

# KERNEL METHODS - PRACTICAL SESSION № 1

AMMI 2020

18/05/2020

## Linear Algebra Recap

### Exercise 1

Let  $A \in \mathbb{R}^{m \times n}$  be a real matrix. For each function  $f$ , specify its codomain and compute  $\nabla f$

(a)  $f : \begin{cases} \mathbb{R} \rightarrow ? \mathbb{R}^{\text{sym}} \\ x \rightarrow xA \end{cases}$

(b)  $f : \begin{cases} \mathbb{R}^n \rightarrow ? \mathbb{R}^m \\ x \rightarrow Ax \end{cases}$

(c)  $f : \begin{cases} \mathbb{R}^{m \times n} \rightarrow ? \mathbb{R} \\ X \rightarrow \text{Tr}(A^T X) \end{cases}$

### Exercise 2

Let  $X \in \mathbb{R}^{n \times p}$ ,  $\lambda > 0$ .

**Prove that  $M = X^T X + \lambda I_p$  is invertible.**

*Hint: Prove that the eigenvalues of  $M$  are larger than  $\lambda$*

## Linear Regression

Using notations from the slides:

- $Y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$  the vector of outcomes
- $X = (x_1, \dots, x_n)^T \in \mathbb{R}^{n \times p}$  the matrix ( $n$  rows=samples,  $p$  columns=features)
- $\beta \in \mathbb{R}^p$  the linear model's parameters

### Exercise 3

- State the loss function for ordinary least squares (OLS)
- Formulate OLS as a minimization problem
- Compute the solution  $\hat{\beta}^{OLS}$
- What happens if  $X^T X$  is singular (not invertible)?

## **Exercise 4 - Bias and Variance of OLS**

slide 16

## **Exercise 5 - Optimality of OLS**

slide 17

## **Ridge Regression**

### **Exercise 6 - Bias and Variance**

slide 33

### **Exercise 7 - Performance**

slide 34

## **Ridge Logistic Regression**

### **Exercise 8**

slide 46

## EXERCISE 01

①  $A \in \mathbb{R}^{m \times m}$

a)  $f: \begin{cases} \mathbb{R} \rightarrow \mathbb{R}^{m \times m} \\ x \mapsto xA \end{cases}$

$$\nabla f_{ij} = f'_{ij}(x) = A_{ij} \quad ; \quad f_{ij}(x) = A_{ij}x \Rightarrow \boxed{\nabla f = A}$$

b)  $f: \begin{cases} \mathbb{R}^m \rightarrow \mathbb{R}^m \\ x \mapsto Ax \end{cases}$

$$\nabla f_{ij} = \frac{\partial f_i}{\partial x_j}$$

$$\begin{aligned} f_i(x) &= A_i^T x \\ &= \sum_k A_{ik} x_k \\ &= A_{ij} x_j \end{aligned}$$

$$\begin{aligned} i \in (1, m) \\ j \in (1, m) \end{aligned} = A_{ij}$$

$$\Rightarrow \boxed{\nabla f = A}$$

c)  $f: \begin{cases} \mathbb{R}^{m \times m} \rightarrow \mathbb{R} \\ X \mapsto \text{Tr}(A^T X) \end{cases}$

$\underbrace{\mathbb{R}^{m \times m} \times \mathbb{R}^{m \times m}}_{\mathbb{R}^{m \times m}}$

$\text{Tr}(M) = \sum$  of diagonal elements of  $M$ .

$$\nabla f = \frac{\partial f}{\partial x_{ij}} \quad \text{with } f(x) = \sum A_{ij} x_{ij}$$

$$\frac{\partial f}{\partial x_j} = A_{ij} \Rightarrow \boxed{\nabla f = A}$$

If we have  $f: \mathbb{R}^m \rightarrow \mathbb{R}$

$$\begin{aligned} f(x) &= f(x_0) + (x - x_0)^T \nabla f(x_0) + O(\|x - x_0\|) \\ &\approx f(x_0) + f'(x_0)(x - x_0) + O(x - x_0) \end{aligned}$$

$$\nabla' f_{ij} = \frac{\partial^2 f}{\partial x_{ij}}$$

## EXERCISE 2

Let  $X \in \mathbb{R}^{n \times p}$ ,  $\lambda > 0$ . Prove that  $M = X^T X + \lambda I_p$

is invertible.

LEMMA

$M' = X^T X$  eigen values  $\mu$

$$M' = \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix} \quad \det A = \prod_{\lambda \in \text{Eig}(A)} \lambda \quad \text{of } M' \text{ are } \geq 0.$$

$\mu \in \text{Eig}(M')$ ,  $y$  eigen Vector

$$M'y = \mu y$$

$$y^T M' y = y^T \mu y, \mu \geq 0$$

$$y^T M' y = \mu y^T y$$

$$y^T M' y = y^T X^T X y$$

$$= (Xy)^T Xy = \|Xy\|^2$$

$$\Rightarrow \mu = \frac{\|Xy\|^2}{\|y\|^2} \geq 0$$

$$y^T y = \sum y_i^2 \\ = \|y\|^2$$

$$\mu = \frac{\|Xy\|^2}{\|y\|^2}$$

$\Rightarrow$  we show that  $\mu$  are not neg.

$\mu \in \text{Eig}(M') \Leftrightarrow \mu + \lambda \in \text{Eig}(M)$

if  $\lambda' \in \text{Eig}(M)$ ,  $\lambda' \geq \lambda > 0$

All the eigenvalues of  $M$  are  $\geq 0 \Rightarrow M$  is invertible

## EXERCISE 3

a) State the loss function for ordinary least squares (OLS)

$$y = \beta^T x + \epsilon$$

$$= \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

$$\epsilon = y - \beta^T x$$

$$\text{The criteriom: } \text{MSE}(\beta) = \frac{1}{m} \sum (y_i - \beta^T x_i)^2$$

$$= \frac{1}{m} (y - X\beta)^T (y - X\beta)$$

b) Formulate OLS as a minimization Problem.

Ordinary Least Square

$$\min_{\beta} \text{MSE}(\beta), \text{MSE}: \mathbb{R}^p \rightarrow \mathbb{R}$$

c) Compute the solution  $\hat{\beta}^{\text{OLS}}$

$$\nabla \text{MSE}(\hat{\beta}) = 0$$

$$\nabla \text{MSE} = \frac{1}{m} [X^T(Y - X\beta) \cdot (Y - X\beta)^T X] = \frac{2}{m} X^T(X\beta - Y)$$

$$\nabla \text{MSE}_i = \frac{\partial \text{MSE}}{\partial \beta_i}$$

OLS Convex quadratic reaches optimum when  $\nabla \text{MSE}(\beta^*)$

$$X^T(X\beta^* - Y) = 0$$

$$X^T X \beta^* = X^T Y$$

If  $X^T X$  non singular  $\beta^* = (X^T X)^{-1} X^T Y$

$$A^T B = B^T A$$

$$\text{MSE} = -\frac{1}{m} [X^T(Y - X\beta) + (Y - X\beta)^T X] = -\frac{2}{m} X^T(Y - X\beta)$$

d) What happens if  $X^T X$  is singular (not invertible)

When  $X^T X$  is singular

$\Rightarrow$  Pseudo inverse  $\rightarrow$  Find Solution when there is no solution.

$\Rightarrow$  Regularize with small  $\lambda$ ; Ridge regression no closed form / unstable

EXERCISE 4 - BIAS AND VARIANCE OF OLS (SLIDE 16)

Show that  $\hat{\beta}^{\text{OLS}} = (X^T X)^{-1} X^T Y$

Satisfies  $\begin{cases} E(\hat{\beta}^{\text{OLS}}) = \beta^* \\ \text{Var}(\hat{\beta}^{\text{OLS}}) = E(\hat{\beta}^{\text{OLS}} - \beta^*)(\hat{\beta}^{\text{OLS}} - \beta^*)^T \end{cases}$

$$\text{Var}(A) = E((A - E(A))(A - E(A))^T)$$

1-dimension  $\text{Var}(A) = E(A - E(A))$

$$\text{Var}(A)_{ij} = \text{Cov}(A_i, A_j) = E(A_i A_j) = E(A_i) E(A_j)$$

$$y = X\beta^* + \varepsilon \quad \text{with } E(\varepsilon) = 0; \text{Var}(\varepsilon) = \sigma^2 I$$

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \quad \text{if } i \neq j$$

$$\text{Var}(\varepsilon_i) = \sigma^2$$

$$\begin{aligned}\hat{\beta}^{\text{OLS}} &= (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T X \beta^* + (X^T X)^{-1} X^T \varepsilon \\ &= \beta^* + \underbrace{(X^T X)^{-1} X^T \varepsilon}_A\end{aligned}$$

$\rightarrow$  E linearit

$$E[\hat{\beta}^{\text{OLS}}] = \beta^* + \underbrace{A E[\varepsilon]}_{E(A\varepsilon)} = \beta^* \Rightarrow \hat{\beta}^{\text{OLS}} \text{ is unbiased}$$

$$\text{Var}(A\varepsilon) = E((A\varepsilon)^T (A\varepsilon)) = E(A \varepsilon \varepsilon^T A) = A \underbrace{E(\varepsilon \varepsilon^T)}_{\sigma^2 I} A^T = \sigma^2 A A^T$$

$$A = (X^T X)^{-1} X^T$$

$$A A^T = (X^T X)^{-1} (X^T X) (X^T X)^{-1} = (X^T X)^{-1}$$

$$\text{Var}(\hat{\beta}) = \text{Var}(A\varepsilon) = \sigma^2 (X^T X)^{-1}$$

### EXERCISE 5

Gauss Markov theorem

Least squares estimator  $\hat{\beta}^{\text{OLS}}$  is a Best Linear Unbiased Estimator.

### Assumption

$$y = X\beta^* + \varepsilon$$

$$E\varepsilon = 0 \quad \text{Var}(\varepsilon) = \sigma^2 I$$

$\varepsilon$  is a matrix

if  $\tilde{\beta} = \tilde{C}y$  unbiased ( $E\tilde{\beta} = \beta^*$ ) then  $\text{Var}(\hat{\beta}^{\text{OLS}}) \leq \text{Var}(\tilde{\beta})$

$$\Delta \neq \hat{\beta}_{ij}^{\text{OLS}} \leq \tilde{\beta}_{ij}$$

$M = \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}^{\text{OLS}})$  is a positive semi-definite matrix

$\Leftrightarrow \forall x_0, x_0^T M$  $\hat{B} = Cy$  candidate unbiased estimator

$$D = C - (X^T X)^{-1} X^T$$

$$E(\hat{B}) = B^* = E(Cy) = C E(Y) = C X B^*, \quad y \in X B^* + \varepsilon \text{ and } E(y) = E(X B^*) +$$

 $\Rightarrow \hat{B}_{OLS}$  is unbiased ( $E(\hat{B})$ ) $\text{is Best Linear Unbiased}$ 

$$\begin{aligned} E(\varepsilon) &= 0 \\ &= X B^* \end{aligned}$$

$$i = [1, n]$$

$$y_i \in \mathbb{R}$$

$$y = \mu + \varepsilon$$

$$y_i = \mu + \varepsilon_i$$

$$E(\varepsilon_i) = 0$$

 $\text{Find an unbiased estimator of } \mu$ 

$$\Rightarrow \hat{y}_i = \frac{1}{m} \sum y_i$$

$$y_i = \mu + \varepsilon_i$$

$$E[y_i] = \mu$$

Unbiased

$$E(\hat{y}) = \frac{1}{m} \sum E(y_i) = \frac{1}{m} \sum_{i=1}^m \mu = \mu$$

 $\hat{y} = y_i$  is also an unbiased estimator

$$E(\hat{y}) = E(y_i) = \mu \quad \text{Now we have 2 unbiased estimator, which one is better}$$

Their Variance are  
different

$$\text{Var}(y_i) \leq \text{Var}(\hat{y}_i)$$

$$\text{Var}(y_i) \subset \text{Var}(\varepsilon)$$

$$\text{Var}(\hat{y}_i) = E((\hat{y}_i - E(\hat{y}_i))^2) = E\left(\frac{1}{m} \sum (y_i - \mu)^2\right)$$

$$\text{Var}(\hat{y}_i) = \frac{1}{m} \text{Var}(\varepsilon) \leq \text{Var}(y_i)$$

 $\hat{y}_i$  is better (smaller overall error), Error = bias + variance

## Back to (Ex 5)

$\hat{\beta}^{OLS}$  is unbiased we want to know if it is the best linear unbiased estimator

We need to take  $\tilde{\beta} = Cy$  linear unbiased estimator.

Show that  $\tilde{\beta}^{OLS}$  is better than  $\tilde{\beta}$ ,  $\text{Var}(\hat{\beta}^{OLS}) \leq \text{Var}(\tilde{\beta})$

$$\tilde{\beta} = Cy \quad D := C - (X^T X)^{-1} X^T$$

$$\text{unbiased } E(\tilde{\beta}) = \beta^*$$

$$\begin{aligned} E(\tilde{\beta}) &= E(Cy) = C E(y), \quad y = X\beta^* + \varepsilon \\ &= CX\beta^* + C E(\varepsilon) \\ &= CX\beta^* \end{aligned}$$

$$\beta^* = CX\beta^*, \quad (CX - I)\beta^* = 0$$

$$Dx = CX - (X^T X)^{-1} X^T X = CX - I$$

$$Dx \beta^* = 0$$

$\tilde{\beta}$  unbiased for all  $\beta^* \Rightarrow Dx = 0$

$$\text{Var}(\tilde{\beta}) = E((\tilde{\beta} - \beta^*)(\tilde{\beta} - \beta^*)^T)$$

$$\tilde{\beta} = Cy = (D + (X^T X)^{-1} X^T)(X\beta^* + \varepsilon) = \beta^* + D\varepsilon + (X^T X)^{-1} X^T \varepsilon$$

$$\text{Var}(\tilde{\beta}) = E((\tilde{\beta} - \beta^*)(\tilde{\beta} - \beta^*)^T) = D\varepsilon + (X^T X)^{-1} X^T \varepsilon$$

$$\begin{aligned} &= E(D\varepsilon (D\varepsilon)^T + 2(X^T X)^{-1} X^T \varepsilon \varepsilon^T X (X^T X)^{-1} + D\varepsilon \varepsilon^T X (X^T X)^{-1} + (X^T X)^{-1} X^T \varepsilon \varepsilon^T D) \\ &= D E(\varepsilon \varepsilon^T) D^T (X^T X)^{-1} X^T E(\varepsilon \varepsilon^T) X (X^T X)^{-1} + D E(\varepsilon \varepsilon^T) X (X^T X)^{-1} + \\ &\quad (X^T X)^{-1} X^T E(\varepsilon \varepsilon^T) D^T \\ &= \sigma^2 D D^T + \sigma^2 (X^T X)^{-1} + 0 + 0 \end{aligned}$$

$\text{Var}(\tilde{\beta}) = \sigma^2 D D^T + \text{Var}(\hat{\beta}^{OLS})$ , we want to show that  $M$  is posit. definite

$\text{Var}(\tilde{\beta}) = \text{Var}(\hat{\beta}^{OLS}) = \sigma^2 D D^T$ ;  $\text{Var}(\tilde{\beta}) \leq \text{Var}(\hat{\beta}^{OLS}) \Leftrightarrow X_0^T M X_0 \geq 0$

$$X_0^T M X_0 = \sigma^2 X_0^T D D^T X_0 = \sigma^2 (D X_0)^T D^T X_0 = \sigma^2 \|DX_0\|^2 \geq 0, \forall x_0$$

$\text{Var}(\tilde{\beta}) \geq \text{Var}(\hat{\beta}^{OLS})$  then  $\hat{\beta}^{OLS}$  is BLUE

# Ridge Regression

## EXERCISE 6

Slides 33

$$\begin{cases} \text{bias}(\hat{\beta}_{\text{Ridge}}) = ? & = \frac{\lambda}{1+\lambda} \beta^* \\ \text{Var}(\hat{\beta}_{\text{Ridge}}) = ? & = \frac{1}{(1+\lambda)^2} \text{Var}(\hat{\beta}_{\text{OLS}}) \end{cases}$$

$$L(\beta) = \text{MSE}(\beta) + \lambda \|\beta\|^2 = \frac{1}{m} \sum (y_i - \beta^T x_i)^2 + \lambda \|\beta\|^2$$

Ridge regression loss

$$\nabla L(\beta) = 0$$

$$\hat{\beta} = (x^T x + \lambda m I)^{-1} x^T y$$

$\lambda > 0$ ,  $x^T x + \lambda m I$  is invertible.

$$\nabla L(\beta) = \frac{2}{m} x^T (x\beta - y) + 2\lambda \beta$$

$$\nabla = 0 \Rightarrow (x^T x + \lambda m I) \hat{\beta} = x^T y$$

$$\begin{aligned} E(\hat{\beta}_{\text{Ridge}}) &= (x^T x + \lambda m I)^{-1} x^T E(y) \\ &= (x^T x + \lambda m I)^{-1} \underbrace{x^T x}_{X^T X} \beta^* \\ &\quad x^T x + \lambda m I = \lambda m I \\ &= \beta^* - \lambda m (x^T x + \lambda m I)^{-1} \beta^* \end{aligned}$$

Assumption :  $x^T x = m I$

$$E(\hat{\beta}_{\text{Ridge}}) = \beta^* - \frac{\lambda}{1+\lambda} \beta^*$$

$$\text{bias}(\hat{\beta}_{\text{Ridge}}) = -\frac{\lambda}{1+\lambda} \beta^*$$

$\Rightarrow$  When  $\lambda = 0$  (no regularization),  $\text{bias} = 0$ ,  $\hat{\beta}_{\text{Ridge}} = \hat{\beta}_{\text{OLS}}$

When  $\lambda = +\infty$ ,  $\text{bias} = -\beta^*$  because  $\hat{\beta}_{\text{Ridge}} = 0$

## Assumption

$$X^T X = mI$$

$$X_i^T X_j = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

$$\hat{\beta}_\lambda^{\text{Ridge}} = \frac{1}{m(1+\lambda)} X^T Y$$

$$\text{Var}(\hat{\beta}_\lambda^{\text{Ridge}}) = \frac{1}{m^2(1-\lambda)^2} \text{Var}(X^T Y) = \frac{1}{m^2(1-\lambda)^2} \frac{\text{Var}(X^T Z)}{m\sigma^2}$$

$$\text{Var}(X^T Z) = E(X^T \epsilon \epsilon^T X) = X^T \text{Var}(\epsilon) X = \sigma^2 X^T X = m\sigma^2$$

$$\text{Var}(\hat{\beta}_\lambda^{\text{Ridge}}) = \frac{\sigma^2}{m(1+\lambda)^2} ; \text{Var}(\hat{\beta}_{\text{OLS}}) = (X^T X)^{-1} \sigma^2 = \frac{\sigma^2}{m}$$

$$\text{Var}(\hat{\beta}_\lambda^{\text{Ridge}}) = \frac{1}{(1+\lambda)^2} \text{Var}(\hat{\beta}_{\text{OLS}})$$

when  $\lambda \rightarrow 0$  : Same Variable (Same estimator)

$\lambda \rightarrow +\infty$  : Variable  $\rightarrow 0$  (estimator  $\hat{\beta}^T X = 0$ )

## EXERCISE 7.

$$f(x) \triangleq \mathbb{E}[\text{bias}^2(X^T \hat{\beta}_\lambda^{\text{Ridge}}) + \text{Var}(X^T \hat{\beta}_\lambda^{\text{Ridge}})]$$

$$\text{1 minimum for } \lambda \in \mathcal{X}^* = \frac{\sigma^2 p}{m \|\beta^*\|^2}$$

(Compute  $f'(x)$ )

$f(x)$  = error function  $\exists$  optimal  $\lambda^*$  that gives the best fit  
cannot compute it directly. Find with cross-validation

$$f(\lambda^*) \leq \min(f(0), f(+\infty))$$

Using  $\lambda^*$  is better than not using regularization.

