

CSC 580 Principles of Machine Learning

# 08 Kernel methods

**Chicheng Zhang**

**Department of Computer Science**



\*slides credit: built upon CSC 580 Fall 2021 lecture slides by Kwang-Sung Jun

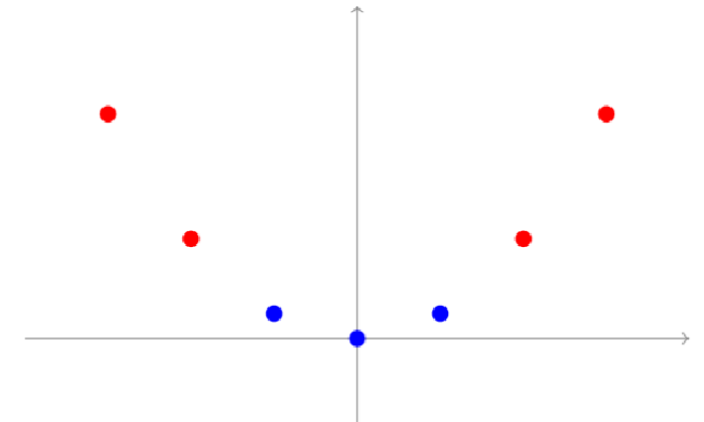
# Recall: Beyond linearity

- Recall: for 1d problem, we embedded the feature:  $x' = (x, 1) \in \mathbb{R}^2$
- Actually, the embedding trick is stronger.
  - $(x^2, x, 1)$ : 2<sup>nd</sup> order polynomial
  - $(x^d, x^{d-1}, \dots, 1)$ : d-th order polynomial (= degree d)
- Higher order => strictly larger class of function  $\mathcal{F}$

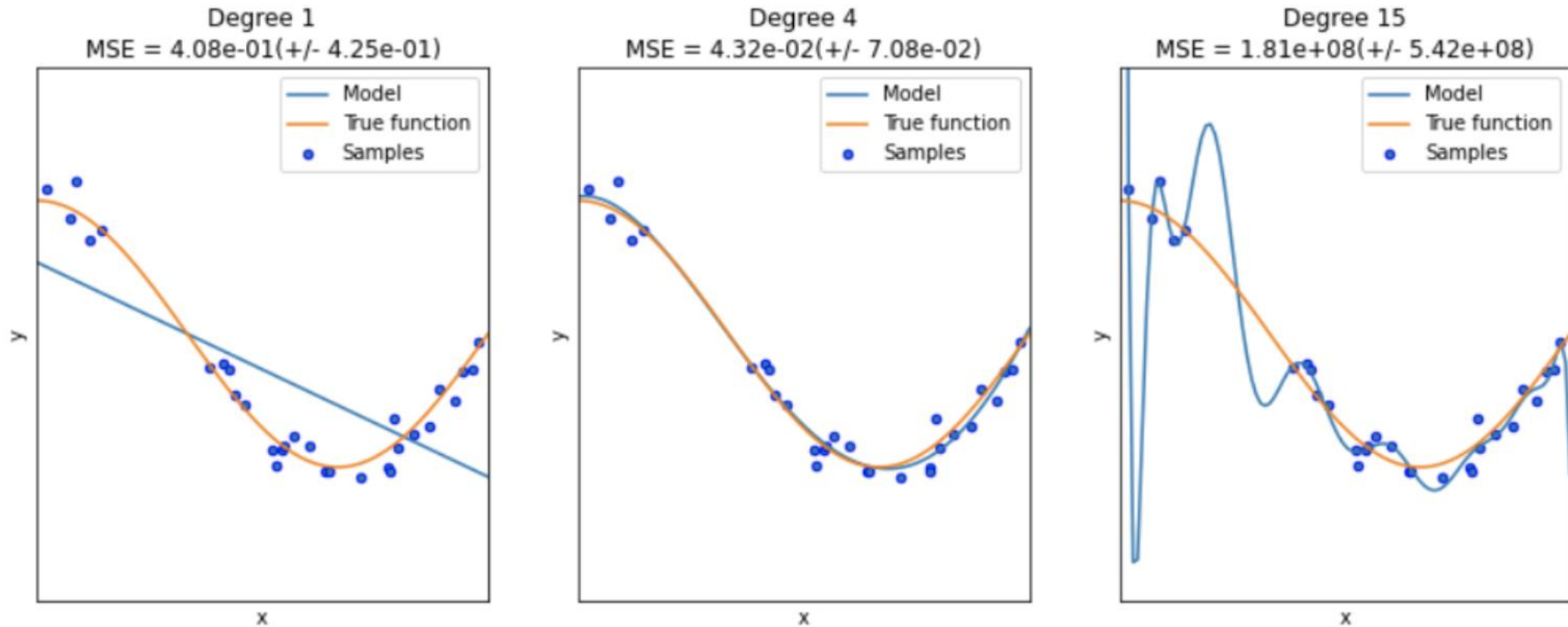
- The following samples in  $\mathbb{R}$  are not separable by halfspaces



- But, if we map  $x \rightarrow (x, x^2)$  it is separable by halfspaces



# Recall: Feature embedding trick



- overfitting vs underfitting
- bias-variance tradeoff.

$$\text{err}(\hat{f}) = [\text{err}(\hat{f}) - \min_{f^* \in \mathcal{F}} \text{err}(f^*)] + \min_{f^* \in \mathcal{F}} \text{err}(f^*)$$

# Kernel trick: high-level idea

- Given (possibly nonlinear) basis functions  $\phi(x): \mathbb{R}^d \rightarrow \mathbb{R}^D$ , where  $D$  is huge or infinite
- Would like to learn a model from class  $\mathcal{F}_\phi = \{h: h(x) = \langle w, \phi(x) \rangle, \text{ for some } w \in \mathbb{R}^D\}$  with running time independent of  $D$
- Computational Challenge:
  - a naïve application of existing algorithms (e.g. SGD, Perceptron) has running time  $\Omega(D)$
- Key structural assumption on  $\phi$ :
  - its induced *kernel function*  $K(x, x') := \langle \phi(x), \phi(x') \rangle$  can be evaluated in time independent of  $D$
- How can we utilize this structure to address the challenge?

# Kernel function: an example

- Let  $|x| < 1$
- $\phi(x) = (1, x, x^2, x^3 \dots) = (x^n)_{n=1}^{\infty}$ 
  - Impossible to write down explicitly
- Induced Kernel function:

$$K(x, y) := \langle \phi(x), \phi(y) \rangle = \sum_{n=1}^{\infty} (x \cdot y)^n = \frac{1}{1 - xy}$$

- Takes  $O(1)$  time to calculate

# How to use kernels with SVM

**Primal**

$$\min_{w,b} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

s.t.  $y_i(w^\top \phi(x_i) + b) \geq 1 - \xi_i, \forall i$

**Dual**

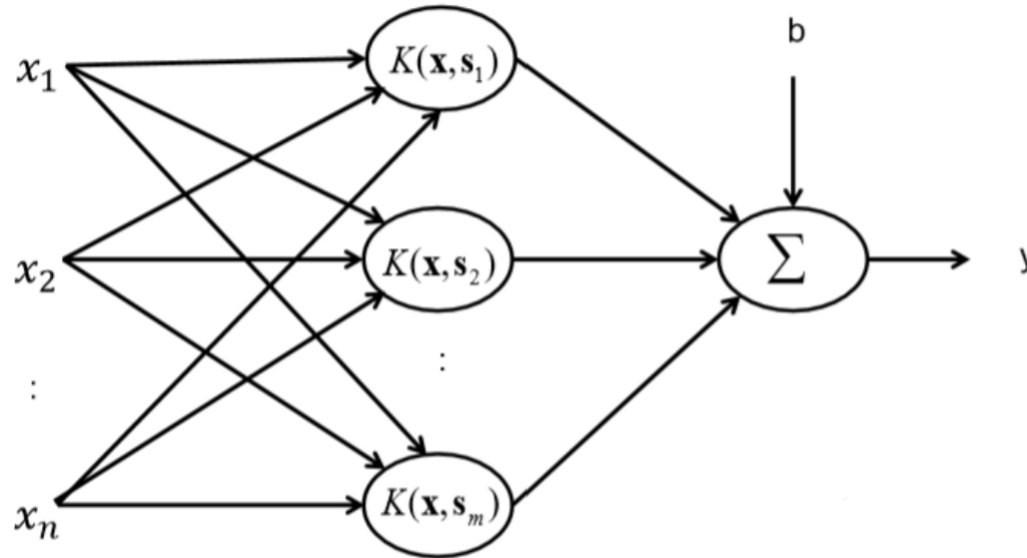
$$\max_{0 \leq \alpha_i \leq C} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle \phi(x_i), \phi(x_j) \rangle$$

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

- To make predictions
  - Recall optimal  $w^* = \sum_i \alpha_i^* y_i \phi(x_i)$
  - $\langle w^*, \phi(x) \rangle = \sum_i \alpha_i^* y_i \langle \phi(x_i), \phi(x) \rangle = \sum_i \alpha_i^* y_i K(x_i, x)$
  - Take  $\text{sign}(\langle w^*, \phi(x) \rangle + b^*)$
  - Interpretation:  $\alpha_i^*$ : weights,  $y_i$ : votes,  $K(x_i, x_*)$ : similarity

# How to use kernels with SVM (cont'd)

$$h(x) = \text{sign}(\langle w^*, \phi(x) \rangle + b^*) = \text{sign}(\sum_i \alpha_i^* y_i K(x_i, x) + b^*)$$



- Summary

- training: compute  $\alpha_i^*$ 's and  $b^*$
- test: compute the kernel functions  $K(x_i, x)$  and then take weighted combination  $\sum_i \alpha_i^* y_i K(x_i, x) + b^*$

Cf. weighted k-NN

$$\hat{y}(x) = \arg \max_y \sum_{i \in N(x_*)} w_i 1\{y_i = y\}$$

$$= \text{sign}(\sum_{i \in N(x_*)} w_i y_i),$$

where e.g.,  $w_i = \exp(-\beta d(x_i, x)^2)$

# Polynomial kernels

- $K(x, x') = (1 + \langle x, x' \rangle)^k \Rightarrow$  Q: is this a valid kernel?

- E.g., if  $x = (x_1, x_2)$  and  $z = (z_1, z_2)$ ,

$$(1 + x_1 z_1 + x_2 z_2)^2 = 1 + x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 + 2x_2 z_2 + 2x_1 x_2 z_2 z_1$$

$$= \langle \phi(x), \phi(z) \rangle$$

$$\text{where } \phi(x) = (x_1^2, x_2^2, \sqrt{2} x_1, \sqrt{2} x_2, \sqrt{2} x_1 x_2, 1)$$



# Polynomial kernels

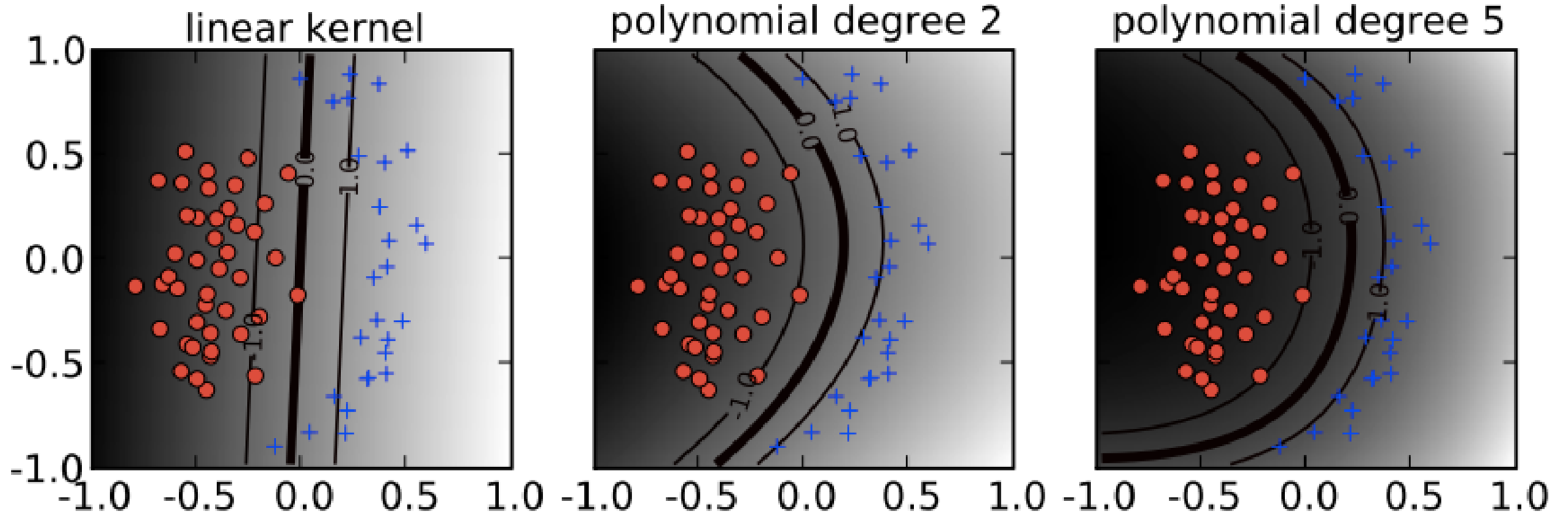


Figure from Ben-Hur & Weston,  
*Methods in Molecular Biology* 2010

# Gaussian/RBF kernels

- $K(x, x') = \exp\left(-\frac{\|x-x'\|^2}{2\sigma^2}\right)$  (often parameterized as  $\exp(-\gamma\|x-x'\|^2)$ )
- How can we show that this is a valid kernel?  $\Rightarrow$  We should find  $\phi(x)$  that results in  $K(x, x') = \langle \phi(x), \phi(x') \rangle$

Assume  $x \in R^1$  and  $\gamma > 0$ .

$$\begin{aligned} e^{-\gamma\|x_i-x_j\|^2} &= e^{-\gamma(x_i-x_j)^2} = e^{-\gamma x_i^2 + 2\gamma x_i x_j - \gamma x_j^2} \\ &= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 + \frac{2\gamma x_i x_j}{1!} + \frac{(2\gamma x_i x_j)^2}{2!} + \frac{(2\gamma x_i x_j)^3}{3!} + \dots\right) \\ &= e^{-\gamma x_i^2 - \gamma x_j^2} \left(1 \cdot 1 + \sqrt{\frac{2\gamma}{1!}} x_i \cdot \sqrt{\frac{2\gamma}{1!}} x_j + \sqrt{\frac{(2\gamma)^2}{2!}} x_i^2 \cdot \sqrt{\frac{(2\gamma)^2}{2!}} x_j^2 \right. \\ &\quad \left. + \sqrt{\frac{(2\gamma)^3}{3!}} x_i^3 \cdot \sqrt{\frac{(2\gamma)^3}{3!}} x_j^3 + \dots\right) = \phi(x_i)^T \phi(x_j), \end{aligned}$$

where

$$\phi(x) = e^{-\gamma x^2} \left[1, \sqrt{\frac{2\gamma}{1!}} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \dots\right]^T.$$

recall how we make **predictions**:

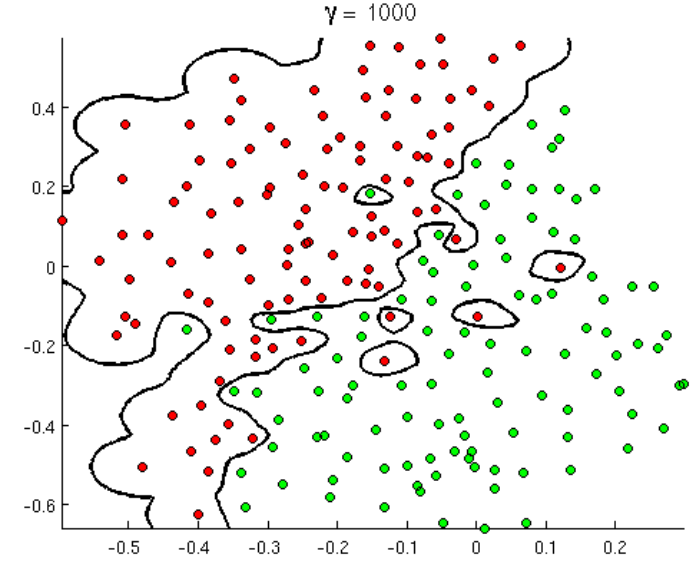
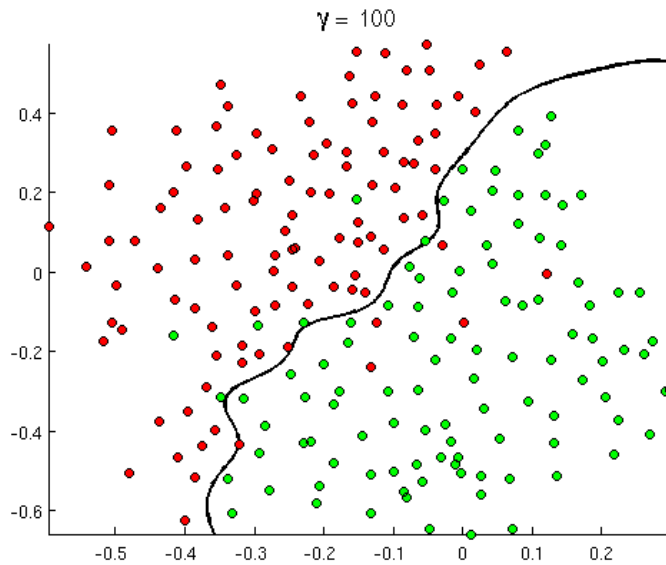
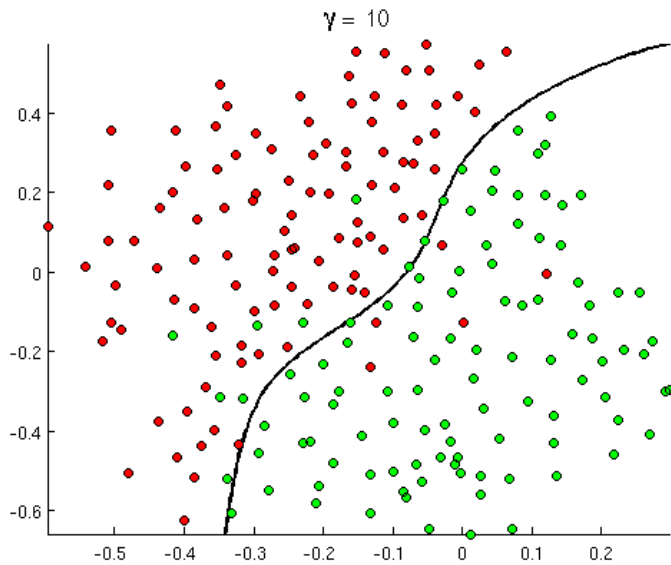
$$w^T x_* = \sum_i \alpha_i y_i x_i^T x_* = \sum_i \alpha_i y_i K(x_i, x_*)$$

**weighted k-NN**

- $\arg \max_y \sum_{i \in N(x_*)} w_i 1\{y_i = y\}$
- e.g.,  $w_i = \exp\left(-\beta \cdot (d(x_i, x_*))^2\right)$

# Gaussian kernel

- $\gamma = \frac{1}{2\sigma^2}$



- Larger  $\gamma \Rightarrow$  smaller  $\sigma^2 \Rightarrow$  more likely to overfit
- A heuristic in practice: choose  $\sigma = \text{median}(\|x_i - x_j\|, i \neq j)$

# How to recognize a valid kernel

- Two methods

(1) Find an explicit feature representation  $\phi(x)$

(2) check Mercer's condition

- (Def) Let  $K(x, x')$  be a kernel. Let  $S = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^d$   
Then,  $G := [K(x_i, x_j)]_{ij} \in \mathbb{R}^{n \times n}$  is called the **Gram matrix** of  $S$

cf. covariance matrix

- (Thm) Let  $K(x, x')$  be a symmetric function. Then,

$K(x, x') = \langle \phi(x), \phi(x') \rangle$  for some  $\phi(x) \iff K$  satisfies *Mercer's condition*: for any  $\forall m \geq 1, \forall x_1, \dots, x_m$ , the **Gram matrix** of  $S = \{x_1, \dots, x_m\}$  is PSD.

# Example

- Is  $K(x, y) = \max(x, y)$  a valid kernel?
- After some failed trials of constructing  $\phi$ , you may want to disprove that  $K$  is a kernel
- Suffices to show that  $K$  fails Mercer's condition, i.e. exists some dataset  $S$  whose Gram matrix is not PSD
- Guess  $S = \{-1\} \Rightarrow G = (-1)$  not PSD
  
- What if we restrict the inputs  $x, y \geq 0$ ?
- Guess  $S = \{0, 2\} \Rightarrow G = \begin{pmatrix} K(0,0) & K(0,2) \\ K(2,0) & K(2,2) \end{pmatrix} = \begin{pmatrix} 0 & 2 \\ 2 & 2 \end{pmatrix}$ 
  - $G$  is not PSD. Why?
  - Method 1: find  $v$  such that  $v^T G v < 0$
  - Method 2: check that some eigenvalues of  $G$  are  $< 0$

# Building Kernels from simpler ones

- Kernels are closed under
  - Positive scaling
  - sum/product
  - composition with a positive power series:  $\sum_{i=1}^{\infty} a_i (K(x, x'))^i$ , where  $a_i \geq 0$  for all  $i$

# Building Kernels from simpler ones (cont'd)

## kernel composition

$$k(\mathbf{x}, \mathbf{v}) = k_a(\mathbf{x}, \mathbf{v}) + k_b(\mathbf{x}, \mathbf{v})$$

$$k(\mathbf{x}, \mathbf{v}) = \gamma k_a(\mathbf{x}, \mathbf{v}), \quad \gamma > 0$$

$$k(\mathbf{x}, \mathbf{v}) = k_a(\mathbf{x}, \mathbf{v})k_b(\mathbf{x}, \mathbf{v})$$

$$k(\mathbf{x}, \mathbf{v}) = \mathbf{x}^\top A \mathbf{v}, \quad A \text{ is p.s.d.}$$

$$k(\mathbf{x}, \mathbf{v}) = f(\mathbf{x})f(\mathbf{v})k_a(\mathbf{x}, \mathbf{v})$$

## mapping composition

$$\phi(\mathbf{x}) = (\phi_a(\mathbf{x}), \phi_b(\mathbf{x}))$$

$$\phi(\mathbf{x}) = \sqrt{\gamma} \phi_a(\mathbf{x})$$

$$\phi_l(\mathbf{x}) = \phi_{ai}(\mathbf{x})\phi_{bj}(\mathbf{x})$$

$$\phi(\mathbf{x}) = L^\top \mathbf{x}, \quad \text{where } A = LL^\top$$

$$\phi(\mathbf{x}) = f(\mathbf{x})\phi_a(\mathbf{x})$$

# Kernelized Perceptron algorithm

- How to combine the Perceptron algorithm with a nonlinear feature mapping  $\phi: \mathcal{X} \rightarrow \mathbb{R}^D$ ?
- Recall the Perceptron algorithm:

---

**Algorithm 29** PERCEPTRONTRAIN( $\mathbf{D}$ , *MaxIter*)

---

```
1:  $w \leftarrow 0, b \leftarrow 0$  // initialize weights and bias
2: for iter = 1 ... MaxIter do
3:   for all  $(x, y) \in \mathbf{D}$  do
4:      $a \leftarrow w \cdot \phi(x) + b$  // compute activation for this example
5:     if  $ya \leq 0$  then
6:        $w \leftarrow w + y \phi(x)$  // update weights
7:        $b \leftarrow b + y$  // update bias
8:     end if
9:   end for
10: end for
11: return  $w, b$ 
```

---

- Suppose  $\phi$  is associated with a kernel  $K$
- Is it possible to implement this without ever explicitly computing  $\phi$ ?



# Kernelized Perceptron algorithm

- Key observation: throughout the run of the Perceptron algorithm,  $w$  always lies in  $\text{span}(\phi(x_1), \dots, \phi(x_n))$ , i.e.

$w$  always has the form  $\alpha_1 \phi(x_1) + \dots + \alpha_n \phi(x_n)$

- Key algorithmic idea: instead of maintaining  $w \in \mathbb{R}^D$ , we maintain its linear combination coefficient  $(\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$ !

---

**Algorithm 30** KERNELIZEDPERCEPTRONTRAIN( $\mathbf{D}$ ,  $MaxIter$ )

---

```
1:  $\alpha \leftarrow \mathbf{0}$ ,  $b \leftarrow 0$  // initialize coefficients and bias
2: for  $iter = 1 \dots MaxIter$  do
3:   for all  $(x_n, y_n) \in \mathbf{D}$  do
4:      $a \leftarrow \sum_m \alpha_m \phi(x_m) \cdot \phi(x_n) + b$  // compute activation for this example
5:     if  $y_n a \leq 0$  then
6:        $\alpha_n \leftarrow \alpha_n + y_n$  // update coefficients
7:        $b \leftarrow b + y$  // update bias
8:     end if
9:   end for
10: end for
11: return  $\alpha, b$ 
```

---

# Kernelized ridge regression

- Recall ridge regression:  $\hat{w} = \arg \min_w \|Xw - y\|^2 + \lambda \|w\|^2$

$$\begin{pmatrix} -x_1 - \\ \dots \\ -x_n - \end{pmatrix} \cdot w \approx \begin{pmatrix} y_1 \\ \dots \\ y_n \end{pmatrix}$$

$X: n \times d$                        $y: n \times 1$

- Woodbury matrix identity (matrix inversion lemma)

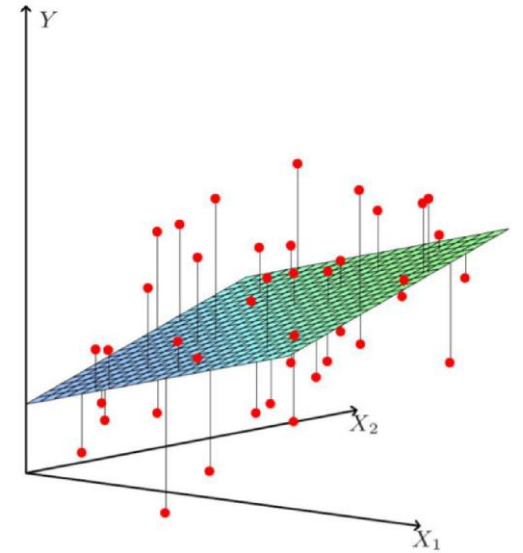
$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

- (Thm)  $\hat{w}$  can be alternatively written as  $\hat{w} = X^T(\lambda I_n + \boxed{X X^T})^{-1}y$

(proof)

Gram matrix

- starting point: recall  $\hat{w} = (\lambda I + \boxed{X^T X})^{-1}X^T y$  => apply the lemma above to  $(\lambda I + X^T X)^{-1}$   
(scaled) covariance matrix
- tip: when you get stuck, try the special case of  $d = 1$  to get a sense.



# Kernelized ridge regression

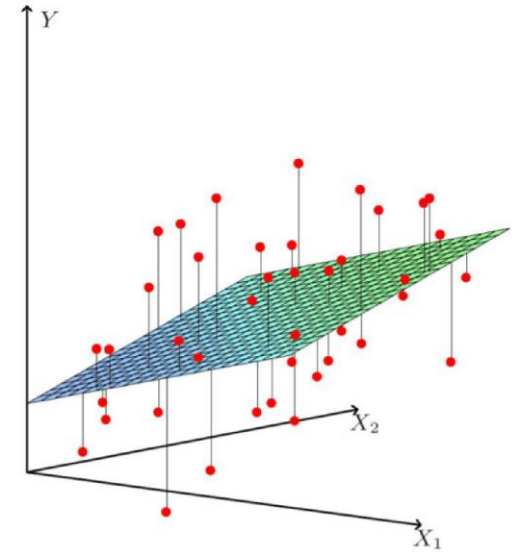
$$\begin{pmatrix} -\phi(x_1) - \\ \dots \\ -\phi(x_n) - \\ \Phi: n \times d \end{pmatrix} \cdot w \approx \begin{pmatrix} y_1 \\ \dots \\ y_n \\ y: n \times 1 \end{pmatrix}$$

- $\hat{w} = \Phi^\top (\lambda I + \Phi \Phi^\top)^{-1} y$

- Recall  $\Phi \Phi^\top = \left( K(x_i, x_j) \right)_{i,j}$  is the Gram matrix

- Prediction for  $x$ : 
$$\begin{aligned} \langle \hat{w}, \phi(x) \rangle &= \phi(x)^\top \hat{w} \\ &= \phi(x)^\top \Phi^\top (\lambda I + \Phi \Phi^\top)^{-1} y \\ &= (K(x, x_1), \dots, K(x, x_n)) \cdot \alpha, \text{ where } \alpha = (\lambda I + \Phi \Phi^\top)^{-1} y \end{aligned}$$

- Again, avoids explicit representation of  $\phi(x_i)$ 's



# Time complexity of (kernelized) Perceptron & ridge regression

- $\phi$ : feature map from  $d$  dimensions to  $p$  dimensions
- $n$ : training set size
- $k$ : the number of operations to evaluate  $K(x, x')$
  
- Test stage:

	Without kernel trick	With Kernel trick
Ridge regression	$O(p)$	$O(nk)$
Perceptron	$O(p)$	$O(nk)$

- Training stage:

	Without kernel trick	With Kernel Trick
Ridge regression	$O(np^2 + p^3)$	$O(n^2k + n^3)$
Perceptron	$O(\#iters \times p)$	$O(\#iters \times nk)$

# Next lecture (10/3)

- Unsupervised learning
- Assigned reading: CIML 3.4 (Review) 11.3