



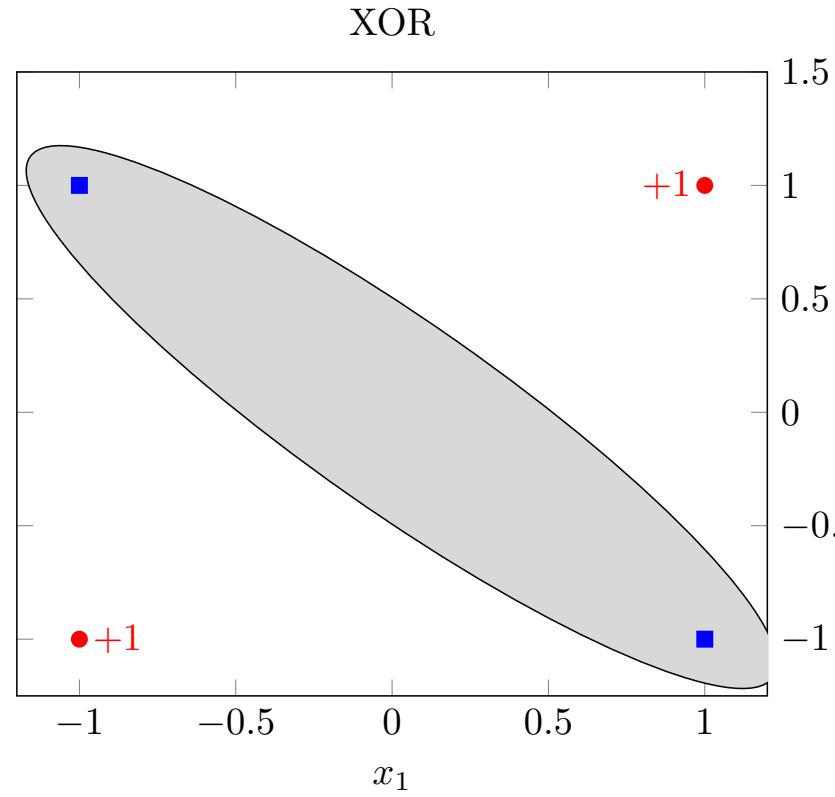
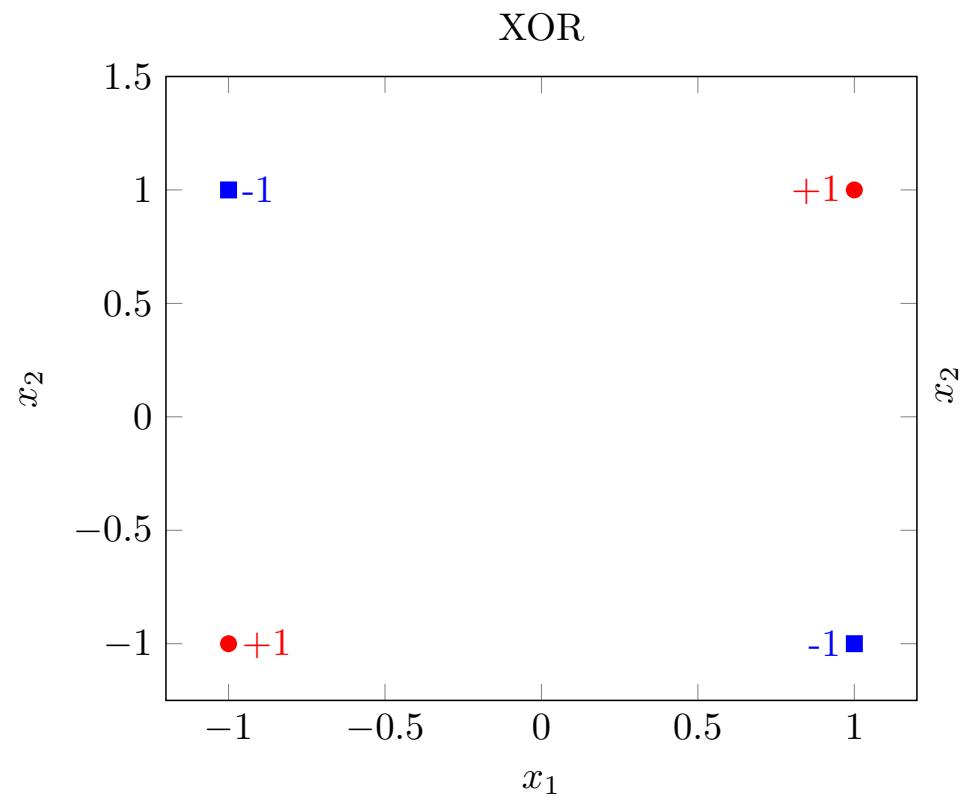
# CS489/698: Intro to ML

## Lecture 08: Kernels

# Outline

- Feature map
- Kernels
- The Kernel Trick
- Advanced

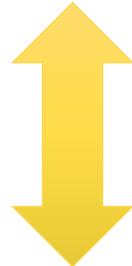
# XOR revisited



# Quadratic classifier

Weights  
(to be learned)

$$\mathbf{x}^\top \boxed{Q} \mathbf{x} + \sqrt{2} \mathbf{x}^\top \boxed{\mathbf{p}} + \boxed{\gamma} \geq 0$$



$$\hat{y} = f(\mathbf{x}) = 1$$

# The power of lifting

$$\mathbf{x}^\top Q \mathbf{x} + \sqrt{2} \mathbf{x}^\top \mathbf{p} + \gamma \geq 0$$

||

$$\mathbf{w}^\top \phi(\mathbf{x}) \geq 0$$

$$\mathbb{R}^d \rightarrow \mathbb{R}^{d^*d+d+1}$$

Feature map

$$\phi(\mathbf{x}) = \begin{bmatrix} \overrightarrow{\mathbf{x}\mathbf{x}}^\top \\ \sqrt{2}\mathbf{x} \\ 1 \end{bmatrix} \quad \text{with } \begin{bmatrix} \overrightarrow{\mathbf{x}_1\mathbf{x}} \\ \overrightarrow{\mathbf{x}_2\mathbf{x}} \\ \vdots \\ \overrightarrow{\mathbf{x}_d\mathbf{x}} \end{bmatrix}$$

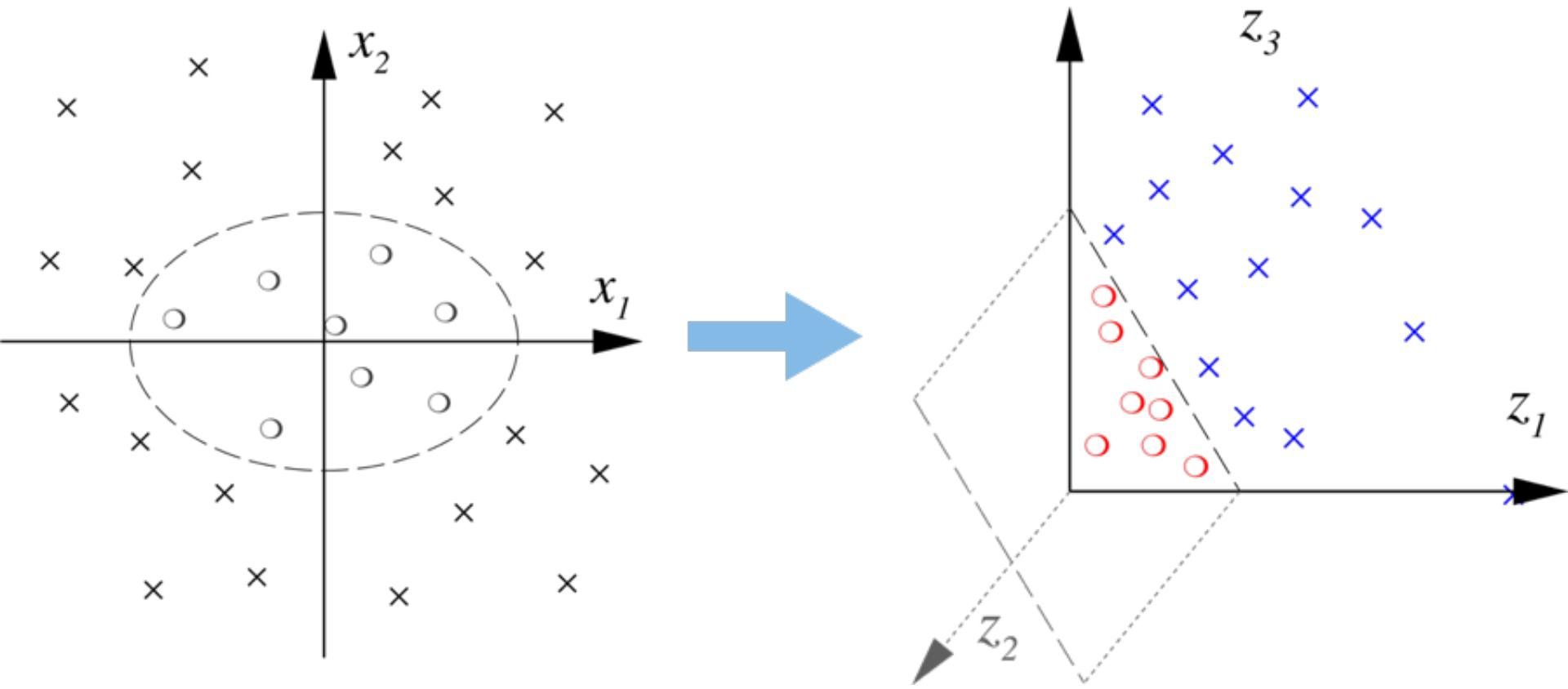
$$\mathbf{w} = \begin{bmatrix} \overrightarrow{Q} \\ \mathbf{p} \\ \gamma \end{bmatrix}$$

# Example

$$\phi(\mathbf{x}) = [x_1^2, \sqrt{2}x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1]$$

$$\phi(\mathbf{x}) = [x_1^2, x_1x_2, x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1]$$

# Does it work?



# Curse of dimensionality?

$$\phi : \mathbf{R}^d \rightarrow \mathbf{R}^{d^2+d+1}$$

computation in this space now

$$\phi(\mathbf{x}) = \begin{bmatrix} \xrightarrow{\mathbf{x}\mathbf{x}^\top} \\ \sqrt{2}\mathbf{x} \\ 1 \end{bmatrix}$$

But, all we need is the dot product !!!

$$\begin{aligned}\phi(\mathbf{x})^\top \phi(\mathbf{x}') &= (\mathbf{x}^\top \mathbf{x}')^2 + 2\mathbf{x}^\top \mathbf{x}' + 1 \\ &= (\mathbf{x}^\top \mathbf{x}' + 1)^2\end{aligned}$$

- This is still computable in  $O(d)$ !

# Feature transform

$$\phi : \mathbf{R}^d \rightarrow \mathbf{R}^h$$

- NN: learn  $\varphi$  simultaneously with  $w$
- Here: choose a nonlinear  $\varphi$  so that for some  $f : \mathbf{R} \rightarrow \mathbf{R}$

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = \boxed{f(\mathbf{x}^\top \mathbf{x}')}}$$

save computation

# Outline

- Feature map
- Kernels
- The Kernel Trick
- Advanced

# Reverse engineering

- Start with some function  $k : \mathbf{R}^d \times \mathbf{R}^d \rightarrow \mathbf{R}$ , s.t. exists feature transform  $\varphi$  with

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = k(\mathbf{x}, \mathbf{x}')$$

- As long as  $k$  is efficiently computable, don't care the dim of  $\varphi$  (could be infinite!)
- Such  $k$  is called a (reproducing) **kernel**.

# Examples

- Polynomial kernel  $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}')^p$   
 $k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^\top \mathbf{x}' + 1)^p$
- Gaussian Kernel  
 $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2/\sigma)$
- Laplace Kernel  
 $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2/\sigma)$
- Matérn Kernel  
$$\frac{1}{2^{\nu-1}\Gamma(\nu)} \left( \frac{2\sqrt{\nu}\|\mathbf{x} - \mathbf{x}'\|_2}{\theta} \right)^\nu H_\nu \left( \frac{2\sqrt{\nu}\|\mathbf{x} - \mathbf{x}'\|_2}{\theta} \right)$$

Verifying a kernel  $\kappa: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\exists \varphi$  s.t.  
 $\forall x, x', \kappa(x, x') = \varphi(x)^T \varphi(x')$

For any  $n$ , for any  $x_1, x_2, \dots, x_n$ , the **kernel matrix  $K$**  with

$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$$

is symmetric and positive semidefinite ( $K \in \mathbb{S}_+^d$ )

*necessity*:  $\sum K_{ij} = \varphi(x_i)^T \varphi(x_j) = \varphi(x_j)^T \varphi(x_i) = K_{ji}$

- Symmetric:  $K_{ij} = K_{ji}$  (2)  $\sum_i \sum_j \alpha_i \alpha_j \varphi(x_i)^T \varphi(x_j) = \left\| \sum_i \alpha_i \varphi(x_i) \right\|^2$
- Positive semidefinite (PSD): for all  $\boldsymbol{\alpha} \in \mathbb{R}^n$

$$\boldsymbol{\alpha}^\top K \boldsymbol{\alpha} = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j K_{ij} \geq 0$$

Yao-Liang Yu

# Kernel calculus

$$\varphi(x)^T \varphi(x')$$

$$\Rightarrow \varphi$$

$$\cancel{\varphi} = \sqrt{\lambda} \varphi$$

- If  $k$  is a kernel, so is  $\lambda k$  for any  $\lambda \geq 0$

$$(x)^T (\sqrt{\lambda} q(x))$$

$$K \text{ is psd} \Rightarrow \forall \alpha \in \mathbb{R} \quad \bar{\alpha}^T K \alpha \geq 0$$

- If  $k_1$  and  $k_2$  are kernels, so is  $k_1 + k_2$

$$\bar{\alpha}^T (\lambda K)$$

$$\lambda \cdot \bar{\alpha}^T K \alpha > 0$$

- If  $k_1$  and  $k_2$  are kernels, so is  $k_1 k_2$

$$\sum_i \gamma_i K_i \text{ is a kernel if } K_i \text{ is}$$

and  $\gamma_i \geq 0$

# Outline

- Feature map
- Kernels
- The Kernel Trick
- Advanced

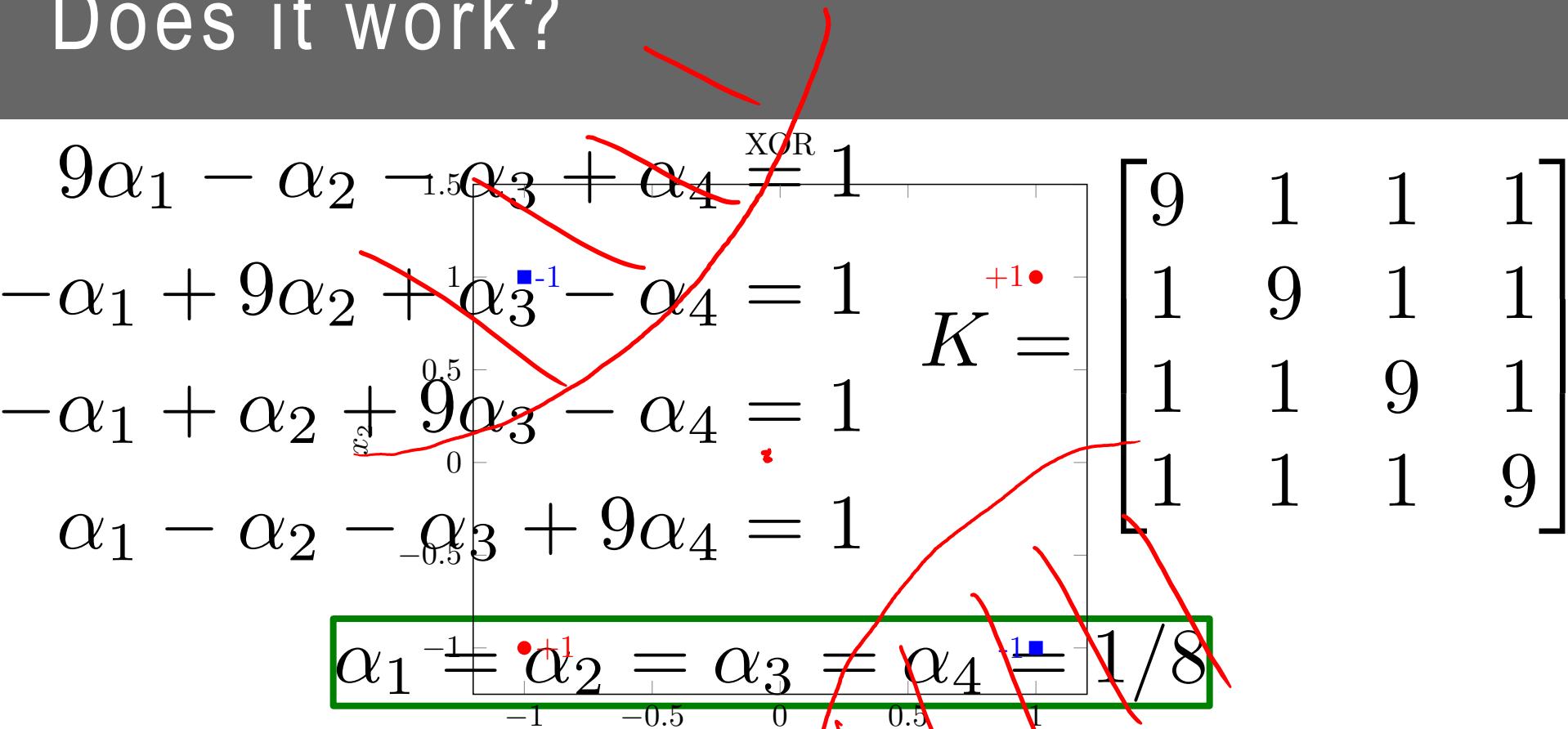
# Kernel SVM (dual)

$$\min_{C \geq \alpha \geq 0} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K_{ij} - \sum_{i=1}^n \alpha_i$$

$$\text{s.t. } \sum_i \alpha_i y_i = 0$$

With  $\alpha$ ,  $\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \phi(\mathbf{x}_i)$  but  $\phi$  is implicit...

# Does it work?



$$\mathbf{w} = [0, -1/\sqrt{2}, 0, 0, 0, 0]$$

$$\mathbf{w}^\top \phi(\mathbf{x}) = -x_1 x_2$$

$$\phi(\mathbf{x}) = [x_1^2, \sqrt{2}x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1]$$

# Testing

- Given test sample  $\mathbf{x}'$ , how to perform testing?

$$\mathbf{w}^\top \phi(\mathbf{x}') = \sum_{i=1}^n \alpha_i y_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}')$$

No explicit access  
to  $\varphi$ , again!

$$= \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}')$$

dual variables

training set

kernel

# Tradeoff

- Previously: training  $O(nd)$ , test  $O(d)$
- Kernel: **training  $O(n^2)$ , test  $O(n)$**
- Nice to avoid explicit dependence on  $h$  (could be inf)
- But if  $n$  is also large... ~~(maybe later)~~

# Learning the kernel (Lanckriet et al.'04)

$$\min_{C \geq \alpha \geq 0} \max_{\zeta \geq 0} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \left[ \sum_{s=1}^t \zeta_s K_{ij}^{(s)} \right] - \sum_{i=1}^n \alpha_i$$

s.t.  $\sum_i \alpha_i y_i = 0$

- Nonnegative combination of  $t$  pre-selected kernels, with coefficients  $\zeta$  simultaneously **learned**

# Logistic regression revisited

$$\min_{\mathbf{w} \in \mathbb{R}^d} \sum_i \log(1 + e^{-y_i \mathbf{w}^\top \mathbf{x}_i}) + \lambda \|\mathbf{w}\|_2^2$$

↓ kernelize

$$\min_{\mathbf{w} \in \mathbb{R}^h} \sum_i \log(1 + e^{-y_i \mathbf{w}^\top \phi(\mathbf{x}_i)}) + \lambda \|\mathbf{w}\|_2^2$$

$\mathcal{X} = \text{Span}\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$

$\mathbf{w}^\top \mathbf{x}_i = w_1^\top \mathbf{x}_i$

$\|\mathbf{w}\|_2^2 = ((w_1^2 + \dots + w_h^2))^{1/2}$

**Representer Theorem** (Wabha, Schölkopf, Herbrich, Smola, Di Nunzio)

The optimal  $\mathbf{w}$  has the following form:

$$\sum_j \alpha_j \mathbf{x}_{ij} = (\mathbf{K}\boldsymbol{\alpha})_i = K_{ii} \cdot \boldsymbol{\alpha}$$

$$\mathbf{w} = \sum_j \alpha_j \phi(\mathbf{x}_j)$$

$$\boldsymbol{\alpha}^\top \mathbf{K} \boldsymbol{\alpha}$$

# Kernel Logistic Regression (primal)

$f(\alpha_1, \alpha_2, \dots, \alpha_n)$

$$\min_{\alpha \in \mathbf{R}^n} \sum_i \log(1 + e^{-y_i \alpha^\top K_{:i}}) + \lambda \alpha^\top K \alpha$$

$$\min_{w \in \mathbf{R}^d} \sum_i \log(1 + e^{-y_i w^\top x_i}) + \lambda w^\top w$$

$$\alpha_{t+1} \leftarrow \alpha_t - \eta_t [\nabla^2 f(\alpha_t)]^{-1} \cdot \nabla f(\alpha_t)$$

$$\nabla^2 f(\alpha_t) = \sum_i p_i (1 - p_i) K_{:i} K_{:i}^\top + 2\lambda K$$

$$\nabla f(\alpha_t) = K^\top (p - \frac{y+1}{2}) + 2\lambda K \alpha_t$$

$$p_i = \frac{1}{1 + e^{-\alpha_t^\top K_{:i}}}$$

uncertain predictions get  
bigger weight

# Outline

- Feature map
- Kernels
- The Kernel Trick
- Advanced

# Universal approximation (Micchelli, Xu, Zhang'06)

**Universal kernel.** For any compact set  $Z$ , for any continuous function  $f: Z \rightarrow \mathbb{R}$ , for any  $\epsilon > 0$ , there exist  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  in  $Z$  and  $\alpha_1, \alpha_2, \dots, \alpha_n$  in  $\mathbb{R}$  such that

$$\max_{\mathbf{x} \in Z} \left| f(\mathbf{x}) - \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) \right| \leq \epsilon$$

decision boundary      kernel methods

**Example.** The Gaussian kernel.

$$\exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2/\sigma^2)$$

# Kernel mean embedding (Smola, Song, Gretton, Schölkopf, ...)

$$\mathbf{P} \mapsto \mu_{\mathbf{P}} := \mathbf{E}(\phi(X)), \text{ where } X \sim \mathbf{P}$$



feature map of some kernel

- **Characteristic kernel:** the above mapping is 1-1
- Completely preserve the information in the distribution  $\mathbf{P}$
- Lots of applications

# Questions?

