

# DSAS105 Lec2

Recall

- ① Linear Reg
- Linear Basis Regress

$$g: x \in \mathbb{R}^d \mapsto g(x) \in \mathbb{R}^p$$

- ② Regularization
  - $\ell_1$ -norm: LASSO
  - $\ell_2$ -norm: Ridge

- ③ Apply linear model to classification task

① change output activation function

Softmax

② change Loss function

CROSS-ENTROPY LOSS

Architecture

$$x \xrightarrow{g(\cdot)} \phi(x) \in \mathbb{R}^M \xrightarrow{W \in \mathbb{R}^{k \times M}} W \cdot \phi(x)$$

Softmax

probability

★ Equivalent

can view as  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^M$  + projection to  $\mathbb{R}^k$

Feature Maps

and  $\phi_j: x \mapsto \phi_j(x) \in \mathbb{R}$   
 $w_j \in \mathbb{R}^k$

→ For Reg, we use  $\mathcal{H}_1 = \{f: f(x) = \sum w_j \phi_j(x)\}$

→ For Classification, we use  $\mathcal{H}_2 = \{f: f(x) = g(\sum_{j=1}^M w_j \phi_j(x))\}$

$g(\cdot) \Rightarrow$  Softmax function to generate a prob. dist.

output of the model is a K-class probability

$$\text{in general, } g(z_k) = \frac{h(z_k)}{\sum h(z_k)} \quad h > 0$$

For classification task

Model Training

① Loss function:  $L(y'; y) = - \sum_{k=1}^K y_k \log y'_k$

empirical risk minimization

$$\text{ERM} \rightarrow \min_{W \in \mathbb{R}^{M \times K}} R_{\text{emp}}(W) = \min_W \frac{1}{N} \sum_{i=1}^N L(g(W^T \Phi)_i, y_i)$$

Notation

$$\mathbb{R}^{M \times N} \ni \Phi = \begin{pmatrix} \phi_1(x_1), \dots, \phi_1(x_N) \\ \vdots \\ \phi_M(x_1), \dots, \phi_M(x_N) \end{pmatrix}$$

$$W = \begin{pmatrix} w_1^T \\ \vdots \\ w_M^T \end{pmatrix} \in \mathbb{R}^{M \times K}$$

the probability prediction for

## ② Training Acc vs Testing Acc

$i$ -th data

$$L(g(W^T \phi(x_i)), y_i)$$

Model training objective is based on training set

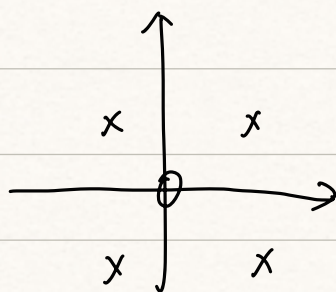
min: objective (training set)

§§ auxiliary?

make all prediction of training sample correct

Today

### ① Importance of feature Map

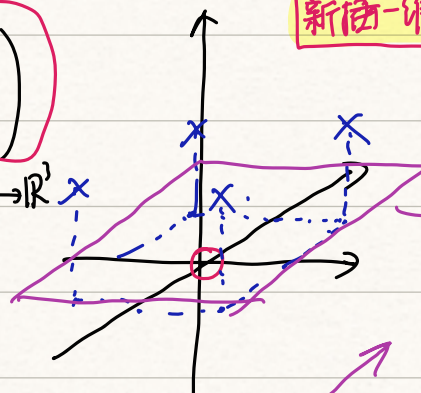


$$\phi_1 = \begin{pmatrix} x_1 \\ x_2 \\ \|x\| \end{pmatrix}$$

$$\phi_1(\cdot): \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

works for Radial-shaped data

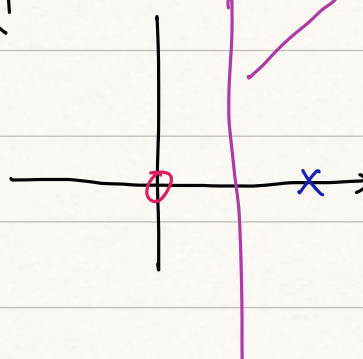
新增一维分层 (idea)



hyperplane to separate 2 classes (2 clusters)

$$\phi_2(\cdot): \mathbb{R}^2 \rightarrow \mathbb{R}^1$$

$$\phi_2 = \|\cdot\|_2$$



② feature map  $\rightsquigarrow$  similarity measure  $\leftarrow$  **KERNEL**

$$\langle \phi_i(x), \phi_j(x) \rangle$$

### ③ **Ridge Regression**

$$\text{Model: } f(x) = \sum_{j=0}^{M-1} w_j \phi_j(x)$$

$$\phi_j: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^M$$



$$\hat{w} = (\phi^T \phi + \lambda I_M)^{-1} \phi^T y$$

$O(N)$  Computational Cost

$$y = \Phi w$$

$$\Phi = \begin{pmatrix} \phi^T(x_1) \\ \vdots \\ \phi^T(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}$$

→ Punchline  $\Rightarrow$  Reformulation of Ridge Reg

$$\hat{w} = (\Phi^T \Phi + \lambda I_M)^{-1} \Phi^T y \leftarrow \text{solution of Ridge Reg}$$

$$\text{Model: } \hat{f}(x) = w^T \phi(x)$$

Let's show

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

$$\Leftrightarrow \Phi^T (\Phi \Phi^T + \lambda I_N)^{-1} y = (\Phi^T \Phi + \lambda I_M)^{-1} \Phi^T y$$

$$\Leftrightarrow (\Phi^T \Phi + \lambda I_M) \Phi^T (\Phi \Phi^T + \lambda I_N)^{-1} y = \Phi^T y$$

$$\longrightarrow \Phi^T \Phi \Phi^T + \lambda I_M \Phi^T$$

$$= \Phi^T \Phi \Phi^T + \lambda \Phi^T I_N$$

$$= \Phi^T (\Phi \Phi^T + \lambda I_N)$$

$$\Leftrightarrow \Phi^T y = \Phi^T y \quad (\checkmark)$$

$$\Phi^T \in \mathbb{R}^{M \times N}$$

$$\Phi = \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M}$$

$$\text{Thus } \hat{w} = \Phi^T (\Phi \Phi^T + \lambda I_N)^{-1} y$$

$$\rightarrow \hat{f}(x) = \hat{w}^T \phi(x)$$

$$= \phi(x)^T \hat{w}$$

$$= \phi(x)^T \Phi^T (\Phi \Phi^T + \lambda I_N)^{-1} y$$

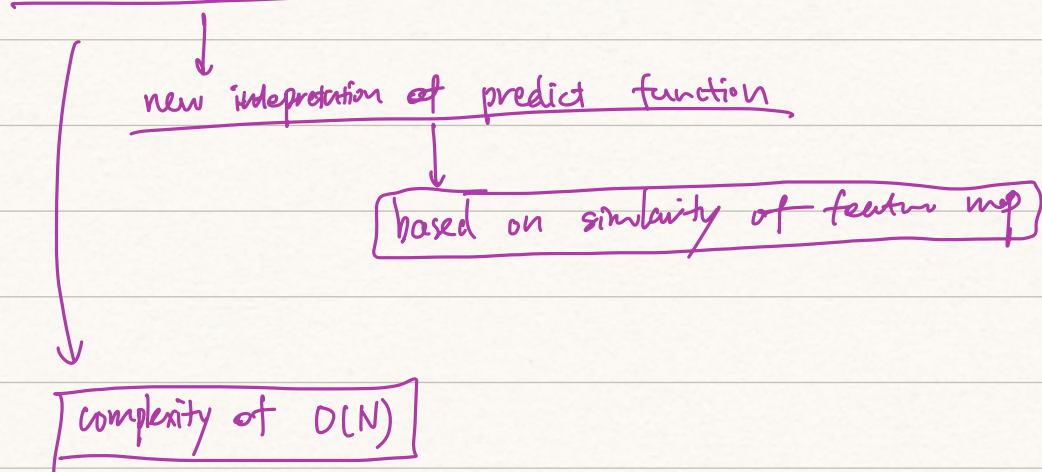
$$= (\Phi \phi(x))^T \alpha$$

$$\text{where } \alpha = (\Phi \Phi^T + \lambda I_N)^{-1} y$$

$$= \left[ \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{pmatrix} \cdot \phi(x) \right]^T \alpha \quad \boxed{\alpha \in \mathbb{R}^N}$$

$$= \begin{bmatrix} \phi(x_1)^T \phi(x) \\ \vdots \\ \phi(x_N)^T \phi(x) \end{bmatrix}^T \alpha$$

$$= \sum_{i=1}^N \alpha_i \langle \phi(x_i), \phi(x) \rangle \Rightarrow \boxed{\text{Dual Form of SVM}}$$



recall: standard form of Ridge Reg  $\rightarrow O(M)$

$$\hat{f}(x) = \sum_{j=0}^{M-1} \hat{w}_j \phi_j(x) \rightarrow \begin{cases} \text{no need of data } x_1, \dots, x_N \\ O(M) \end{cases}$$

$$\hat{f}(x) = \sum_{i=1}^N \alpha_i \langle \phi(x_i), \phi(x) \rangle \begin{cases} \text{need } x_1, \dots, x_N \text{ to compute similarity} \\ O(N) \end{cases}$$

$$\alpha = (\Phi \Phi^T + \lambda I_N)^{-1} y \quad \leftarrow \text{depends on feature map}$$

$$\Phi = \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{pmatrix} \in \mathbb{R}^{N \times M} \quad \text{ONLY THROUGH } \boxed{\phi(x_j)^T \phi(x)}$$

$$\text{and } \Phi \Phi^T := G$$

$$G_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$$

define  $K(x, y)$  directly

$$\hat{f}(x) = \sum \alpha_i K(x_i, x) \\ \alpha = (G + \lambda I_N)^{-1} y$$

$\Rightarrow$  can define directly

$K(x_j, x) \rightsquigarrow$  kernel function

$\downarrow$   
similarity measure



$$\hookrightarrow G_{ij} = k(x_i, x_j)$$

How to choose a valid kernel?  $k(\cdot, \cdot)$

① Necessary Condition

1. symmetry

~~2. Non-negative  $k(x, x) \geq 0$~~

3. PSD for arbitrary choice of  $(x_1, \dots, x_n)$

★

② Mercer's THM tells us this is also Sufficient:

if  $k(\cdot, \cdot)$  is SPD, then there exists  $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$

$$k(x, y) = \langle \phi(x), \phi(y) \rangle$$

IDEA: if we want to do Linear BASIS Ridge Reg.

we may not need to define feature map  $\phi(\cdot)$ ,

instead, we select a VALID kernel  $k(\cdot, \cdot)$

and do:  $\hat{f}(x) = \sum \alpha_i k(x, x_i)$

$$\alpha = (\Phi \Phi^T + \lambda I_N)^{-1} y \quad !!!$$

Kernels & Feature Maps

① polynomial kernel

$$k(x, x') = (1 + x^T x')^2$$

$$\phi(x) = \begin{pmatrix} 1 \\ \sqrt{2}x \\ x^2 \end{pmatrix}$$

② RBF kernel / Gaussian kernel

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{\alpha}\right) \quad x, y \in \mathbb{R}^d$$

Taylor:

$$\exp(z) = \sum_{k=0}^{\infty} \frac{1}{k!} z^k$$



$$\exp\left(\frac{2xy}{\alpha}\right)$$

$$= \exp\left(-\frac{x^2}{\alpha}\right) \exp\left(-\frac{y^2}{\alpha}\right) \exp\left(\frac{2xy}{\alpha}\right)$$

$$= \exp\left(-\frac{x^2}{\alpha}\right) \exp\left(-\frac{y^2}{\alpha}\right) \sum_{k=0}^{\infty} \frac{2^k}{\alpha^k k!} (x)^k (y)^k$$

$$= \sum_{k=0}^{\infty} \frac{2^k}{\alpha^k k!} \left( \exp\left(-\frac{x^2}{\alpha}\right) x^k \right) \left( \exp\left(-\frac{y^2}{\alpha}\right) y^k \right)$$

$$\sum_{k=0}^{\infty} \frac{1}{k!} \left( \frac{2xy}{\alpha} \right)^k$$

$$= \sum_{k=0}^{\infty} \frac{2^k}{k! \alpha^k} (x)^k (y)^k$$

$$= \phi(x)^T \phi(y)$$

$$\phi(x) = \begin{pmatrix} \phi_0(x) \\ \vdots \\ \phi_n(x) \end{pmatrix}$$

where  $\phi_i(x) = \sqrt{\frac{2^i}{\alpha^i i!}} \exp\left(-\frac{x^2}{\alpha}\right) \cdot x^i$

Gaussian feature map

$$\mathcal{F}_{\text{GOF}}: \mathbb{R}^{d'} \rightarrow \mathbb{R}^{\infty}$$

## Kernel Recap

① Why we need kernel?

→ Model (Kernel Ridge Regression)

$$\hat{f}(x) = \hat{w}^T \phi(x)$$

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

$$= \Phi^T (\Phi \Phi^T + \lambda I_N)^{-1} y$$

$$= \Phi^T (G + \lambda I_N)^{-1} y$$

$$\hat{f}(x) = (\Phi^T (G + \lambda I_N)^{-1} y)^T \phi(x)$$

$$= \phi(x)^T \Phi^T \cdot (G + \lambda I_N)^{-1} y$$

$$= \begin{pmatrix} \phi(x_1)^T \phi(x) \\ \vdots \\ \phi(x_n)^T \phi(x) \end{pmatrix}^T \cdot (G + \lambda I_N)^{-1} y$$

$$= \sum_{i=1}^N \alpha_i \phi(x_i)^T \phi(x)$$

$$\underline{\alpha} = (G + \lambda I_N)^{-1} y$$



$$:= \sum_{i=1}^n \alpha_i k(x_i, x)$$

② What can construct a valid kernel  $k(x, y)$  ?

[E.g.]  $k(x, y) = (1 + x^T y)^k$  if  $x, y \in \mathbb{R}^d$

$$(1 + x_1 y_1 + \dots + x_d y_d)^k = \sum_{k_0 + k_1 + \dots + k_d = k, k_i \geq 0} \binom{k}{k_0, \dots, k_d} \prod_{i=1}^d (x_i y_i)^{k_i}$$

$$= \begin{pmatrix} \frac{1}{\sqrt{k}} x_1 \\ \vdots \\ \frac{1}{\sqrt{k}} x_d \\ \vdots \\ \sqrt{\binom{k_0 + \dots + k_d}{k}} \prod_{i=1}^d x_i^{k_i} \\ \vdots \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{k}} y_1 \\ \vdots \\ \frac{1}{\sqrt{k}} y_d \\ \vdots \\ \sqrt{\binom{k_0 + \dots + k_d}{k}} \prod_{i=1}^d y_i^{k_i} \\ \vdots \end{pmatrix} \left\{ \begin{pmatrix} d \\ k+d \end{pmatrix} \right\}$$

$$\Rightarrow \phi(x) = \begin{pmatrix} \vdots \\ \sqrt{\binom{k_0 + \dots + k_d}{k}} \prod_{i=1}^d x_i^{k_i} \\ \vdots \end{pmatrix} \in \mathbb{R}^{\binom{d}{k+d}}$$

★

Mercer's theorem

$k(x, y)$   $\rightarrow$  SPD kernel  $\iff$  the existence of feature map s.t.  $K(x, y) = \langle \phi(x), \phi(y) \rangle$

Defn of SPD kernel  $k(x, y) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

$\Rightarrow$  consider  $\begin{cases} k(x, y) \\ \{x_1, \dots, x_n\} \end{cases} \Rightarrow$  matrix  $K \in \mathbb{R}^{n \times n}$  and  $K_{ij} = \langle \phi(x_i), \phi(x_j) \rangle$

$[\forall n \text{ and } \forall x_1, \dots, x_n]$

satisfy  $\begin{cases} ① \text{ symmetric} \\ ② \text{ PSD} \end{cases}$

Kernel trick

③ Avoid introducing feature map, we can define a SPD kernel  $k(x, y)$  to apply Kernel Ridge Regression Method and so on.

Question: How can we verify  $k(x, y)$  is SPD ?

④ Verify SPD Kernel

a) through definition: [for simple SPD kernel]

→ for  $\forall n$  &  $\forall x_1, \dots, x_n$ , consider  $K_{ij} = k(x_i, x_j)$

check  $K$  is  $\begin{cases} \text{symmetric} \\ \text{PSD} \end{cases}$

(RBF kernel)

b)  $\star$  use SPD Kernel Closure Property: [for complicate SPD kernel]

- ① scaling property:  $K(x, y) = \lambda k_1(x, y)$
- ② addition property:  $K(x, y) = k_1(x, y) + k_2(x, y)$
- ③ normalization property:  $K(x, y) = g(x) k_1(x, y) g(y)$
- ④ limit property:  $K(x, y) = \lim_{n \rightarrow \infty} k_n(x, y)$
- ⑤ product property:  $K(x, y) = k_1(x, y) k_2(x, y)$