

MA4270. Lecture 7.

Last time: Kernel.

(Affinely)

$\underline{x} \in \mathbb{R}^d$: Original feature vector / may not be linearly separable in \mathbb{R}^d .

One Way is SVM with slackness
 Other way \rightarrow Kernel. \downarrow Hinge loss

$d=1$ Map $x \in \mathbb{R}^d \xrightarrow{\phi}$ feature map $\underline{\phi}(x)$: e.g. polynomial kernel.
 $\left\{ \begin{array}{l} \underline{\phi}(x) = [1, \sqrt{2}x, x^2] \rightarrow \text{quadratic} \\ \underline{\phi}(x) = [1, \sqrt[3]{x}, \sqrt[3]{x^2}, x^3] \rightarrow \text{cubic} \end{array} \right.$

yield non-linear regression function, e.g.

$$y = \theta_0 + \theta_1 \sqrt{4}x + \theta_2 \sqrt{6}x^2 + \theta_3 \sqrt{4}x^3 + \theta_4 x^4.$$

When $d=2$. $\underline{x} = (x_1, x_2)^T \in \mathbb{R}^2 \longrightarrow \underline{\phi}(x) = \begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}$

$$\begin{aligned} \langle \underline{\phi}(x), \underline{\phi}(z) \rangle &= \begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix}^T \begin{bmatrix} 1 \\ \sqrt{2}z_1 \\ \sqrt{2}z_2 \\ \sqrt{2}x_1 z_2 \\ z_1^2 \\ z_2^2 \end{bmatrix} \\ &\uparrow \quad \text{High dimension } \in \mathbb{R}^6 \end{aligned}$$

$$\begin{aligned} &(1 + x_1 z_1 + x_2 z_2)^3 \\ &\downarrow \quad \phi(x) = \begin{pmatrix} 1 \\ \sqrt{3}x_1 x_2 \\ \sqrt{3}x_1^2 x_1 \\ \sqrt{3}x_1^2 x_2 \\ x_1^3 \\ x_2^3 \end{pmatrix} \quad \begin{pmatrix} 1 \\ \sqrt{3}x_1 \\ \sqrt{3}x_1 \\ \sqrt{3}x_2 \\ \sqrt{3}x_2 \\ \sqrt{3}x_1^2 \end{pmatrix} \end{aligned}$$

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

$$= (1 + x_1 z_1 + x_2 z_2)^2 = (1 + \langle \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \rangle)^2.$$

$$= (1 + \langle \underline{x}, \underline{z} \rangle)^2 \xrightarrow{P} \text{takes } O(d+p) \text{ time.}$$

\xrightarrow{d}

Why do we consider Inner Product?

→ Consider a kernelized linear regression model. \Rightarrow Just need the calculation of Inner product among $\{\Phi(\mathbf{x}_t)\}_{t=1}^n \subseteq \mathbb{R}^D$.

$$y = \underline{\theta}^T \Phi(\mathbf{x}) + \varepsilon$$

feature map

Goal: Show that estimation of $\hat{\underline{\theta}}$ using MLE [ONLY] involves the inner product among feature vectors $\Phi(\mathbf{x}_t)$

→ Regularized Least Square Prob.

$$J(\underline{\theta}) = \sum_{t=1}^n (y_t - \underline{\theta}^T \Phi(\mathbf{x}_t))^2 + \lambda \|\underline{\theta}\|_2^2$$

Penalized log-likelihood Criterion

important -

$$\underline{\alpha} = \lambda (\lambda I + K)^{-1} \mathbf{y}$$

$$\hat{\underline{\theta}} = \frac{1}{\lambda} \sum_{t=1}^n \alpha_t \Phi(\mathbf{x}_t)$$

prediction

$$\mathbf{y}' = \hat{\underline{\theta}}^T \Phi(\mathbf{x}')$$

$$= \left[\frac{1}{\lambda} \sum_{t=1}^n \alpha_t \Phi(\mathbf{x}_t) \right]^T \Phi(\mathbf{x}')$$

$$= \frac{1}{\lambda} \sum_{t=1}^n \alpha_t \langle \Phi(\mathbf{x}_t), \Phi(\mathbf{x}') \rangle$$

$$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$$

$$D \gg d$$

$$\Rightarrow \hat{\underline{\theta}} = \frac{1}{\lambda} \sum_{t=1}^n \alpha_t \cdot \Phi(\mathbf{x}_t) \in \mathbb{R}^D \rightsquigarrow \text{high dimension}$$



In order to determine $\hat{\underline{\theta}}$ in \mathbb{R}^d , we only have

to specify n numbers α_t $t=1, 2, \dots, n$



$$\hat{\underline{\theta}} \in \text{span} \{ \Phi(\mathbf{x}_t) \}_{t=1}^n$$

$\hat{\underline{\theta}}$ is determined only by $\{\mathbf{x}_t\}_{t=1}^n$ instead of $\dim(\Phi) \in \mathbb{R}^D$

($D \gg d, D \gg n$)

① Q: How to find α_t ? \Rightarrow Answer { kernel matrix } y .

$$\alpha_t = y_t - \underbrace{\Phi^T \Phi}_{\text{prediction difference}} (\underline{\alpha}) = y_t - \left[\frac{1}{n} \sum_{s=1}^n \alpha_s \langle \Phi(x_s), \Phi(x_t) \rangle \right]^T \Phi(x_t)$$

$$= y_t - \frac{1}{n} \sum_{s=1}^n \alpha_s \langle \Phi(x_s), \Phi(x_t) \rangle$$

Note: α_t : depends on y_t as well as the inner product

Def [Kernel matrix] - Given $\{x_t\}_{t=1}^n \subseteq \mathbb{R}^d$ $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^P$

$$K = \begin{bmatrix} \langle \Phi(x_1), \Phi(x_1) \rangle & \cdots & \langle \Phi(x_1), \Phi(x_n) \rangle \\ \langle \Phi(x_2), \Phi(x_1) \rangle & \ddots & \vdots \\ \vdots & \ddots & \vdots \\ \langle \Phi(x_n), \Phi(x_1) \rangle & \cdots & \langle \Phi(x_n), \Phi(x_n) \rangle \end{bmatrix}_{n \times n}$$

\Downarrow
 $\Phi^T \Phi$

$$[K]_{st} = \langle \Phi(x_s), \Phi(x_t) \rangle$$

Ex: Prove that K is positive semi-definite. \rightarrow definitive.

$$\text{Now, } \alpha_t = y_t - \frac{1}{n} \sum_{s=1}^n \alpha_s \langle \Phi(x_t), \Phi(x_s) \rangle$$

$$= y_t - \frac{1}{n} \sum_{s=1}^n \alpha_s K_{st} = \boxed{K_{ts}}$$

$$\underline{\alpha} = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} \quad \underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

then we have:

$$\underline{\alpha} = \underline{y} - \frac{1}{n} K \underline{\alpha} \Rightarrow (I + \frac{1}{n} K) \underline{\alpha} = \underline{y}$$

positive definite

$$\Rightarrow \hat{\alpha} = (\underbrace{I + \frac{1}{\lambda} K}_{\text{Invert}})^{-1} \underline{y}$$

$$= \lambda (N I + K)^{-1} \underline{y}$$

Punchline:

To solve $\underline{\alpha}$, we just need

$\left\{ \begin{array}{l} \underline{y} \rightarrow \text{target value} \\ K \rightarrow \text{kernel matrix} \end{array} \right.$

\Downarrow
only involves Inner Product

among $\{\Phi(\underline{x}_t)\}_{t=1}^n$

②.

After finding the $\{\alpha_t\}_{t=1}^n$, how do we predict a target given a test sample \underline{x}' ?

$$y = \hat{\theta}^\top \Phi(\underline{x}') = \left(\frac{1}{\lambda} \sum_{t=1}^n \alpha_t \Phi(\underline{x}_t) \right)^\top \Phi(\underline{x}')$$

$\hat{\theta}^\top$

$$= \frac{1}{\lambda} \sum_{t=1}^n \alpha_t \langle \Phi(\underline{x}_t), \Phi(\underline{x}') \rangle$$

prediction

$$= \frac{1}{\lambda} \sum_{t=1}^n \alpha_t K(\underline{x}_t, \underline{x}')$$

Calculate $K(\underline{x}_t, \underline{x}')$ for all $t \in \{1, 2, \dots, n\}$

$$= \langle \Phi(\underline{x}'), \Phi(\underline{x}_t) \rangle$$

Conclusion: { finding $\{\alpha_t\}_{t=1}^n$

prediction given sample \underline{x}'

\Rightarrow Both need $\langle \Phi(\cdot), \Phi(\cdot) \rangle$

Def: A (valid) kernel function $K: \underbrace{\mathbb{R}^d \times \mathbb{R}^d}_{\text{can be abstract!}} \rightarrow \mathbb{R}$.

is a function st \exists a feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$
s.t. $\forall x, x' \in \mathbb{R}^d$.

$$K(x, x') = \langle \phi(x), \phi(x') \rangle$$

if this value is large, then x & x' are similar!

Actually we don't care this!

[Rmk]: While we have written the original and kernelized feature space as \mathbb{R}^d & \mathbb{R}^D , neither of the spaces need to be finite dimension.

Prop: Let K_1, K_2 be valid kernel functions. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$.

be an arbitrary func. Then the following are valid kernel func.

- i) $K(x, x') = f(x) K_1(x, x') f(x')$
- ii) $K(x, x') = K_1(x, x') + K_2(x, x')$
- iii) $K(x, x') = K_1(x, x') K_2(x, x')$

[Intuition]

$$\left\{ \begin{array}{l} \phi^{(1)}: \mathbb{R}^d \rightarrow \mathbb{R}^{D_1} \\ \phi^{(2)}: \mathbb{R}^d \rightarrow \mathbb{R}^{D_2} \\ \langle \phi^{(1)}(x), \phi^{(1)}(x') \rangle = \langle \phi^{(2)}(x), \phi^{(2)}(x') \rangle \\ \langle \phi(x), \phi(x') \rangle = \sum_{D_1} \phi^{(1)}(x) \cdot \phi^{(1)}(x') + \sum_{D_2} \phi^{(2)}(x) \cdot \phi^{(2)}(x') \end{array} \right.$$

Pf: Let $\phi^{(1)}, \phi^{(2)}$ be the feature map associated to K_1 & K_2
 $\Leftrightarrow \forall i=1,2, K_i(x, x') = \langle \phi^{(i)}(x), \phi^{(i)}(x') \rangle$

i). Let $\phi(x) = \underbrace{\phi^{(1)}(x) f(x)}_{\text{scaler}}: \mathbb{R}^d \rightarrow \mathbb{R}^D$

\Downarrow
scaler.

Goal: prove ϕ is the feature map of K

$$\langle \phi(x), \phi(x') \rangle = \langle \phi^{(1)}(x), \phi^{(1)}(x') \rangle \cdot f(x) \cdot f(x')$$

$$= f(x) < \phi^{(1)}(x), \phi^{(1)}(x') > f(x')$$

$$= f(x) K_1(x, x') f(x')$$

$$= K(x, x')$$

\Rightarrow K is a valid kernel function.

ii) Let $\phi(x) = \begin{bmatrix} \phi^{(1)}(x) \\ \phi^{(2)}(x) \end{bmatrix} : \mathbb{R}^d \rightarrow \mathbb{R}^{2d}$

$$\langle \phi(x), \phi(x') \rangle = \left\langle \begin{bmatrix} \phi^{(1)}(x) \\ \phi^{(2)}(x) \end{bmatrix}, \begin{bmatrix} \phi^{(1)}(x') \\ \phi^{(2)}(x') \end{bmatrix} \right\rangle$$

$$= \langle \phi^{(1)}(x), \phi^{(1)}(x') \rangle + \langle \phi^{(2)}(x), \phi^{(2)}(x') \rangle$$

$$= K_1(x, x') + K_2(x, x')$$

iii) $K(x, x') = K_1(x, x') K_2(x, x')$ is a (valid) kernel function

$$\phi(x) = \phi^{(1)}(x) \phi^{(2)}(x)^T : \mathbb{R}^d \rightarrow \mathbb{R}^{D_1 \times D_2}$$

$$\underbrace{\langle \phi(x), \phi(x') \rangle}_{\text{matrix inner product}} = \sum_{i,j} \phi_{ij}(x) \cdot \phi_{ij}(x')$$

$$= \sum_{ij} \phi_i^{(1)}(x) \phi_j^{(1)}(x') \phi_i^{(2)}(x) \phi_j^{(2)}(x')$$

$$= \left(\sum_i \phi_i^{(1)}(x) \phi_i^{(1)}(x') \right) \left(\sum_j \phi_j^{(2)}(x) \phi_j^{(2)}(x') \right)$$

$$= K_1(x, x') K_2(x, x')$$

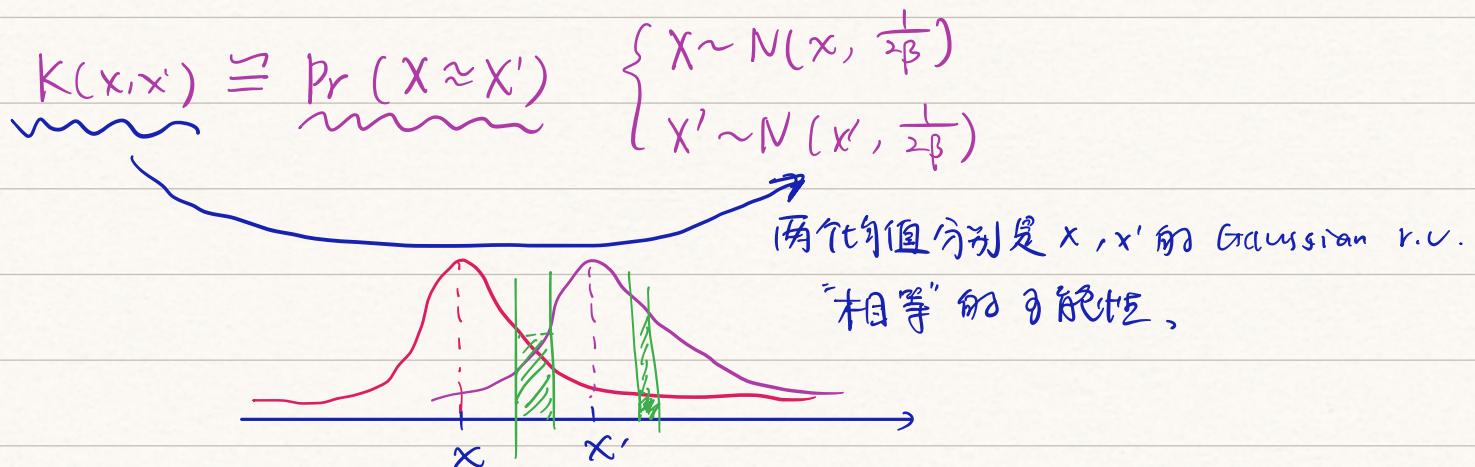
Claim: The function $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ given by:

$$K(x, x') = \exp\left(-\frac{\beta}{2}\|x-x'\|^2\right) \text{ is a kernel.}$$

Ex: Use the above Prop. to prove this!

Rmk: This is known as the Radial basis kernel.

Roughly Speaking,



$$= \lim_{\Delta \rightarrow 0} \int_{\mathbb{R}} \Pr(X \in [z, z+\Delta], X' \in [z, z+\Delta]) dz$$

Consider the functions $\phi_z(x) = C(\beta) N(z; x, \frac{1}{2\beta}) \quad z \in \mathbb{R}$

$$\begin{aligned} \phi : \mathbb{R}^d &\longrightarrow \text{Function Space} \\ x &\longmapsto \phi_z(x) \Rightarrow \text{a function of } z. \end{aligned}$$

$$\langle \phi(x), \phi(x') \rangle = \int_{\mathbb{R}} \phi_z(x) \phi_z(x') dz = \exp\left(-\frac{1}{2\beta}(x-x')^2\right).$$

$$\langle f, g \rangle = \int_{\mathbb{R}} f(t) g(t) dt$$

$$= \int_{\mathbb{R}} C(\beta)^2 \underbrace{N(z; x, \frac{1}{2\beta})}_{\Pr(X \approx z)} \underbrace{N(z; x', \frac{1}{2\beta})}_{\Pr(X' \approx z)} dz$$

$$\propto \int_{\mathbb{R}} \exp(-\beta(z-x)^2) \exp(-\beta(z-x')^2) dz$$

NEED TO
CHECK

$$= \exp(-\beta(x^2 + (x')^2)) \int_{\mathbb{R}} \exp(-\beta z^2 + 2\beta z(x+x') - \beta z^2) dz$$

$$= \exp(-\beta(x^2 + (x')^2)) \int_{\mathbb{R}} \exp(-\beta(2z^2 - 2z(x+x'))) dz$$

$$= \exp\left(-\frac{\beta}{2}(x-x')^2\right) \propto K(x, x')$$

The feature maps for the RBK $\exp\left(-\frac{\beta}{2}\|x-x'\|^2\right)$
are normal prob. density function $N(\cdot; x, \frac{1}{2\beta})$

Conclusion:

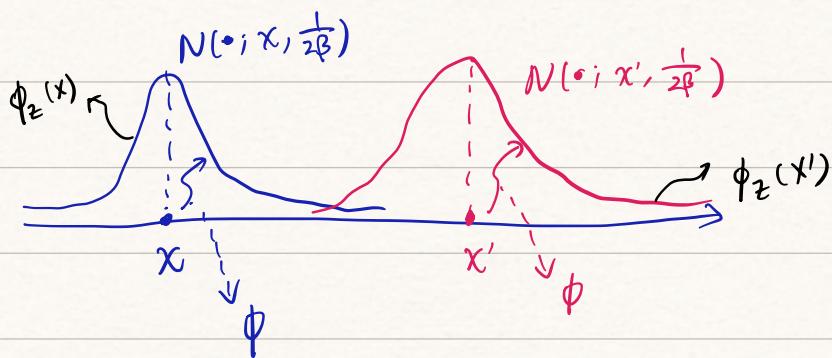
$$\text{Origin: } \langle \phi(x'), \phi(x) \rangle = \int_{\mathbb{Z}} \phi_z(x') \phi_z(x) dz$$

Now: we have $K(x, x') = \exp\left(-\frac{\beta}{2}\|x-x'\|^2\right)$

$\phi : \mathbb{R}^d \rightarrow \text{Function Space}$

$x \mapsto \{N(z; x, \frac{1}{2\beta}) : z \in \mathbb{R}\} \triangleq \{\phi_z(x)\}$

this is the variable of
function $\phi_z(x)$



String Kernels (Lodhi et al., JMLR 2012)

Training samples "car", "cat", "bat" & "bar"

Map each sample to a length-8 feature vector

Fix $\lambda \in (0, 1)$, weighting factor: account for compactness of substring in a string

e.g. "car" is present in "castard" and "card" but is more compact in "card"

Substrings $k=2$. (Feature map)

	ca	ct	at	ba	bt	cr	ar	br
$\phi(\text{cat})$	λ^2	λ^3	λ^2	0	0	0	0	0
$\phi(\text{car})$	λ^2	0	0	0	0	λ^3	λ^2	0
$\phi(\text{bat})$	0	0	λ^2	λ^2	λ^3	0	0	0
$\phi(\text{bar})$	0	0	0	λ^2	0	0	λ^2	λ^3

$$K(\text{car}, \text{cat}) = \langle \phi(\text{car}), \phi(\text{cat}) \rangle = \lambda^4$$

$$\hat{K}(\text{car}, \text{cat}) = \frac{K(\text{car}, \text{cat})}{\sqrt{K(\text{car}, \text{car})} \sqrt{K(\text{cat}, \text{cat})}} = \frac{\lambda^4}{2\lambda^4 + \lambda^6} = \frac{1}{2 + \lambda^2}$$

↓ also kernel

String Subsequence kernel (SSK)

Σ : finite alphabet $\{a, b, \dots, z\}$

s.t.: strings - consisting of letters from Σ

$|s|$ ($|t|$): length of strings

st : string formed by $s \& t$.

u : subseq of s if $\exists (i_1, \dots, i_{|u|}) = \underline{i}$

$1 \leq i_1 \leq i_2 \leq \dots \leq i_{|u|} \leq |s|$ s.t $u_j = s_{i_j} \forall j = 1, 2, \dots, |u|$

$u = s[\underline{i}]$

$\ell(\underline{i}) = i_{|u|} - i_1 + 1$

[E.g.:] $s = \text{segmentation}$ $u = \text{set}$

$\underline{i} = (1, 2, 7)$ $\ell(\underline{i}) = 7 - 1 + 1 = 7$

$$F_n = \mathbb{R}^{|\Sigma^n|}$$

$$\phi_u(s) = \sum_{\underline{i}: u = s[\underline{i}]} \lambda^{\ell(\underline{i})}$$

Cartesian product

[E.g.:] $s = \text{hono lulu}$ $u = \text{lu}$

$$\underbrace{\phi_u(s)}_{= 2\lambda^2}$$

$$K_n(s, t) = \sum_{u \in \Sigma^n} \langle \phi_u(s), \phi_u(t) \rangle$$

$$= \boxed{\sum_{u \in \Sigma^n}} \sum_{i: s[i] = u} \sum_{j: t[j] = u} \lambda^{\ell(i) + \ell(j)}$$



$\underbrace{O(|\Sigma|^n)}$

[Rmk].

This exponential dependence on n . can be reduced to linear

in n by using dynamic programming

Trick : $\underline{K_n = f(K_{n-1})}$