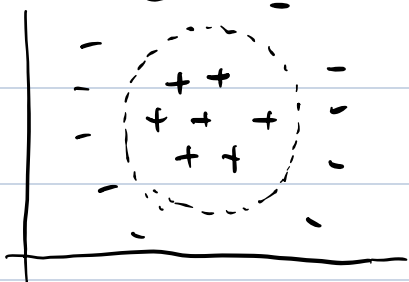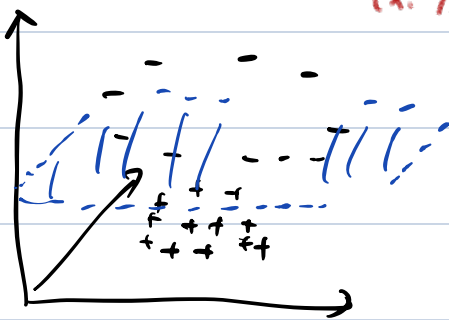# Reproducing Kernel Hilbert Space

§ Feature:

$(x, y)$ data

$(x, y, x^2+y^2)$ feature



linear      Tokyo $-$ Japan $+$ U.S $=$ D.C.

① Construct My feature by hand.

② Save all my feature. (If $\infty$ feature, the computation is hard)

$-$ Primal: size # data.                Dual: size # feature.

§ Revisiting Ridge Regression From Dual Side.

$$\min_w \quad \|Xw - Y\|_2^2 + \lambda \|w\|^2.$$

$\Rightarrow$ introduce $z = Xw$                    "Operator Splitting"

$$\min_w \quad \|z - Y\|_2^2 + \lambda \|w\|_2^2$$

$$s.t. \quad z = Xw$$

$$L(w, z, \alpha) = \|z - Y\|^2 + \lambda \|w\|_2^2 + \alpha^T(z - Xw)$$

$\underset{\text{Constant}}{}$

$- \quad \nabla_w L = 0 \qquad \underline{2\lambda w = X^T \alpha}$

$$\Rightarrow \quad w = X^T \alpha \quad (\Delta)$$

put $(\Delta)$ back to the objective.     $\min \|XX^T\alpha - Y\|_2^2 + \lambda \|X^T\alpha\|^2$

$$\Rightarrow \alpha = (XX^T + \lambda I)^{-1} Y$$

$$\min_w \quad \| Xw - y \|_2^2 + \lambda \| w \|^2$$

$\underline{\text{Dual}}: \quad w = X^T \alpha. \quad \text{where}. \quad \alpha = (XX^T + \lambda I)^{-1} y \quad$ #data

$\underline{\text{Primal}}: \quad w = (X^T \underset{\text{#fecture}}{X} + \lambda I)^{-1} X^T y \quad$ #fecture ... #data

§ Revisit of Polynomial Ression, $\mathbf{X}$

$$Y = \alpha_1 X^2 + \alpha_2 X + \alpha_3 = \underline{(X^2, X, 1)} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix}$$

↑ This is my new data

⇒ we can use linear regression on data $(X^2, X, 1)$ to do poly rgress

$$\begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} = (X^T X + \lambda I)^{-1} X^T y \quad \text{primal}$$

$$= X^T (XX^T + \lambda I)^{-1} y \quad \text{dual}$$

$$\begin{pmatrix} x_1^2 & x_1 & 1 \\ x_2^2 & x_2 & 1 \\ \vdots & & \\ x_n^2 & x_n & 1 \end{pmatrix} \begin{pmatrix} x_1^2 & x_2^2 & \cdots & x_n^2 \\ x_1 & x_2 & & x_n \\ 1 & 1 & & 1 \end{pmatrix}$$

$$(XX^T)_{i,j} = \langle (x_i^2, x_i, 1), (x_j^2, x_j, 1) \rangle$$

" inner product of fecture "

" How similar are the two features "

Kernel
$$\begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_1) & \cdots & \\ & & & \\ k(x_n, x_1) & k(x_n, x_1) & \cdots & k(x_n, x_n) \end{pmatrix}$$

how the data $x_i$ are similar to $x_j$

Idea. We can't compute the $w = X^T (XX^T + \lambda I)^{-1} y$

This can be compute

I still needs to know the feature.

if I have a test data $(x_{test})$

$\langle w, x_{test} \rangle = x_{test} X^T (XX^T + \lambda I)^{-1} y$

$$\left( k(x_{test}, x_1), k(x_{test}, x_2) \cdots k(x_{test}, x_n) \right)$$

- privacy: $x_1 \cdots x_n$ are used in testing
- $(XX^T + \lambda I)^{-1}$ is costly in large data sets.
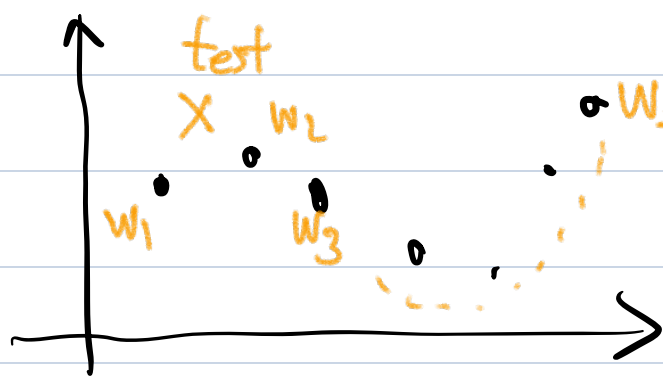- It is hard for online learning.

linear Behaviour of RKHS

$\langle w, x_{test} \rangle = x_{test} X^T \left( XX^T + \lambda I \right)^{-1} y$

a $\mathbb{R}^{\#data}$ vector
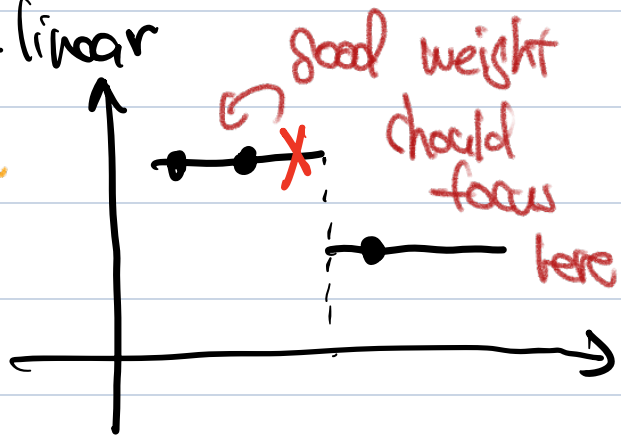only depend on $X$, but not depend on $y$

is a reweighting of $y$

Remark,

linear



non-linear

good weight
should
focus
here

non-parametric Estimator (?)

"Ridge Regression"

"Lasso"

① Why NN is better than Kernel         answer)
    ↳ adapt to jump    (one of the

---

① $W = X^T \alpha = \alpha_1 X_1 + \cdots + \alpha_n X_n$

What does data here mean.

$$f(x) = \langle W. x \rangle$$

X (data) actually means the mapping from
          function → function value.

Reproducing: $f \in H$, there always exists a (bounded)
linear mapping $k_x$, such that $\langle f, k_x \rangle = f(x)$

$f \in L_\infty$,    "RKHS is not a large space".

$\frac{2S}{d+2S}$           ↳ RKHS in $\mathbb{R}^d$, always have $d/2$ smooth

Hilbert Space: linear subspace with inner product.

$$k(x,y) = \langle k_x, k_y \rangle$$

$k = XX^T \Rightarrow k$ should be positive definite ✓

Question. can $k(x,y)$ be arbitrary?

Example.

① $k(x,y) = x^T y$

② $k_1$ is a kernel, $k_2$ is a kernel. then, $k_1 + k_2$ is a kernel

$\phi_1(x)$ is a feature, $\phi_2(x)$ is a feature. $[\phi_1(x), \phi_2(x)]$ is also feature

③ $k_1 \cdot k_2$ is also a kernel. $\phi_1 \cdot \phi_2$ is a feature

④ $\exp(k(x,y))$ is a kernel.

by ②, ③, $k(x,y)$ is a kernel, then

$f(k(x,y))$ is a kernel, if $f$ is a polynomial with positive coeff.

$$\exp(x) = \lim_{i \to \infty} (1 + x + \cdots + \frac{x^i}{i!})$$

⑤ $k(x,y) = \exp\left( -\frac{\|x-y\|^2}{\sigma^2} \right)$ is a kernel.

$$\exp\left(-\frac{\|x\|^2}{\sigma^2}\right) \exp\left(-\frac{\|y\|^2}{\sigma^2}\right) \quad \exp\left(\frac{2x \cdot y}{\sigma^2}\right)$$

⇓

rank-1 ← $k(u,v) = g(u)g(v)$ is a kernel

⑥ Translation-invariant kernel.

$$k(u,v) = f(u-v)$$

$$f(x) = \sum_{n=0}^{\infty} a_n \cos(nx)$$

$$f(u-v) = \underbrace{\sum_{n=0}^{\infty} a_n \cos\left(n(u-v)\right)}_{\sin(nu)\sin(nv) + \cos(nu)\cos(nv)} - (\Delta)$$

rank 1

rank $-1$

Eigen decompose $E$

$\Downarrow$

$\left\{ \sin(nu) : u \geq 1 \right\} \cup \left\{ \cos(nu) : u \geq 1 \right\}$ is the feature

of a translation invariant kernel.

[Random Fourier feature] random sample sin and cos

to construct feature map. is a good approximation.

to the translation-invariant kernel.

---

## Remark.

$(\Delta)$ is actually eigen decomposition.

$$A u_i = \lambda u_i$$

- If A is symmetric. $A = \underbrace{\sum_{i=1}^{n} \underbrace{\lambda u u^T}_{\text{rank-1 matrix}}}$

$k$ | #data

# data

$$\left[ \boxed{\ k\ } \right] \Big\| \text{#data} = \sum_{i=1}^{n} \lambda_i \left\| \begin{array}{c} u(x_1) \\ \ \\ u(x_n) \end{array} \right\| \text{#data} \quad \overbrace{\begin{array}{cc} u(x_1) & u(x_1) \\ \square & \square \end{array}}^{\text{#data}} \longrightarrow u(x_n)$$

#data

$\Downarrow$

$$\begin{vmatrix} u(x_1)u(x_1) & u(x_1)u(x_2) & \cdots \\ u(x_2)u(x_1) & u(x_1)u(x_1) & \cdots \\ \vdots & & \\ u(x_n)u(x_1) & u(x_n)u(x_1) & \cdots \end{vmatrix}$$

# Mercer's theorem.

$$k(x, y) = \sum_{i=1}^{n} \lambda_i \underline{\phi(x_i) \cdot \phi(x_j)}$$

$\Downarrow$

$$\Rightarrow \quad k(x, y) \, \phi(y) = \lambda \phi$$

$\Downarrow$

$$\int k(x, x') \, \phi(x') = \lambda \phi(x')$$