

# CS489/698: Intro to ML

## Lecture 08: Kernels

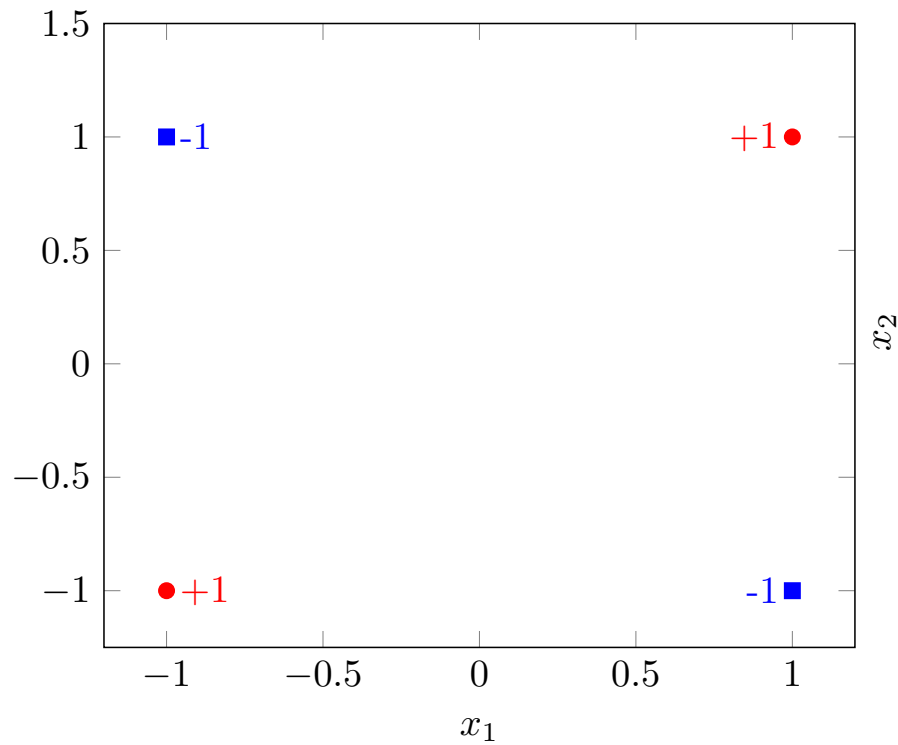


# Outline

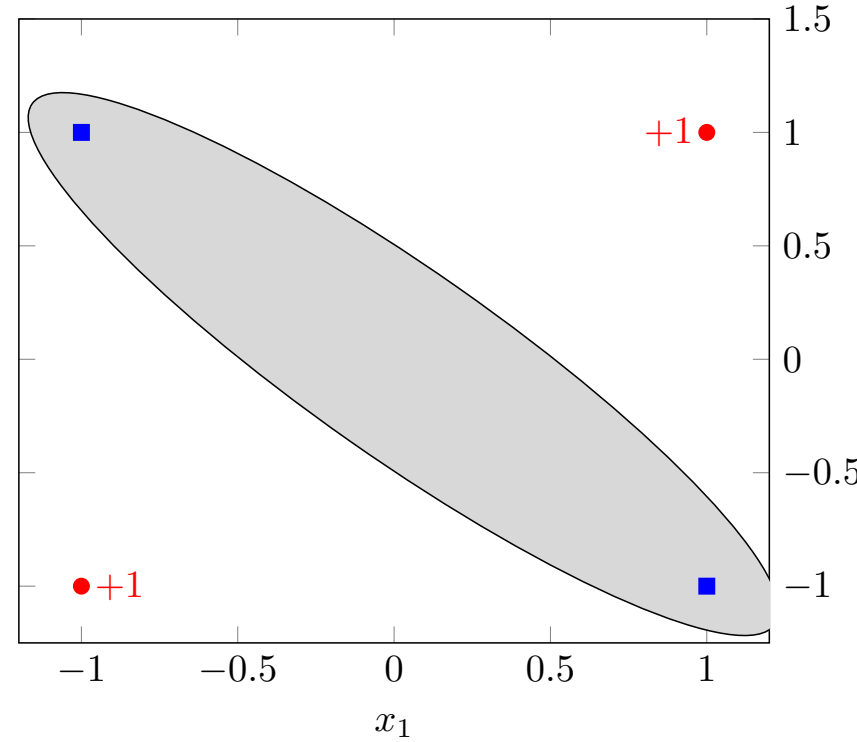
- Feature map
- Kernels
- The Kernel Trick
- Advanced

# XOR revisited

XOR



XOR



# 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

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



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

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

Feature map

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

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

# Example

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

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

# Curse of dimensionality?

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

computation in this space now

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

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

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

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  $\phi$  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  $\phi$  (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

For any  $n$ , for any  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , the **kernel matrix**  $K$  with

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

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

- Symmetric:  $K_{ij} = K_{ji}$
- Positive semidefinite (PSD): for all  $\boldsymbol{\alpha} \in \mathbf{R}^n$

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

Yao-Liang Yu

# Kernel calculus

- If  $k$  is a kernel, so is  $\lambda k$  for any  $\lambda \geq 0$
- If  $k_1$  and  $k_2$  are kernels, so is  $k_1 + k_2$
- If  $k_1$  and  $k_2$  are kernels, so is  $k_1 k_2$

# Outline

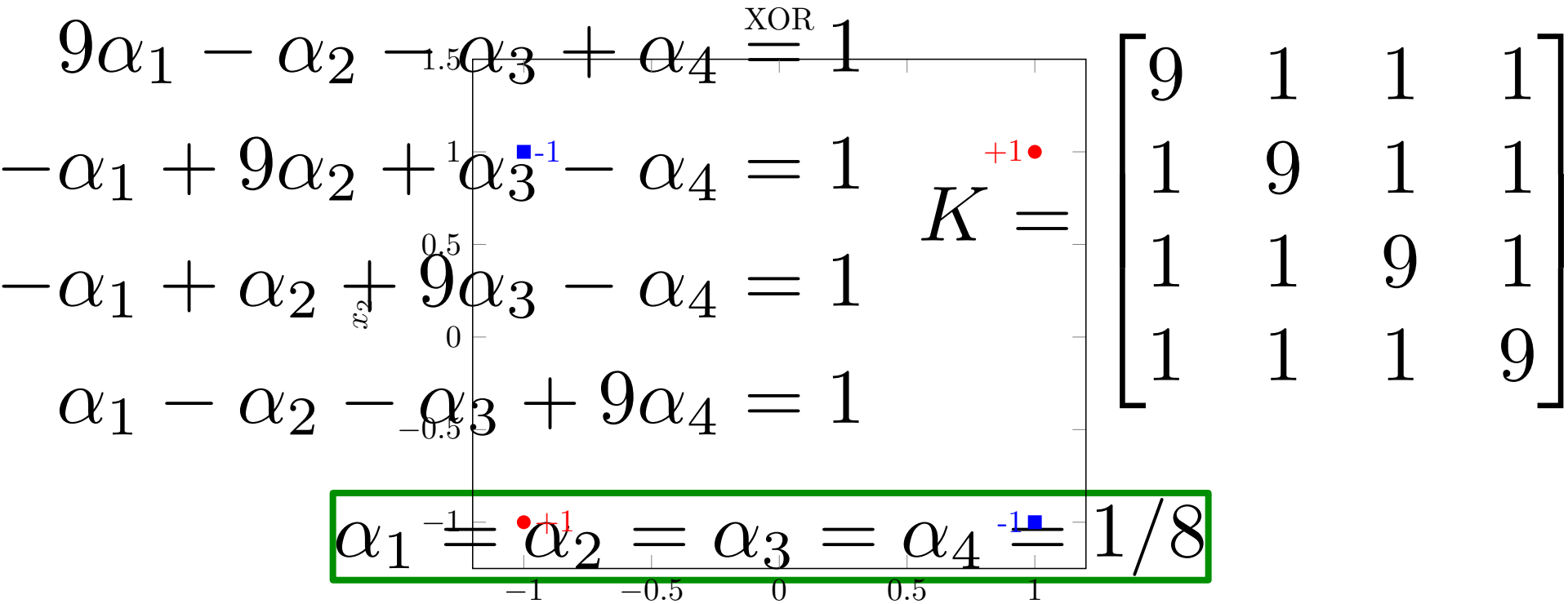
- Feature map
- Kernels
- The Kernel Trick
- Advanced

# Kernel SVM (dual)

$$\begin{aligned} \min_{C \geq \alpha \geq 0} \quad & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K_{ij} - \sum_{i=1}^n \alpha_i \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \end{aligned}$$

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

# Does it work?



$$\mathbf{w} = [0, -1/\sqrt{2}, 0, 0, 0, 0] \quad \mathbf{w}^\top \phi(\mathbf{x}) = -x_1 x_2$$

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



# 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  $\phi$ , 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)

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

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

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

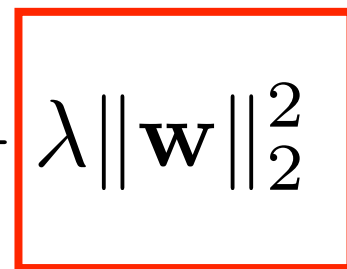
# Logistic regression revisited

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



kernelize

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



**Representer Theorem** (Wabha, Schölkopf, Herbrich, Smola, Dinuzzo, ...).

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

$$\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_i)$$



# Kernel Logistic Regression (primal)

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

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

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

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

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

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  $f: Z \rightarrow \mathbf{R}$ , for any  $\varepsilon > 0$ , there exist  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  in  $Z$  and  $\alpha_1, \alpha_2, \dots, \alpha_n$  in  $\mathbf{R}$  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.

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

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

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



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?

