

Lecture 2 : Linear Regression.

Goal: $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ $x^{(i)} \mapsto y^{(i)}$
 \downarrow \downarrow
 \mathbb{R}^d \mathbb{R} $(x, y) \sim P(x, y)$
 (features, targets,
 samples) (labels)

Ex. 1 [Housing Price Prediction] (Regression)

x : $\begin{bmatrix} \leftarrow \\ \leftarrow \\ \leftarrow \\ \leftarrow \end{bmatrix}$ # bedrooms
bathrooms
sq. footage. y : $\begin{bmatrix} \leftarrow \end{bmatrix}$ Price of House

Ex. 2 Image Classification (Classification)

x : $\begin{bmatrix} \leftarrow \\ \leftarrow \\ \leftarrow \end{bmatrix}$ pixel location + color
 $[0, 255] \subset \mathbb{Z}$.

y : $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$ one-hot vector.
dogs
cats
birds

Q1: We want $\hat{f}(x) \approx y$ for all points $(x, y) \sim P_{(x,y)}$

predictor

What does \approx mean?

Q2: What class of functions

\hat{f} are we considering?

Q3: How do we find a good \hat{f} from this class?

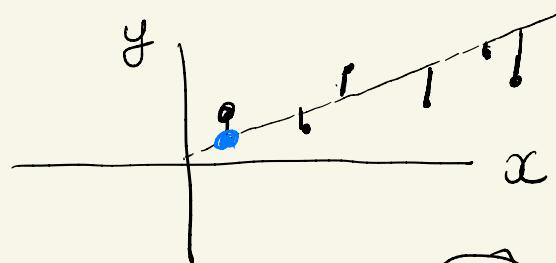
A1: Performance metric : Mean squared error (MSE) \leftarrow

$$\{(y_i, \hat{y}_i)\}_{i=1}^n$$

\uparrow
 R

\uparrow
 R

$$MSE := \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$



A2: $\{f: \mathbb{R}^d \rightarrow \mathbb{R}\}; f(x) = \underbrace{\omega}_{\mathbb{R}^{1 \times d}} \underbrace{x}_{\mathbb{R}^d}$

A3: $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$

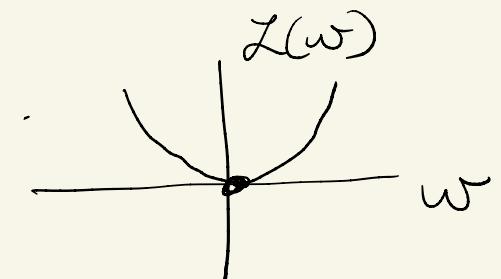
Empirical Risk Minimization

$$\text{minimize } \frac{1}{2} \sum_{i=1}^n (\omega x^{(i)} - y^{(i)})^2$$

Objective: minimize $\frac{1}{2} \sum_{i=1}^n (\omega x^{(i)} - y^{(i)})^2 \Leftrightarrow \frac{1}{2} \|\omega X - y\|_2^2$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(n)} \\ 1 & 1 & 1 & \dots & 1 \end{bmatrix}; \quad y = [y^{(1)} \dots y^{(n)}]$$

Ex. 1 $L(\omega) = \frac{1}{2} \omega^2$; Datapoint $x=1, y=0$.



$$L'(\omega) = \boxed{\omega = 0}.$$

Ex. 2 $L(\omega) = \frac{1}{2} \|\omega X - y\|_2^2$ ← convex.

$$\nabla L(\omega) = (\omega X - y) X^T = 0. \quad (1) \quad \begin{array}{l} XX^T \text{ invertible.} \\ n \gg d. \end{array}$$

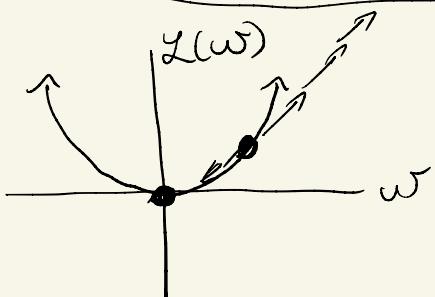
$$\omega \underset{\mathbb{R}^{1 \times d}}{\overset{\mathbb{R}^{d \times n}}{\underset{\mathbb{R}^{n \times d}}{\underset{\mathbb{R}^{1 \times n}}{XX^T}}} = y \underset{\mathbb{R}^{1 \times n}}{\overset{\mathbb{R}^{n \times d}}{\underset{\mathbb{R}^{1 \times d}}{X^T}}} \in \mathbb{R}^{n \times d}$$

$$\rightarrow \omega = \boxed{(y X^T) (X X^T)^{-1}}$$

$$(2) \quad \begin{array}{l} XX^T \text{ is not invertible.} \\ \underbrace{n < d.} \end{array}$$

Gradient Descent: $L(\omega)$; $\omega^{(0)}$ initialization, η step size.

$$\boxed{\omega^{(t+1)} = \omega^{(t)} - \eta \nabla L(\omega^{(t)})}$$



$$L(\omega) = \frac{1}{2} \omega^2 ; \quad \omega^{(0)}$$

$$\text{Step 1: } \underline{\omega^{(1)}} = \omega^{(0)} - \eta \omega^{(0)} = \underbrace{(1-\eta)}_C \underline{\omega^{(0)}}$$

$$\text{Step 2: } \underline{\omega^{(2)}} = (1-\eta) \underline{\omega^{(1)}} = \underbrace{(1-\eta)^2}_C \underline{\omega^{(0)}}$$

$$\boxed{C < 1} : \omega^{(t)} = \underbrace{C^t}_{\uparrow} \underline{\omega^{(0)}} \quad \overset{t \rightarrow \infty}{\omega^{(t)} \rightarrow \underline{0}}$$

$$C < 1 \Rightarrow 1 - \eta < 1$$

$$\boxed{0 < \eta < 2}$$

$$L(\omega) = \frac{1}{2} \|\omega X - y\|^2 ; \text{ Run g.d. step size } n, \omega^{(0)} = 0$$

$$\nabla L(\omega) = (\omega X - y) X^T$$

$$\underline{\text{Step 1: }} \omega^{(1)} = \omega^{(0)} - n (\omega^{(0)} X - y) X^T = \omega^{(0)} \underbrace{[\mathbb{I} - n X X^T]}_A + \underbrace{n y X^T}_B$$

$$\omega^{(1)} = \cancel{(\omega^{(0)} A)} + B$$

$$\underline{\text{Step 2: }} \omega^{(2)} = \omega^{(0)} A + B = (\omega^{(0)} A + B) A + B \\ = \omega^{(0)} A^2 + BA + B \\ = \cancel{\omega^{(0)} A^2} + B [A + \mathbb{I}]$$

$$\omega^{(1)} = B$$

$$\omega^{(2)} = B [A + \mathbb{I}]$$

$$\omega^{(3)} = B [A^2 + A + \mathbb{I}]$$

$$\rightarrow \omega^{(t)} = B [A^{t-1} + A^{t-2} + \dots + A + \mathbb{I}]$$

$$A := Q \Delta Q^T ; A^K = Q \Delta^K Q^T$$

$$\omega^{(t)} = B \left[\underline{Q} \Delta^{t-1} \underline{Q}^T + \underline{Q} \Delta^{t-2} \underline{Q}^T + \dots + \underline{Q} \Delta^1 \underline{Q}^T + \underline{Q} \mathbb{I} \underline{Q}^T \right] Q^T \\ = B Q \left[\underbrace{\Delta^{t-1} + \Delta^{t-2} + \dots + \Delta + \mathbb{I}} \right] Q^T$$

$$\omega^{(t)} = B Q [\Delta^{t-1} + \dots + \Delta + I] Q^T$$

$$A = I - n X X^T$$

$$= Q \Delta Q^T$$

$$B = n y X^T$$

$$= B Q \left[\begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}^{t-1} + \dots + \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} + \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \right] Q^T$$

$$= B Q \left[\frac{1 - \lambda_1^t}{1 - \lambda_1} \cdots \frac{1 - \lambda_n^t}{1 - \lambda_n} \right] Q^T$$

$$\boxed{1 + \lambda_1 + \lambda_1^2 + \dots + \lambda_1^{t-1} = \frac{1 - \lambda_1^t}{1 - \lambda_1}}$$

$$A = I - n X X^T ; X = U \Sigma V^T$$

$$= I - n U \underline{\Sigma^2} U^T ;$$

$$= \cancel{U \underline{I}} - n U \underline{\Sigma^2} U^T$$

$$= U [I - n \underline{\Sigma^2}] U^T$$

$$\omega^{(t)} = n y X^T Q \left[\frac{1 - (1 - n \sigma_1^2)^t}{n \sigma_1^2} \cdots \frac{1 - (1 - n \sigma_n^2)^t}{n \sigma_n^2} \right] Q^T$$

$$\boxed{\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}}$$

$$\lambda_i = 1 - n \sigma_i^2$$

$$\sigma_1 > \sigma_2 > \dots > \sigma_n$$

$$|1 - n \sigma_i^2| < 1$$

$$\boxed{n < \frac{2}{\sigma_1^2}}$$

$$t \rightarrow \infty w^{(t)} = \underbrace{y x^T u \left[\frac{1}{\sigma_1^2} \dots \frac{1}{\sigma_n^2} \right] u^T}_{A}$$

$$A = I - n X X^T ; X = \underline{u \Sigma v^T}$$

$$= \cancel{u} [I - n \cancel{\Sigma^2}] \cancel{u^T}$$

$\stackrel{=}{\circlearrowleft} Q \Delta Q^T$

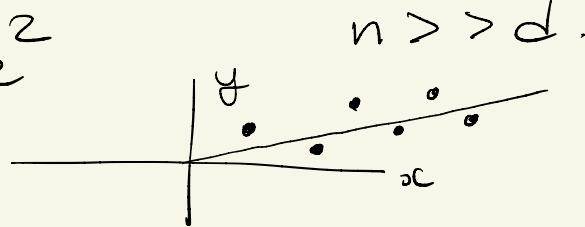
$$w^{(\infty)} = \cancel{y v \cancel{\Sigma^T} u^T} u \left[\frac{1}{\sigma_1^2} \dots \frac{1}{\sigma_n^2} \right] u^T$$

$$= y v \underbrace{\begin{bmatrix} \sigma_1 & \dots \\ & \ddots & \sigma_n \end{bmatrix}}_{\Sigma} \underbrace{\left[\frac{1}{\sigma_1^2} \dots \frac{1}{\sigma_n^2} \right]}_{\Sigma^{+2}} u^T$$

$$= \boxed{y v \left[\frac{1}{\sigma_1} \dots \frac{1}{\sigma_n} \right] u^T} = \boxed{y v \Sigma^+ u^T} = \boxed{y x^T}$$

(1) Underparameterized Regime

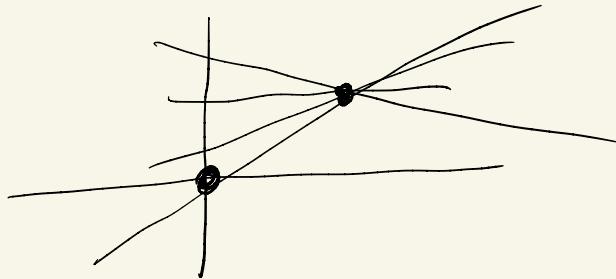
$$\frac{1}{2} \|w^T x - y\|_2^2$$



$$(X X^T)$$

$$w = (y X^T) (X X^T)^{-1}$$
$$= y X^+$$

(2) Overparameterized Regime

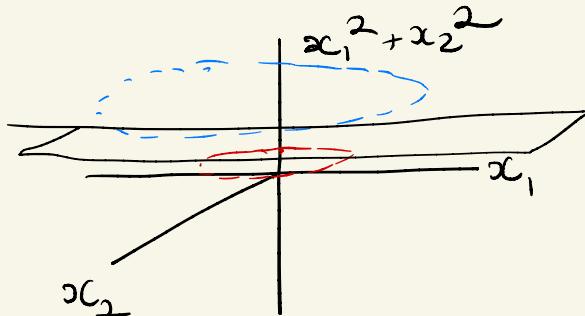
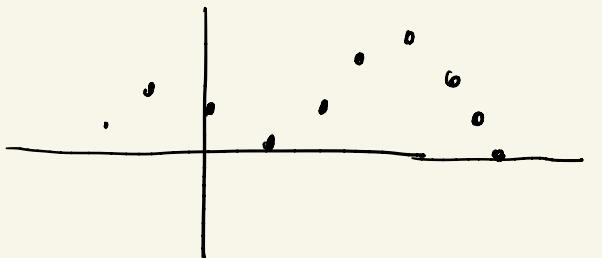


Last time : Linear Regression

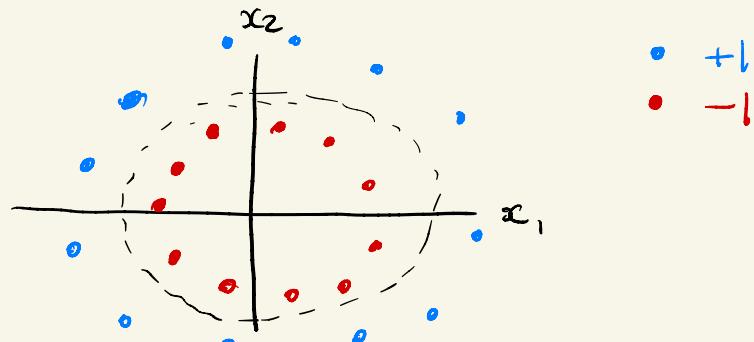
$$\{ f: \mathbb{R}^d \rightarrow \mathbb{R}; f(x) = \omega^\top x \}$$

Solution $\hat{f}(x) = y X^\top \omega$,
given by minimizing $\frac{1}{2} \|y - \omega^\top X\|_2^2$
square loss
or mean squared error.

Today : Kernel Regression



Ex. 1



$$\{x^{(i)}\}_{i=1}^n \rightarrow \{\psi(x^{(i)})\}_{i=1}^n \rightarrow \text{L.R.}$$

$$\begin{aligned}\psi: \mathbb{R}^2 &\rightarrow \mathbb{R}^3 \\ (x_1, x_2) &\mapsto (x_1, x_2, x_1^2 + x_2^2)\end{aligned}$$

$$\mathcal{F} = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} ; f(x) = \underbrace{\omega, \gamma(x)}_{\omega^\top \gamma(x)} ; \underbrace{\gamma: \mathbb{R}^d \rightarrow \mathbb{R}^p} \right\}$$

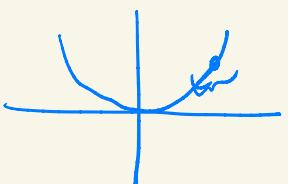
Q1: How do we deal with $p \gg d$?

$$L(\omega) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \omega^\top x^{(i)})^2 = \frac{1}{2} \|y - \omega X\|_2^2$$

$$\omega^{(t+1)} = \omega^{(t)} + \eta \underbrace{(y - \omega^{(t)} X)}_{B^{(t)}} X^\top \leftarrow$$

$B^{(t)}$
 $[B_1^{(t)}, B_2^{(t)}, \dots, B_n^{(t)}]$

$$\begin{bmatrix} -x^{(1)} \\ -x^{(2)} \\ \vdots \\ -x^{(n)} \end{bmatrix}$$



$$= \omega^{(t)} + n \sum_{j=1}^n \underbrace{B_j^{(t)} x^{(j)\top}}_{\omega^{(t)} \quad \quad \quad}$$

$$\omega^{(0)} = 0$$

$$\omega^{(1)} = n \sum_{j=1}^n B_j^{(0)} x^{(j)\top}$$

Guess: $\omega = \underbrace{\sum_{i=1}^n \alpha_i x^{(i)\top}}$

Guess: $\omega = \underbrace{\sum_{j=1}^n \alpha_j x^{(j) \top}}_{\text{; }} ; \quad L(\omega) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \omega x^{(i)})^2$

$$L(\omega) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \underbrace{\sum_{j=1}^n \alpha_j \underbrace{x^{(j) \top} x^{(i)}}_{\langle x^{(j)}, x^{(i)} \rangle}}_{\langle \alpha, \dots, \alpha_n \rangle})^2$$

$x^{(i)} \in \mathbb{R}^d$
 $y(x^{(i)}) \in \mathbb{R}^P$

$$\langle \underline{y(x^{(i)})}, \underline{y(x^{(i)})} \rangle$$

$$\begin{bmatrix} \alpha_1 & \dots & \alpha_n \end{bmatrix} \begin{bmatrix} \langle x^{(1)}, x^{(i)} \rangle \\ \vdots \\ \langle x^{(n)}, x^{(i)} \rangle \end{bmatrix}$$

$$\alpha \in \mathbb{R}^{1 \times n}$$

$$K(X, x^{(i)}) \in \mathbb{R}^n$$

$$= \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \alpha \underbrace{K(X, x^{(i)})}_{\mathbb{R}^n})^2$$

$$\alpha \in \mathbb{R}^{1 \times n} \quad \begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} K(X, x^{(1)}) & K(X, x^{(2)}) & \dots & K(X, x^{(n)}) \end{bmatrix} = y$$

$n \quad \boxed{n \times n}$

$$\boxed{\alpha K(X, x^{(i)}) = y}$$

$1 \times n \quad \omega \in \mathbb{R}^{1 \times d} \quad \boxed{\omega X = y}$

$1 \times d \quad \boxed{d \times n} \quad 1 \times n$

$$\frac{1}{2} \| \omega^T x - y \|_2^2 \Rightarrow \frac{1}{2} \| y - \underbrace{\alpha K(x, x)}_{n \times n} \|_2^2 \quad y = \underbrace{\alpha K(x, x)}_{n \times n}$$

$$w = yx^T \quad K(x, x) = \begin{bmatrix} K(x, x^{(1)}) & K(x, x^{(2)}) & \dots & K(x, x^{(n)}) \\ \vdots & \vdots & & \vdots \end{bmatrix} \quad yK(x, x)^{-1}$$

Step 1: Compute $K(x, x)$;

$$K(x, x)_{ij} = \langle x^{(i)}, x^{(j)} \rangle$$

Step 2: $\alpha = y K(x, x)^{-1}$

$$K(x, x)_{ij} = \langle \underline{y(x^{(i)})}, \underline{y(x^{(j)})} \rangle$$

Define: Kernel as a symmetric, positive definite function

Kernel as a symmetric, positive definite function

$$K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \quad K(x, z) \in \mathbb{R}$$

$$(1) \quad K(x, z) = K(z, x)$$

$$(2) \quad \langle \underline{y(x)}, \underline{y(x)} \rangle = \sum_{i=1}^n y(x_i)^2 \geq 0 \quad \leftarrow$$

$$\forall \{x^{(i)}\}_{i=1}^n, \forall \{c_i\}_{i=1}^n \subset \mathbb{R} \quad \sum_{i=1}^n \sum_{j=1}^n c_i c_j K(x^{(i)}, x^{(j)}) > 0.$$

$$\boxed{\forall \{x^{(i)}\}_{i=1}^n, \quad \forall \{c_i\}_{i=1}^n \in \mathbb{R} \quad \sum_{i=1}^n \sum_{j=1}^n c_i c_j K(x^{(i)}, x^{(j)}) \geq 0.}$$

Ex: $K(x, z) = \langle \psi(x), \psi(z) \rangle$

$$n=2; \quad \{c_1, c_2\}$$

$$\sum_{i=1}^2 \sum_{j=1}^2 c_i c_j K(x^{(i)}, x^{(j)}) = c_1^2 K(x^{(1)}, x^{(1)}) + c_1 c_2 K(x^{(1)}, x^{(2)}) \\ + c_2 c_1 K(x^{(2)}, x^{(1)}) + c_2^2 K(x^{(2)}, x^{(2)})$$

$$= c_1^2 K(x^{(1)}, x^{(1)}) + 2 c_1 c_2 K(x^{(1)}, x^{(2)}) + c_2^2 K(x^{(2)}, x^{(2)})$$

$$= c_1^2 \underbrace{\langle \psi(x^{(1)}), \psi(x^{(1)}) \rangle}_{\| \psi(x^{(1)}) \|_2^2} + 2 c_1 c_2 \underbrace{\langle \psi(x^{(1)}), \psi(x^{(2)}) \rangle}_{\sim (c_1 a + c_2 b)^2} + c_2^2 \| \psi(x^{(2)}) \|_2^2$$

$$= \boxed{\| c_1 \psi(x^{(1)}) + c_2 \psi(x^{(2)}) \|_2^2} \geq 0. \quad (c_1 a + c_2 b)^2$$

Examples: (1) $K(x, z) = \langle x, z \rangle$ (Linear)

$n \ll d$

(2) Laplace Kernel

$$K(x, z) = \exp \left\{ -\frac{\|x - z\|_2}{L} \right\}$$

\uparrow
bandwidth. $L > 0$

$L = 10$

Given: $(X_{\text{train}}, Y_{\text{train}})$, $K(x, z) \rightarrow \mathbb{R}$.

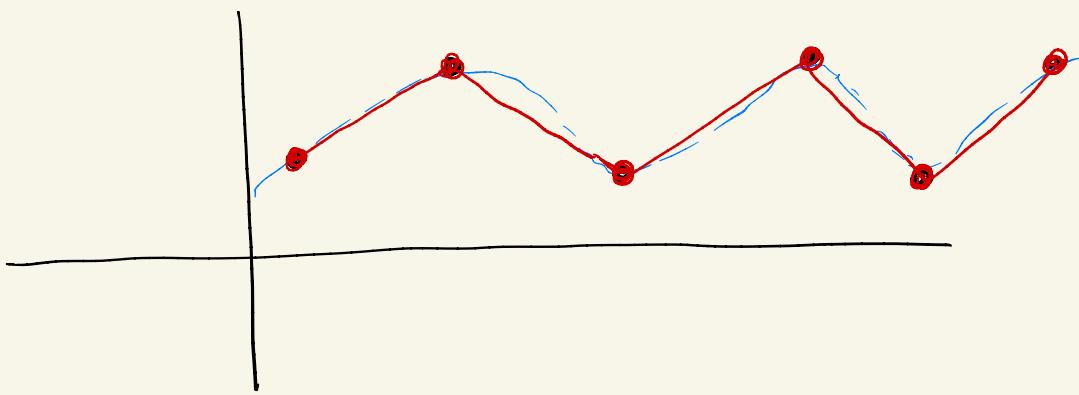
Algorithm:

Step 1: $K_{\text{train}} = K(X_{\text{train}}, X_{\text{train}})$

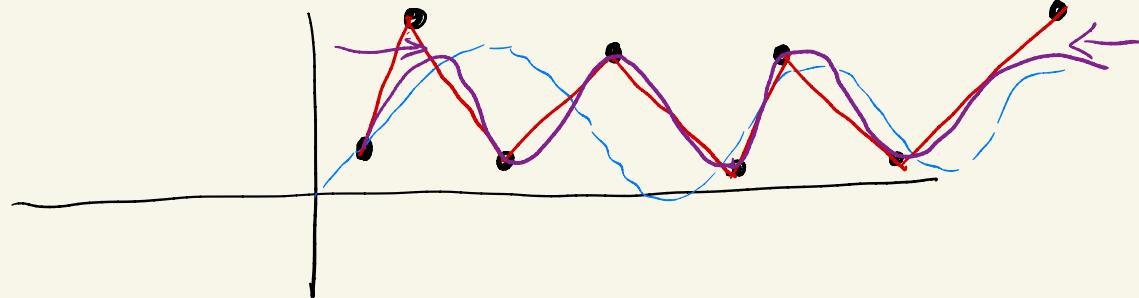
Step 2: $\alpha = Y K_{\text{train}}^{-1}$

Step 3: $K_{\text{test}} = K(X_{\text{train}}, X_{\text{test}})$

test-pred = αK_{test}



$$\alpha = \frac{yK^{-1}}{L=10}$$



Step 1: $K_{train} = K(X_{train}, X_{train})$

Step 2: $\alpha = \frac{y(K + \lambda I)^{-1}}{\lambda \in 10^{-2}, 10^{-3}}$

$$K(x, z) = \exp \left\{ - \frac{\|x - z\|_2}{L} \right\} \quad L \rightarrow \infty ; \lambda \mathbb{I}$$

$$K(X, X) = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = J$$

$$\alpha = y \underbrace{(K(X, X) + \lambda \mathbb{I})^{-1}}_{\overbrace{J}^{\mathbb{I}} \overbrace{A}^{\mathbb{I}}}$$

$$\begin{array}{c} u \rightarrow \mathbb{U} \quad \mathbb{U}^T \leftarrow v^T \\ \boxed{\cdot} \quad \boxed{\cdot \cdot \cdot \cdot} \end{array}$$

$$A^{-1} = \frac{1}{\lambda} \mathbb{I}$$

$$\alpha = y \left[\frac{1}{\lambda} \mathbb{I} - \frac{\frac{1}{\lambda} J}{1 + \frac{1}{\lambda} n} \right]$$

Sherman-Morrison
 $(A + uv^T)^{-1} =$

$$A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u}$$

$$\hat{f}(x) = \alpha \underbrace{K(X, x)}_{[\cdot]} = y \left[\frac{1}{\lambda} \mathbb{U} - \frac{\frac{n}{\lambda} \mathbb{U}}{1 + \frac{1}{\lambda} n} \right] = \sum_{i=1}^n y^{(i)} \left(\frac{1}{\lambda} - \frac{n}{\lambda^2 + \lambda n} \right) C$$