

CS480/680: Intro to ML

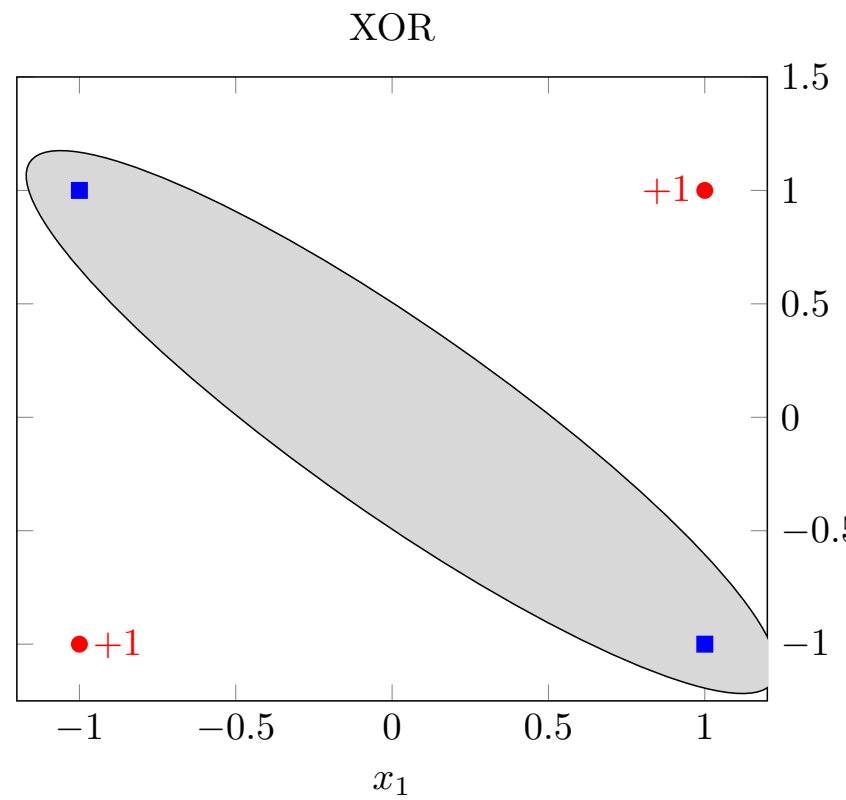
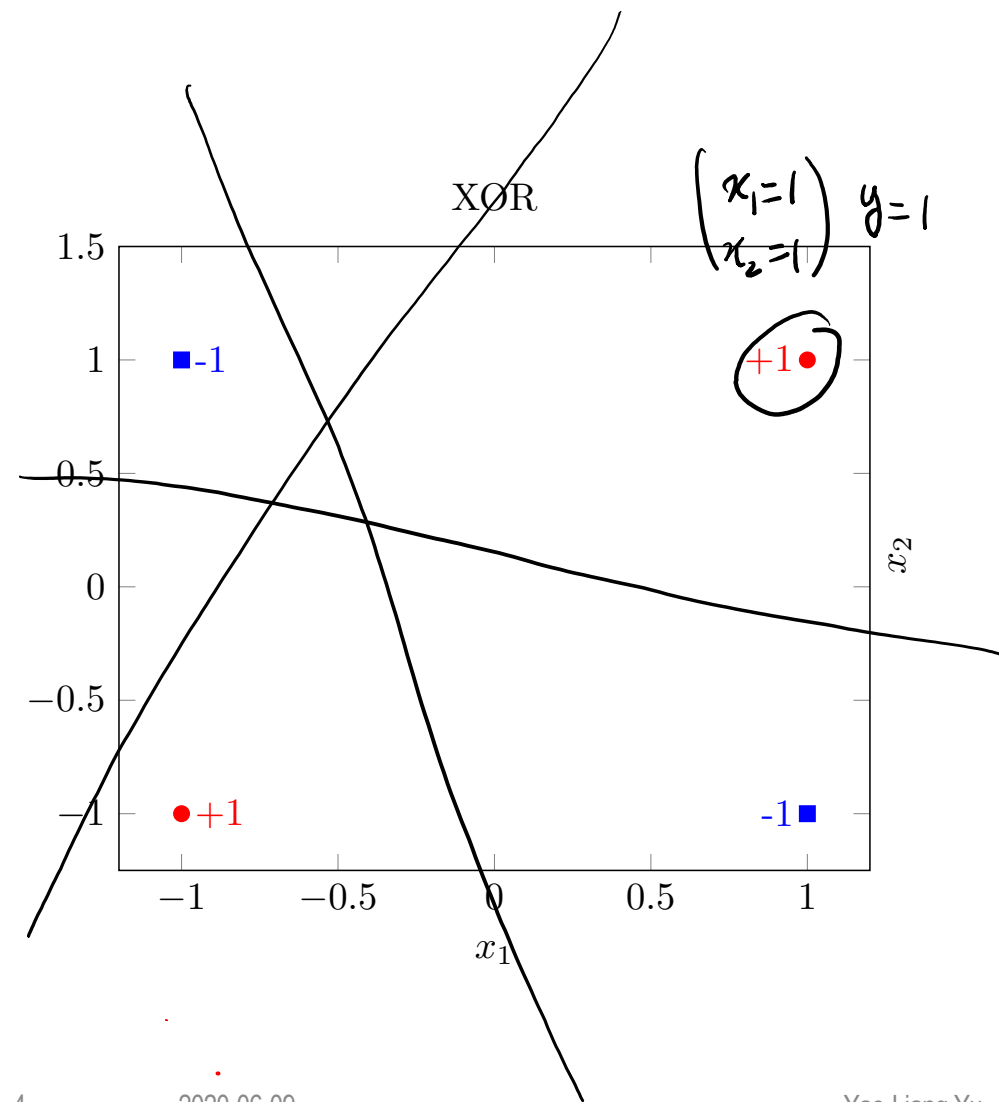
Lecture 09: Kernels



Outline

- Feature map
- Kernels
- The Kernel Trick

XOR



Quadratic classifier

Weights
(to be learned)

$$\sum_{i,j} x_i Q_{ij} x_j$$
$$\mathbf{x}^T \boxed{Q} \mathbf{x} + \sqrt{2} \mathbf{x}^T \boxed{\mathbf{p}} + \boxed{\gamma} \geq 0$$



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

The power of lifting

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\text{tr}(x^T Q x) = \underbrace{x^T Q x + \sqrt{2} x^T p + \gamma}_{\geq 0}$$

$$\begin{matrix} \parallel \\ \text{tr}(Q \cdot x x^T) \\ R^d \Rightarrow R^{d \times d+d+1} \\ \langle Q, x x^T \rangle \end{matrix}$$

Feature map

$$w^T \phi(x) \geq 0$$

$$\phi(x)^T w \geq 0$$

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

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

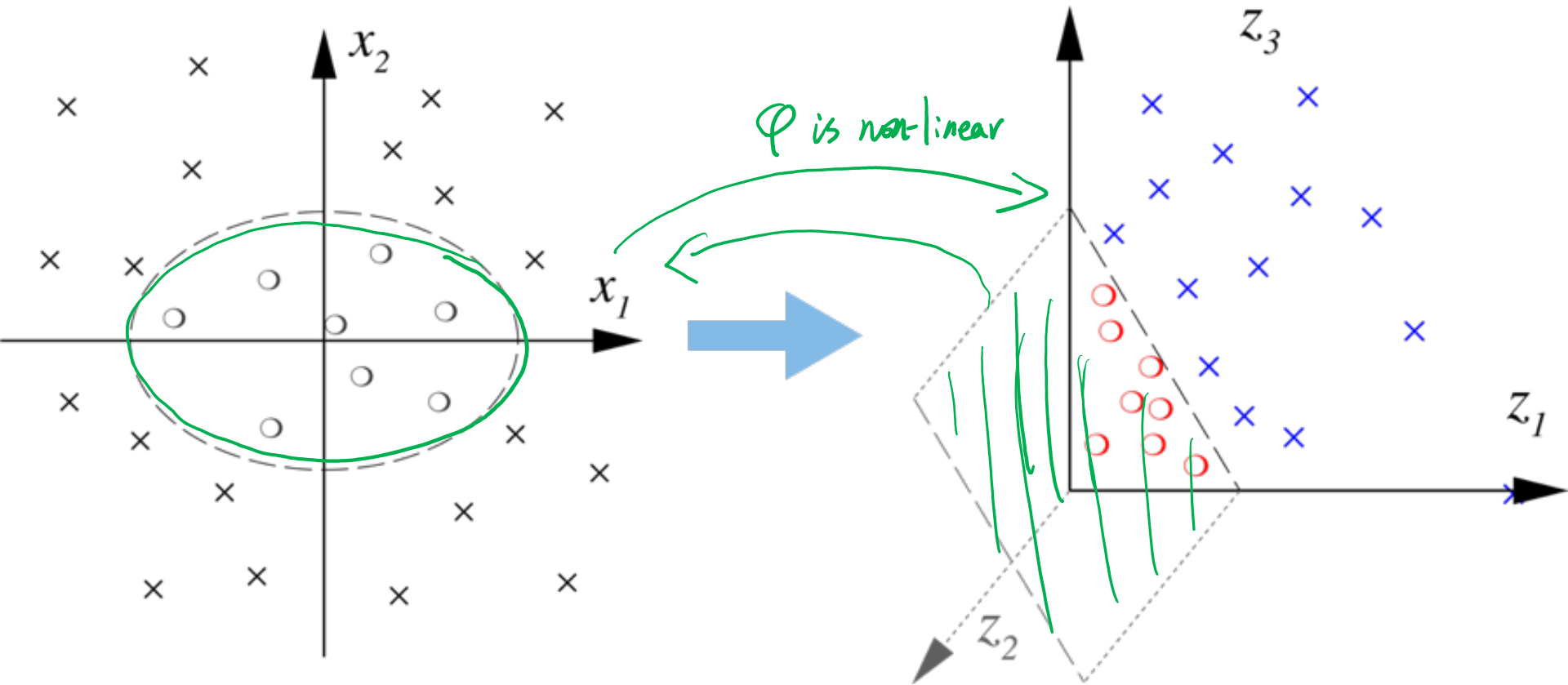
Example

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

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

Feature map may not be unique

Does it work?



Linear hyperplane in the feature space corresponds to nonlinear boundary in the original space!

And (almost) vice versa!!

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 !!!

$$\rightarrow \phi(\mathbf{x})^\top \phi(\mathbf{x}') = \frac{(\mathbf{x}^\top \mathbf{x}')^2 + 2\mathbf{x}^\top \mathbf{x}' + 1}{1}$$

dot prod between \mathbf{x} & $\mathbf{x}' = \frac{(\mathbf{x}^\top \mathbf{x}' + 1)^2}{1}$

• This is still computable in $O(d)$ (LHS: d^2 , RHS: dd)

Feature transform

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

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

$$\phi(\mathbf{x})^\top \phi(\mathbf{x}') = \underbrace{f(\mathbf{x}^\top \mathbf{x}')}_{\text{save computation}}$$

save computation

Outline

- Feature map
- **Kernels**
- The Kernel Trick

Reverse engineering

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

$\sin \theta$
 $\exp(t) = \sum_k \frac{t^k}{k!}$

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2 / \sigma)$$

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

- Laplace Kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2 / \sigma)$$

- Matérn Kernel

$$\frac{1}{2^{\nu-1} \Gamma(\nu)} \left(\frac{2\sqrt{\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\theta} \right)^\nu \underline{H_\nu} \left(\frac{2\sqrt{\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\theta} \right)$$

Verifying a kernel

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

$$K_{ij} = \underline{k}(\mathbf{x}_i, \mathbf{x}_j) \quad n \times n$$

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

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

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

Kernel calculus

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

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

- k_1 with φ_1 , k_2 with φ_2 $\exists \varphi_i: k_i(x, x') = \varphi_i(x)^T \varphi_i(x')$, $i \in \{1, 2\}$
- k_1+k_2 with ??

define $\varphi = \begin{pmatrix} \varphi_1 \\ \varphi_2 \end{pmatrix}$, then $\varphi(x)^T \varphi(x') = k_1(x, x') + k_2(x, x')$

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

Outline

- Feature map
- Kernels
- **The Kernel Trick**

Kernel SVM (dual)

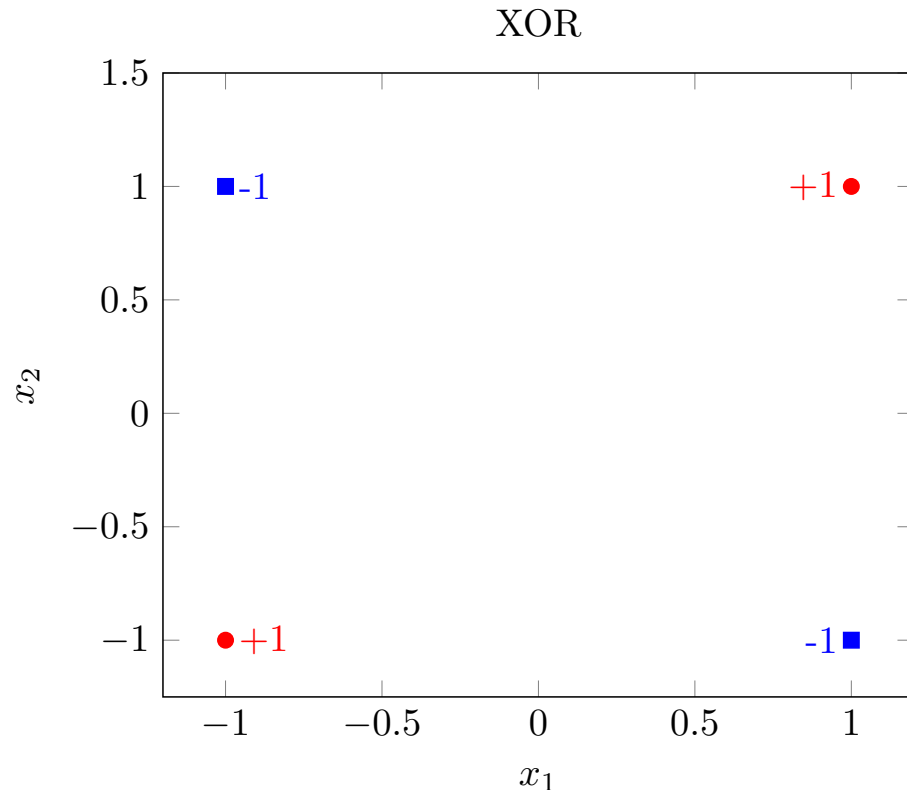
$$\min_{C \geq \alpha \geq 0} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K_{ij} - \sum_{i=1}^n \alpha_i$$

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

$x_i^T x_j \rightarrow \phi(x_i)^T \phi(x_j)$
" "
 $\kappa(x_i, x_j)$
" "
 $\underline{\kappa_{ij}}$

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

Does it work?

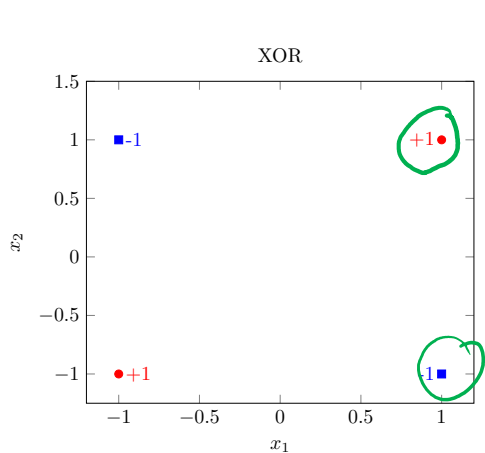


$$\underline{\phi(\mathbf{x})} = [x_1^2, \sqrt{2}x_1x_2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1]$$

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^2$$

Does it work?

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



$$k_{ii} = \left[\begin{pmatrix} 1 \\ 1 \end{pmatrix}^T \begin{pmatrix} 1 \\ 1 \end{pmatrix} + 1 \right]^2 = 9$$

$$K = \begin{bmatrix} 9 & 1 & 1 & 1 \\ 1 & 9 & 1 & 1 \\ 1 & 1 & 9 & 1 \\ 1 & 1 & 1 & 9 \end{bmatrix}$$

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + 1)^2$$

Does it work?

$$9\alpha_1 - \alpha_2 - \alpha_3 + \alpha_4 = 1$$

$$-\alpha_1 + 9\alpha_2 + \alpha_3 - \alpha_4 = 1$$

$$-\alpha_1 + \alpha_2 + 9\alpha_3 - \alpha_4 = 1$$

$$\alpha_1 - \alpha_2 - \alpha_3 + 9\alpha_4 = 1$$

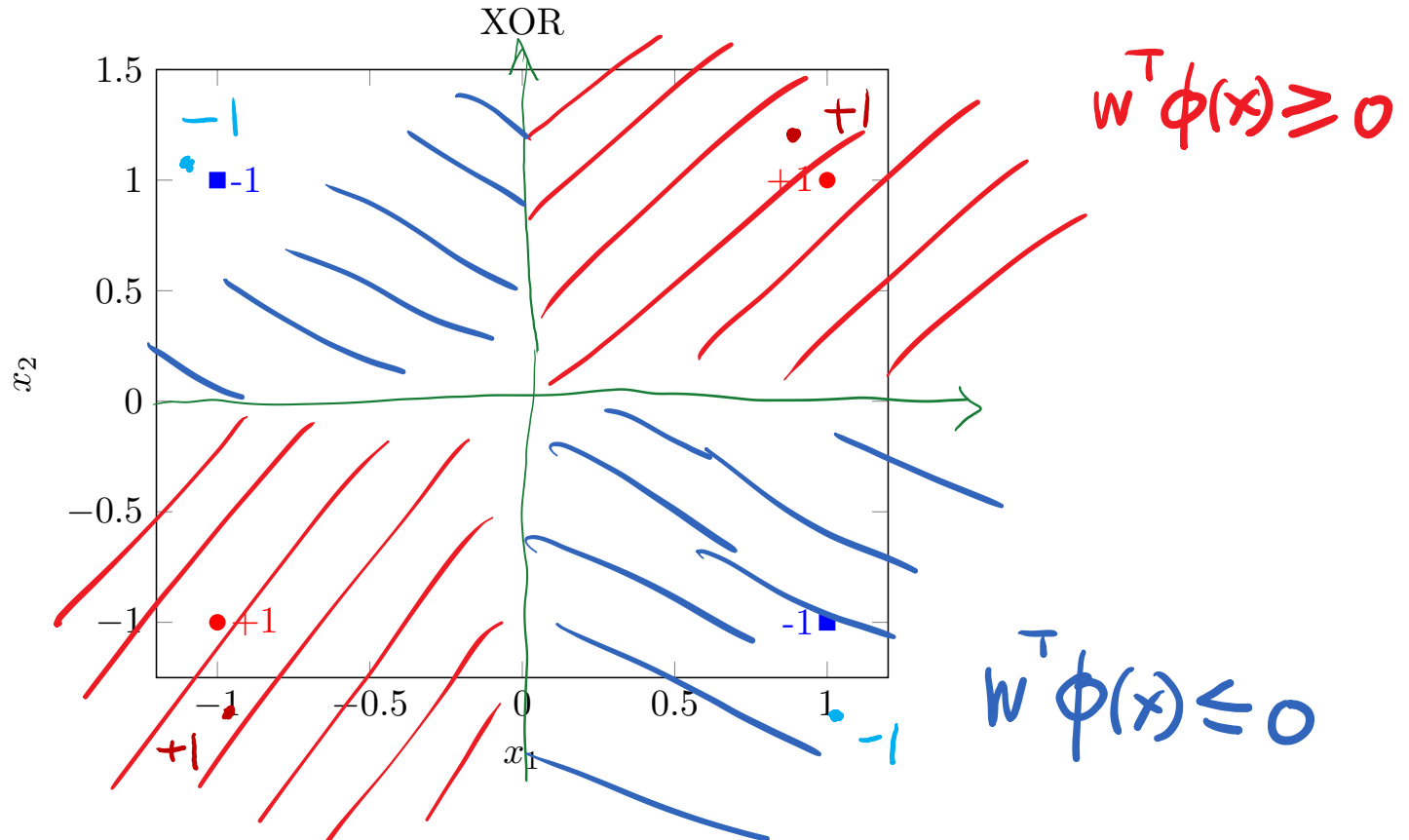
$$W = \sum_i \alpha_i y_i \phi(x_i)$$

$$\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 1/8$$

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

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

Does it work?



$$w^T \phi(\mathbf{x}) = x_1 x_2 \geq 0$$
$$\leq 0$$

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 ϕ , 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 $O(n^2d)$, test $O(nd)$
- 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. $\sum_i \alpha_i y_i = 0$

Annotations: A red box highlights the inner maximization over $\zeta \succeq 0$. A yellow arrow points from the \min operator to this box. A yellow circle highlights the variable t in the summation index. A yellow wavy underline is drawn under the summation $\sum_{s=1}^t \zeta_s K_{ij}^{(s)}$.

- Nonnegative combination of t pre-selected kernels, with coefficients ζ 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

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

$\sum_j \alpha_j y_j \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_i)$
 K_{ji}

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

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

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

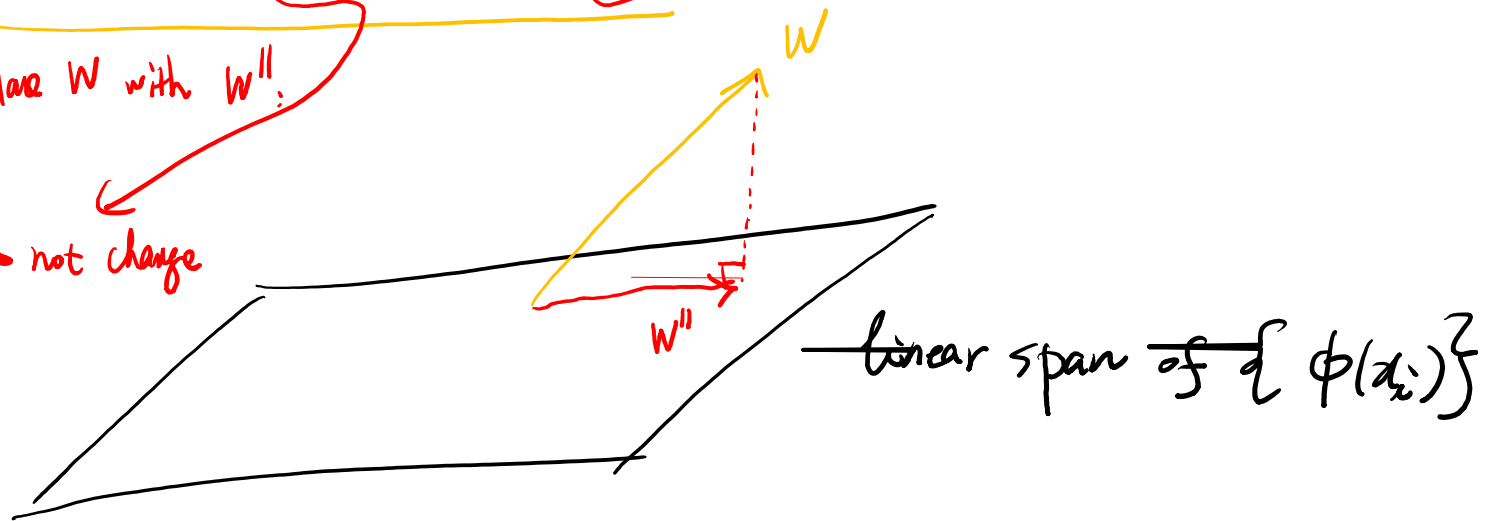
Orthogonal decomposition

$$\min_w \sum_i \ell(y_i w^T \phi(x_i)) + \frac{\lambda}{2} \|w\|_2^2$$

replace w with w''

does not change

decrease



$$w = \sum_i \beta_i \phi(x_i) = \sum_i \underbrace{\alpha_i y_i}_{\beta_i} \phi(x_i)$$

Questions?

