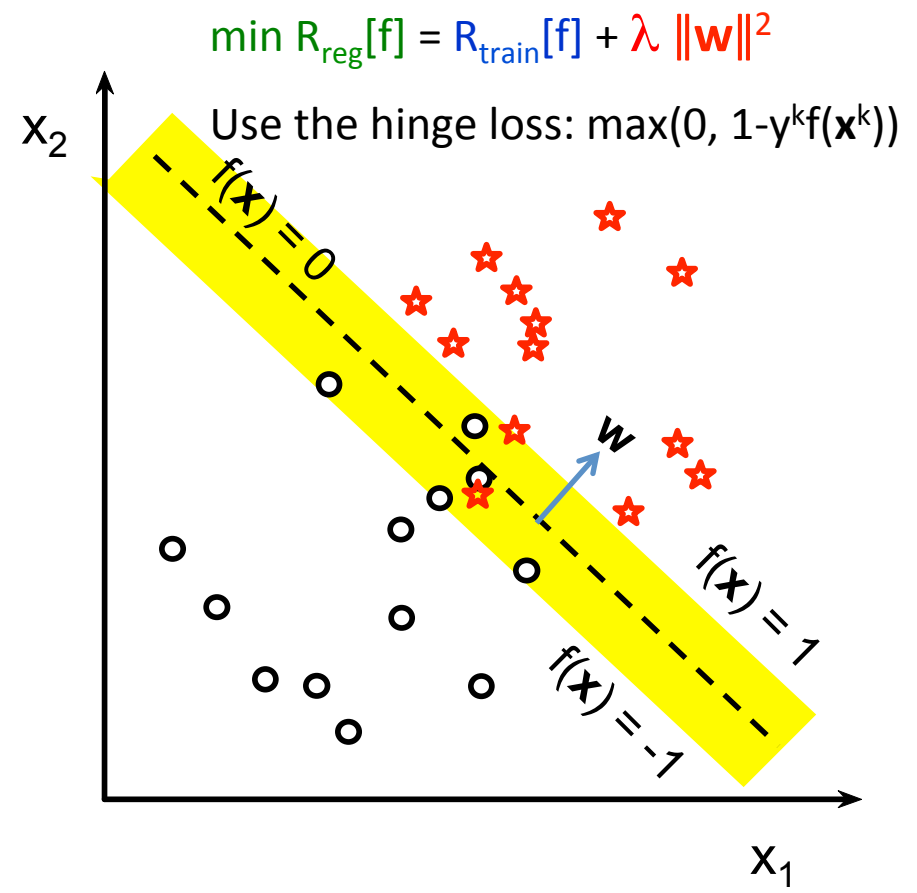


Classification with optimum margin

Hard margin

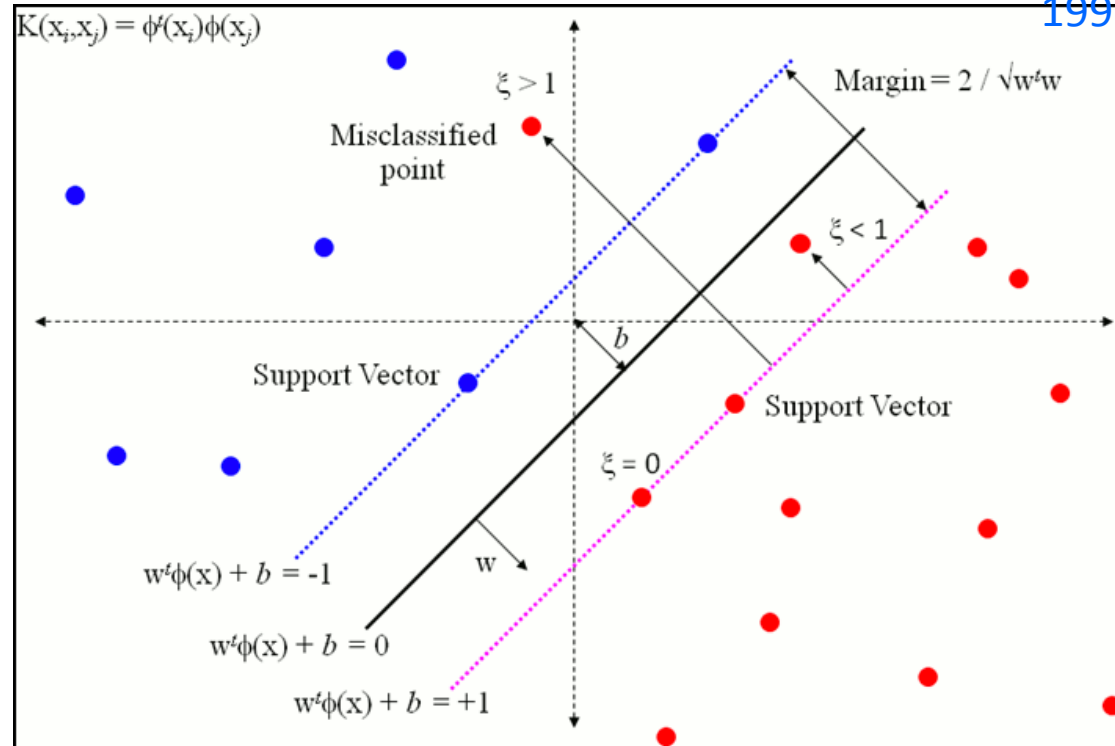


Soft margin



Soft margin

Cortes-Vapnik,
1995



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

Use the hinge loss: $\max(0, 1 - y^k f(x^k))$

$$\min C \sum_{k=1:N} \xi_k + (1/2) \|w\|^2$$

$$1 - y^k f(x^k) \leq \xi_k \quad \xi_k \geq 0$$

Standard quadratic programs

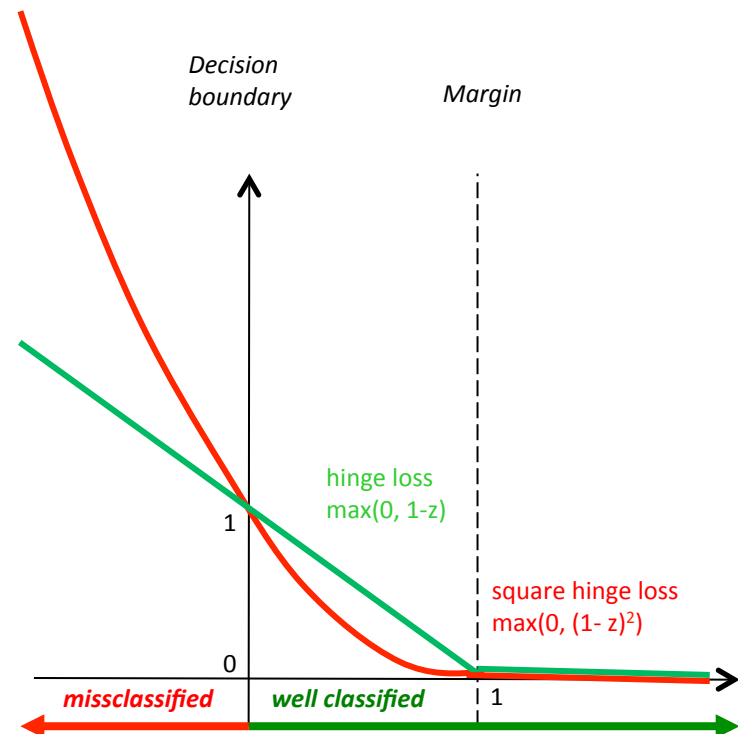
$$f(\mathbf{x}) = \mathbf{w} \cdot \Phi(\mathbf{x}) + b$$

1) Soft margin Perceptron-style
(with hinge loss)

$$\begin{cases} \min_{\mathbf{w}, b, \xi} C \sum_{k=1:N} \xi_k + (1/2) \|\mathbf{w}\|^2 \\ y^k f(\mathbf{x}^k) - 1 + \xi_k \geq 0 \\ \xi_k \geq 0 \end{cases}$$

2) Soft margin RSS-style
(with square hinge loss)

$$\begin{cases} \min_{\mathbf{w}, b, \xi} C \sum_{k=1:N} (\xi_k)^2 + (1/2) \|\mathbf{w}\|^2 \\ y^k f(\mathbf{x}^k) - 1 + \xi_k \geq 0 \\ \xi_k \geq 0 \end{cases}$$



Dual problems (in α space)

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

$$H = [y_k y_h k(\mathbf{x}^k, \mathbf{x}^h)]$$

1) Hinge loss version:

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

$$0 \leq \alpha_k \leq C \text{ box constraint}$$

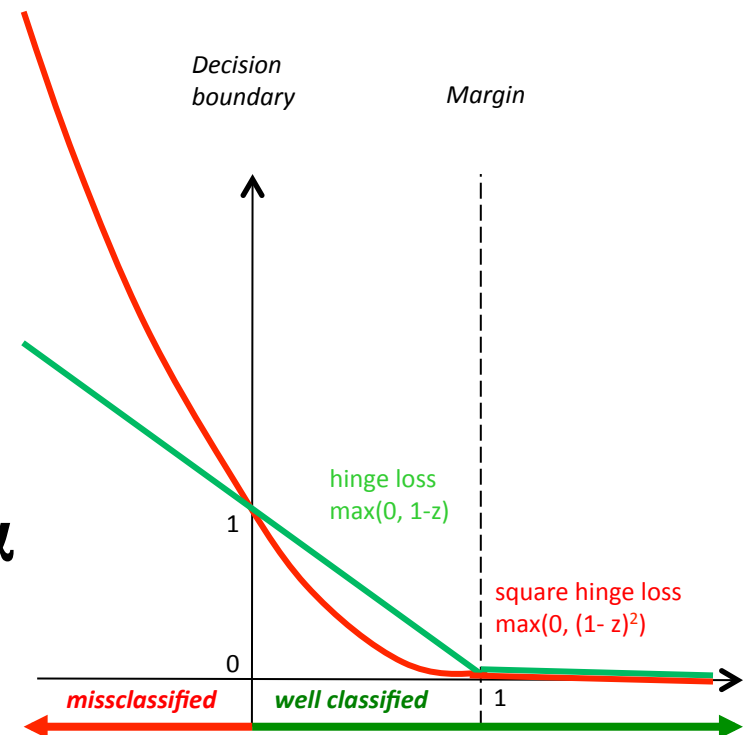
$$\alpha^T \mathbf{y} = 0 \text{ bias constraint}$$

2) Square hinge loss version

$$\min_{\alpha} \alpha^T \mathbf{1} - (1/2) \alpha^T (H + (1/C) I) \alpha$$

$$\alpha_k \geq 0$$

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



Dual problems (in α space)

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

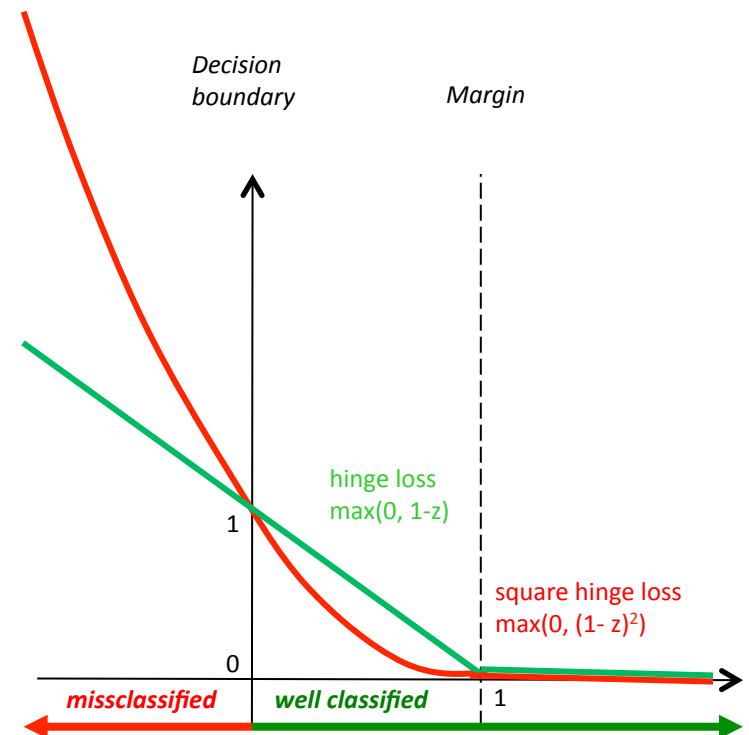
$$K = [k(\mathbf{x}^k, \mathbf{x}^h)]$$

1) Hinge loss version:

$$\begin{cases} \min_{\alpha} \alpha^T \mathbf{y} - (1/2) \alpha^T K \alpha \\ 0 \leq \alpha_k \leq 1/\lambda \quad \text{box constraint} \\ \alpha^T \mathbf{1} = 0 \quad \text{bias constraint} \end{cases}$$

2) Square hinge loss version

$$\begin{cases} \min_{\alpha} \alpha^T \mathbf{y} - (1/2) \alpha^T (K + \lambda I) \alpha \\ \alpha_k \geq 0 \quad \text{ridge} \\ \alpha^T \mathbf{1} = 0 \quad \text{bias constraint} \end{cases}$$



Summary

- $f(\mathbf{x}) = \sum_{k=1:N} \alpha_k k(\mathbf{x}, \mathbf{x}_k)$
- Kernel methods are inspired by non-parameteric example-based methods; each example \mathbf{x}_k votes according to:
 - How similar it is to \mathbf{x} measured by $k(\mathbf{x}, \mathbf{x}_k)$
 - Its weight α_k
- Kernels are “dot products” in a (possibly infinite) Φ space.
- The validity of a kernel can be checked with Mercer’s condition; there are also kernel composition theorems.
- The weights α_k can be optimized by gradient descent (for big data) or convex optimization methods.

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

Next time

Performance evaluation

