

- Underlying model of the data

$X$ : input / covariate      Ex:  $X: 2, 5, 8$

$Y$ : output / response       $Y: 11$

$(X, Y)$  is a joint random vector, drawn from a probability distribution  $P$  (which might be horribly complicated)

- Prediction is not always possible

↳ sometimes there is inherent noise in the data

Ex.  $X$  and  $Y$  are independent

- No learning algorithm is universal

**No-free-lunch theorem:** for any learning algorithm  $A$ ,

there exists a data distribution  $P$  s.t.  $A$  fails.      Ex. Number Seq.

**Consequence:** Any ML algorithm needs to make assumptions on how data is generated.

This course:

- Describe standard assumptions
- Introduce a statistical framework to apply / study performance of the ML algorithms.
- Cover main learning algorithms.

## Machine Learning Paradigms

• A no-free-lunch theorem:

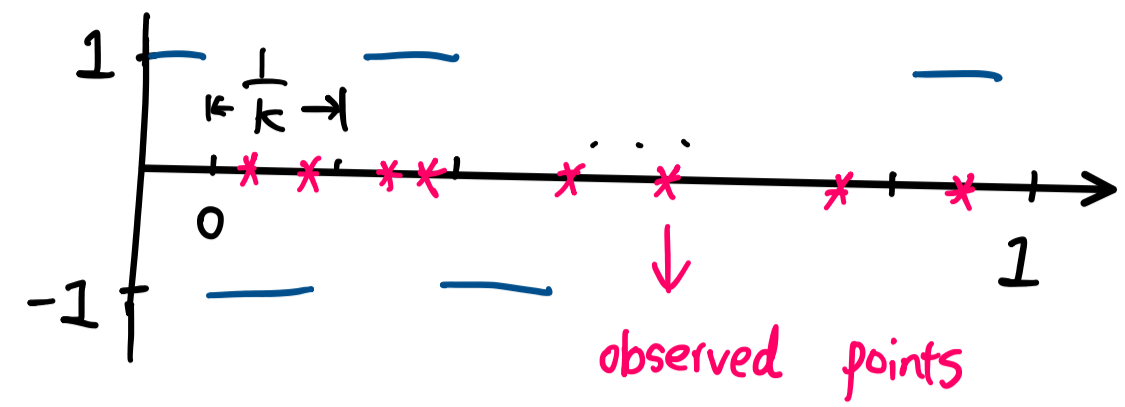
(We will give an instance of an "impossible" learning problem)

Let  $n$ : number of training examples

Consider  $k \gg n$

for  $r_j = \begin{cases} +1 & \text{with prob. } \frac{1}{2} \\ -1 & \text{with prob. } \frac{1}{2} \end{cases}$  & independent  
 $r_j = \pm 1, j = 1, \dots, k$

Target Function (Random)



Now let's consider a training set  $\{(x_i, y_i)\}_{i=1, \dots, n}$

$x_i \sim \text{Unif}([0, 1])$

$y_i = r_{\lfloor kx_i \rfloor} \rightarrow$  lower integral

$\rightarrow$  If  $x$  is drawn s.t.  $\lfloor kx \rfloor \neq \lfloor kx_i \rfloor$  for  $\forall i$

then can any learning algorithm predicts  $r_{\lfloor kx \rfloor}$ ?

$\rightarrow$  Any learning algorithm  $A(x)$  only depends on R.V.  $\{r_j; j \text{ in the training set}\}$

$\Rightarrow$  Prediction for such  $x$  is limited to random guess

$$P(A(x) \neq y | X \text{ is not in a observed bracket}) = \frac{1}{2}$$

Hence,

$$P(A(x) \neq y) = \frac{1}{2} P(X \text{ is not observed}) \geq \frac{1}{2} (1 - \frac{n}{k})$$

for any algorithms  $A(\cdot)$

$\rightarrow$  Interpretation: there is no interaction between training and testing

$\rightarrow$  Contrast with because smoothness of this curve creates dependence between train and test

## Machine Learning Paradigms

① Simplest setting (focus of this course): Supervised Learning (SL)

dataset of labeled examples  $\{(x_i, y_i)\}$

$x_i \in \mathcal{X} \rightarrow$  input feature space  $\mathcal{X} \stackrel{\text{e.g.}}{=} \{\text{natural images}\}$

$\stackrel{\text{e.g.}}{=} \{\text{text sequences}\}$

$y_i \in \mathcal{Y} \rightarrow$  label space  $\mathcal{Y} = \mathbb{R}$  for regression e.g. (predicting temp.)

$\mathcal{Y} = \{1, \dots, k\}$  for classification (category)

$\mathcal{Y} = \mathbb{R}^d$  (protein folding) ["structured prediction"]

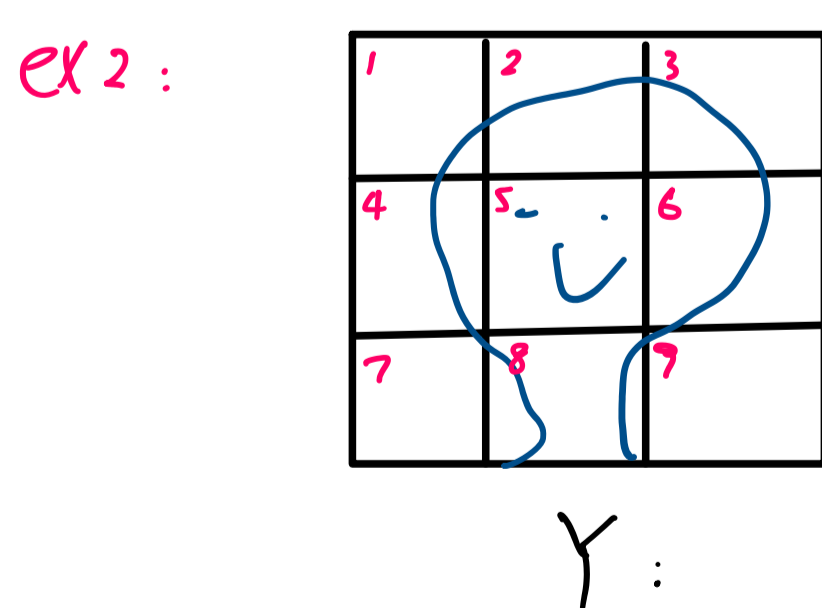
② Important special case of SL: Self-supervised Learning (SSL)

$\rightarrow$  We define the label  $\mathcal{Y}$  ourselves from unlabeled data.

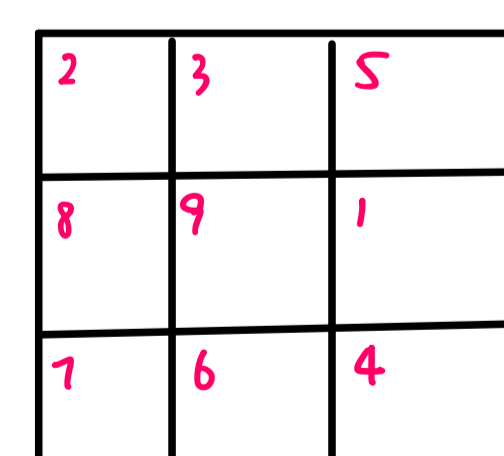
(video)  
 ex1:  $z_1, z_2, z_3, \dots, z_t, \dots$  future

Sequential Data  
 past present  
 $x_t y_t$   
 "input" "label"

"Representation Learning"  
 Condition: Near future being highly dependent on the presence



shuffle  $\rightarrow$



$\mathcal{X}$ : set of 9 patches shuffled

③ Unsupervised Learning

$\rightarrow$  Consider an unlabeled dataset  $\mathcal{D} = \{x_i\}$ . Extracting "original" information out  $\mathcal{D}$

$\rightarrow$  Geometric Encoding  $X \rightarrow \boxed{\Phi} \rightarrow Z \rightarrow \boxed{\Psi} \rightarrow \mathcal{X}$  [autoencoder]  
 Encoder Latent space Decoder

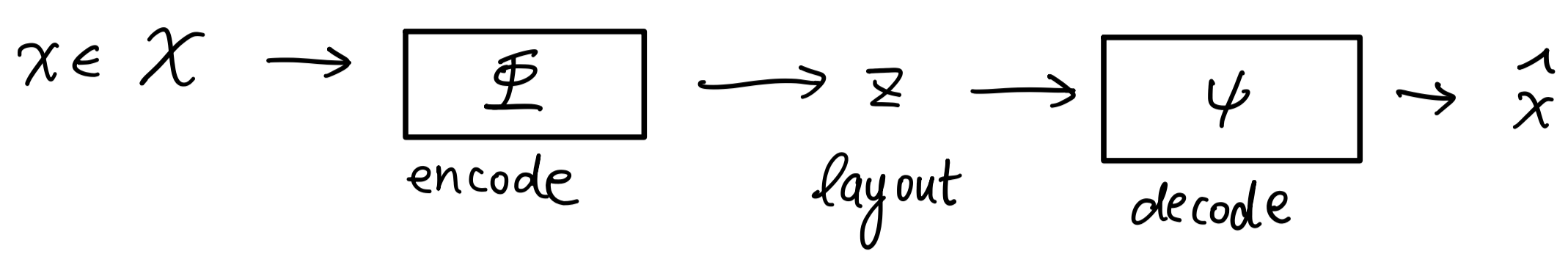
Design encoder/decoder to minimize  $\min_{\Phi, \Psi} E \|X - \Psi(\Phi(X))\|^2$

Canonical Ex.  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$   
 $\psi: \mathbb{R}^k \rightarrow \mathbb{R}^d$   $1 \ll k \ll d$  linear maps [PCA]

FoML Lecture 3: Linear Regression I

- Recap from Lec. 2
  - \* Impossibility of Learning (NFLT)
  - \* Several paradigms of ML
    - Supervised Learning
    - Self-supervised learning
    - Unsupervised learning

Unsupervised Learning  
 "discovering" hidden structure in data



ex. Principal Component Analysis (PCA)  
 where  $\Phi$  and  $\Psi$  are linear maps

Probabilistic view on unsupervised learning

Input  $x_i \in \mathcal{X}$ ,  $i \in \{1, \dots, n\}$  is viewed as  $n$  i.i.d. samples of an unknown prob. distribution  $p$ .

↳ Unsupervised learning to estimate  $\hat{p}$  from samples  $\{x_i\}$

Main Application: generative modeling use  $\hat{p}$  to draw new sample (Dall-E, ChatGPT, ...)

Semi-supervised Learning

Large unlabeled dataset  $D = \{x_i\}_i$

Small labeled dataset  $D' = \{x'_i, y'_i\}_i$

Assumptions:  $x_i$  and  $x'_i$  are drawn from some distribution

Goal: Combine  $D$  and  $D'$  to "propagate" labels

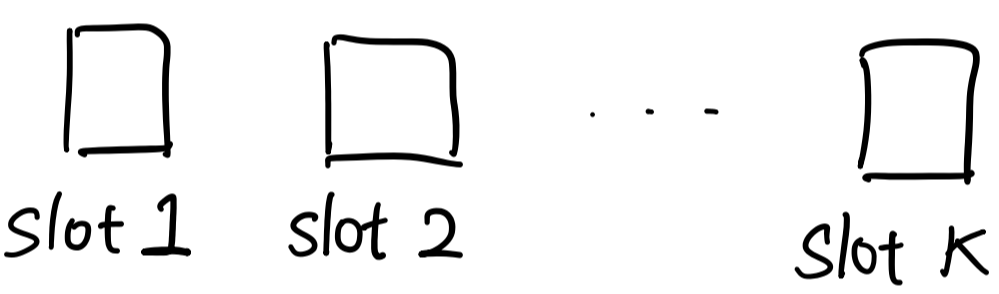
(If  $x_i$  is 'similar' to  $x'_i$ , then  $y_i$  should be similar to  $y'_i$ )

Online and Reinforcement Learning

→ So far learning has been passive

→ Learning is also the ability to act and adapt to changing adversified environment

eg 1: Bandit Problem



Each slot is modeled as a distribution  $\mathcal{U}_k$  with  $r_j \sim \mathcal{U}_j \in [0, +\infty)$ ,  $j = 1, \dots, K$

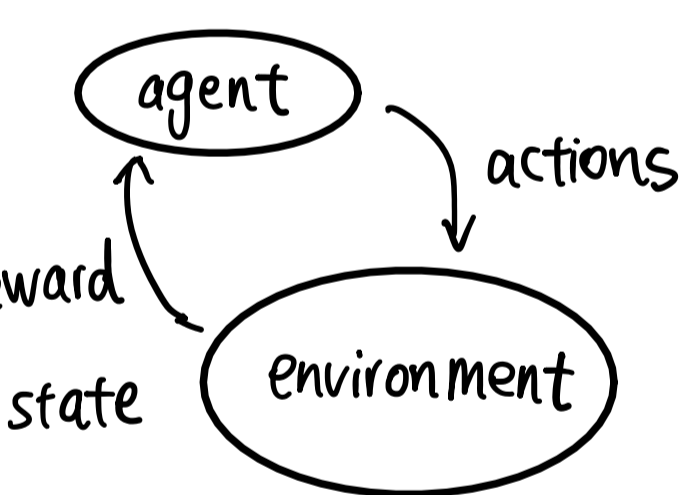
Each round  $t$  the player picks a slot  $S_t \in \{1, \dots, K\}$  (we can pick the same slot repeatedly)

Goal: After  $T$  rounds, maximize reward  $\sum_{t=1}^T r_{S_t}$

Online Learning key aspect: tradeoff between exploration and exploitation

eg → Allocation of Research Budget

→ Clinical Trials



→ Extension to settings when environment depends on an agent  
 Reinforcement Learning Structure

eg: Games (GO, Chess, Mazes)

Robotics

Focus on Supervised Learning

Regression Problems: We are given a dataset  $\{(x_i, y_i)\}_{i=1, \dots, n}$

$x_i \in \mathcal{X}$  inputs (features)  $y_i \in \mathcal{Y}$  output (response)  $(x_i, y_i) \sim P$  (i.i.d.)

still inside the sample space (i.e., we are doing training only)

Goal: Given a new  $x \in \mathcal{X}$  drawn from  $P_x$ , predict  $f(x) \approx y$

Regression:  $\mathcal{Y} = \mathbb{R}$ , measure errors using square loss  $\ell(y, \hat{y}) = (y - \hat{y})^2$

Q: What is the optimum least squares predictor? predictor  $f: \mathcal{X} \rightarrow \mathcal{Y}$

$$\text{find } \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_p [(f(x) - y)^2]$$

eg1: If  $\exists f^*$  s.t.  $f^*(y) = x$ , then the best predictor is  $f^*$  which makes  $\ell = 0$

eg2: If  $y \sim \mathcal{N}(0, 1)$  and indep. of  $x$ , then the best predictor is  $f = 0$

$$\mathbb{E}_{x, y} [(f(x) - y)^2] = \mathbb{E}_x \mathbb{E}_{y|x} [(f(x) - y)^2]$$

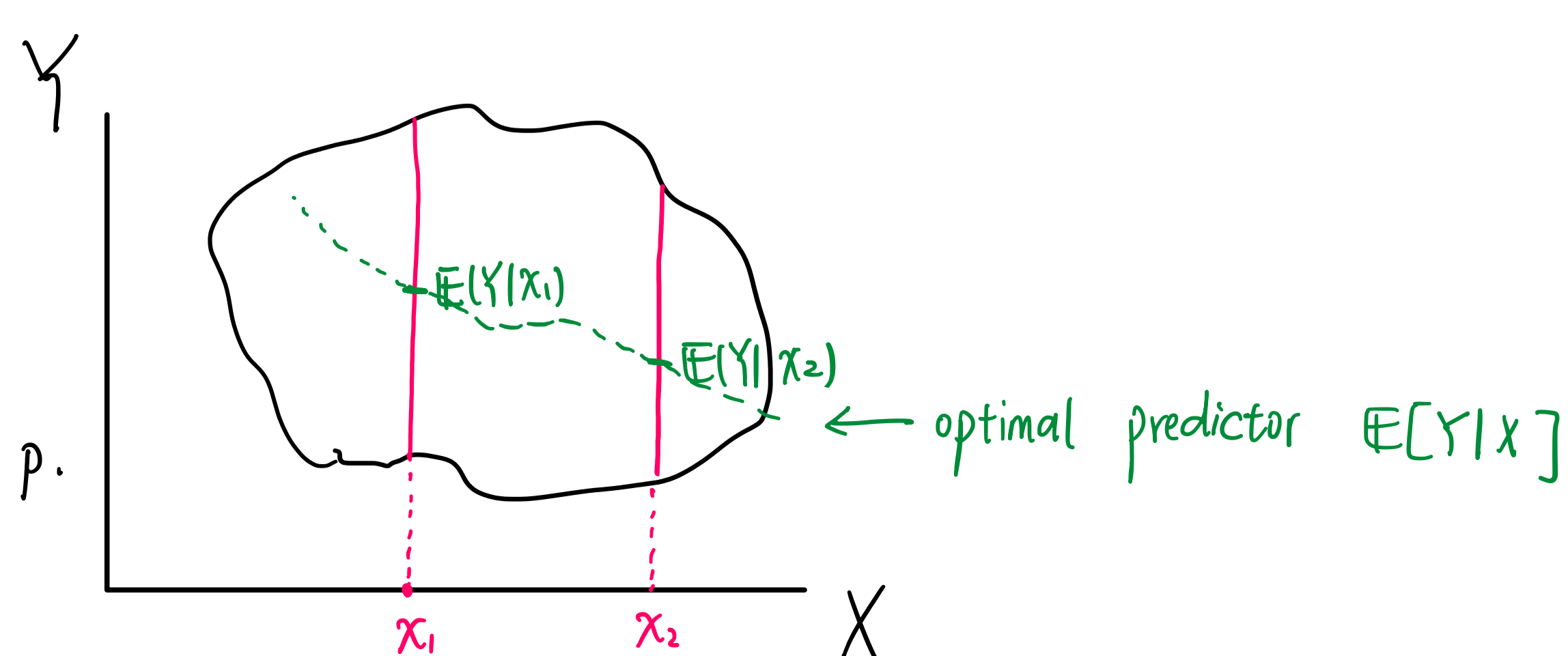
Hence Goal  $\Leftrightarrow \min_{c \in \mathbb{R}} \mathbb{E}_z [(z - c)^2] \Rightarrow$  take  $c = \mathbb{E}z$  because  $\mathbb{E}[(z - c)^2] = \underbrace{(c - \mathbb{E}[z])^2}_{\text{bias}} + \underbrace{\text{Var}(z)}_{\text{variance}}$

Hence we choose our optimum predictor as  $f^*(x) = \mathbb{E}[y|x]$

Conclusion 1: Optimal predictor is  $f_*(x) = \mathbb{E}_p [Y|X = x]$

Conclusion 2: This optimal predictor is unknown in general as we don't know  $p$ .

Instead, we will find the optimal function in the linear span of a set of predetermined predictors



FML Lecture 4 : Linear Regression II

Recap: Regression Problem  $\min_{f: X \rightarrow Y} \mathbb{E}[|f(x) - y|^2]$

→ Optimal solution  $f^*(x) = \mathbb{E}_p[Y|X=x]$

Linear Regression

- $f_1$
  - $f_2$
  - ⋮
  - $f_d$
- candidate solutions

Regression Model :

$f_{\Theta}(x) = \theta_1 f_1(x) + \theta_2 f_2(x) + \dots + \theta_d f_d(x)$ ,  $\Theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{pmatrix} \in \mathbb{R}^d$   
 linear combinations of candidate solutions

→  $f_{\Theta}(x)$  is linear in  $\Theta$ :  $f_{\alpha\Theta + \alpha'\Theta'} = \alpha f_{\Theta} + \alpha' f_{\Theta'}$

→ But  $f_{\Theta}$  is NOT linear w.r.t.  $x$  !! (as  $f_i$  could be nonlinear)

eg. How bitter is an espresso shot ?

$x$ : barista coffee makers  $y$ : acidity level

$f_1(x)$  = temperature of water

$f_2(x)$  = altitude of beans

$f_3(x)$  = pressure

$f_4(x)$  = pressure<sup>2</sup>

⋮

• Given observations  $x_1, \dots, x_n$  and candidate solutions  $f_1, \dots, f_d$

Then linear regression is  $\min_{\Theta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n [y_i - \underbrace{\sum_{j=1}^d \theta_j f_j(x_i)}_{f_{\Theta}(x_i)}]^2 = \min_{\Theta \in \mathbb{R}^d} \hat{R}(\Theta) \rightarrow$  Empirical Risk  
 ( $\hat{R}$  is a random function)

Collect features in a matrix  $\hat{H} = \begin{bmatrix} f_1(x_1) & f_2(x_1) & \dots & f_d(x_1) \\ f_1(x_2) & \dots & \dots & f_d(x_2) \\ \vdots & \dots & \dots & \vdots \\ f_1(x_n) & \dots & \dots & f_d(x_n) \end{bmatrix} \in \mathbb{R}^{n \times d}$   
 labels in a vector  $\hat{y} = (y_1, \dots, y_n) \in \mathbb{R}^n$

Then  $\hat{R}(\Theta) = \frac{1}{n} \|\hat{y} - \hat{H}\Theta\|^2$

- Tasks :
- (1) How to minimize  $\hat{R}(\Theta)$
  - (2) Geometric Interpretation
  - (3) Statistical Analysis

(1) The Normal Equations

Assumption: The matrix  $\hat{H}$  has rank  $d$ . In particular,  $n \geq d$ .

↳ direct consequence: the associated Gram matrix  $\hat{K} = \frac{1}{n} \hat{H}^T \hat{H} \in \mathbb{R}^{d \times d}$   
 where  $\hat{K}_{i,j} = \frac{1}{n} \langle \hat{H}_i, \hat{H}_j \rangle$  → columns

Moreover,  $K^T = K$  and  $K$  is invertible [consequence of SVD]

for  $\forall z \in \mathbb{R}^d$ ,  $z^T \hat{K} z = \frac{1}{n} z^T \hat{H}^T \hat{H} z = \frac{1}{n} \|\hat{H}z\|^2 > 0$  as  $\hat{H}$  has rank  $d$  ( $z \neq 0$ )

Therefore  $\hat{K}$  is positive-definite.

• Now  $\hat{R}(\Theta) = \frac{1}{n} \|\hat{y} - \hat{H}\Theta\|^2$   
 $= \frac{1}{n} \|\hat{y}\|^2 + \frac{1}{n} (\hat{H}\Theta)^T \hat{H}\Theta - \frac{2}{n} \hat{y}^T \hat{H}\Theta$   
 $\quad \quad \quad \Theta^T \hat{K} \Theta$

Claim:  $\hat{R}$  is convex.

Then  $\nabla \hat{R}(\hat{\Theta}) = 2 \hat{K} \hat{\Theta} - \frac{2}{n} (\hat{y}^T \hat{H})^T = 0$  → make dimensions coincide

$\Rightarrow \hat{\Theta} = \frac{1}{n} \hat{K}^{-1} \hat{H}^T \hat{y} = \frac{1}{n^2} (\hat{H}^T \hat{H})^{-1} \hat{H}^T \hat{y}$  [Normal Equations]  
 Legendre, early 19th century

Associated Risk :

$\hat{R}(\hat{\Theta}) = \frac{1}{n} \|\hat{y}\|^2 + \frac{1}{n^2} \hat{y}^T \hat{H} \underbrace{(\hat{K}^{-1})}_{In} \hat{K} \hat{K}^{-1} \hat{H}^T \hat{y} - \frac{2}{n^2} \hat{y}^T \hat{H} \hat{K}^{-1} \hat{H}^T \hat{y}$   
 $= \frac{1}{n} \hat{y}^T \hat{y} - \frac{1}{n^2} \hat{y}^T \hat{H} \hat{K}^{-1} \hat{H}^T \hat{y}$   
 $= \frac{1}{n} \hat{y}^T \left( I_n - \frac{1}{n} \hat{H} \hat{K}^{-1} \hat{H}^T \right) \hat{y}$

FML Lecture 5: Linear Regression (cont'd) Fixed Design

Recap from last lecture:

Dataset  $\{(X_i, Y_i)\}_{i=1, \dots, n}$ ,  $(X_i, Y_i) \sim \mathcal{P}$   
 Linear Regression Model:  $f_{\theta}(x) = \theta^T H(x)$  where  $H(x) = (H_1(x), \dots, H_d(x)) \in \mathbb{R}^d$   
 Risk:  $R(\theta) = \mathbb{E} \|Y - \theta^T H(x)\|^2$   
 Empirical Risk:  $\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n [y_i - \theta^T H(x_i)]^2 = \frac{1}{n} \|\hat{y} - \hat{H}\theta\|^2$  where  $\hat{y} = (y_1, \dots, y_n) \in \mathbb{R}^n$

Feature Matrix:  $\hat{H} = [H_j(x_i)]_{i=1, \dots, n, j=1, \dots, d} \in \mathbb{R}^{n \times d}$  (We assume  $H_1, \dots, H_d$  are l.i., i.e.,  $\hat{H}$  has rank  $d$ )

**Ordinary** why ordinary?  $\hat{\theta} = \frac{1}{n} \hat{K}^{-1} \hat{H}^T \hat{y}$ , where  $\hat{K} = \frac{1}{n} \hat{H}^T \hat{H}$   $\hat{\theta}$ : because the cost function does not have regularization term.

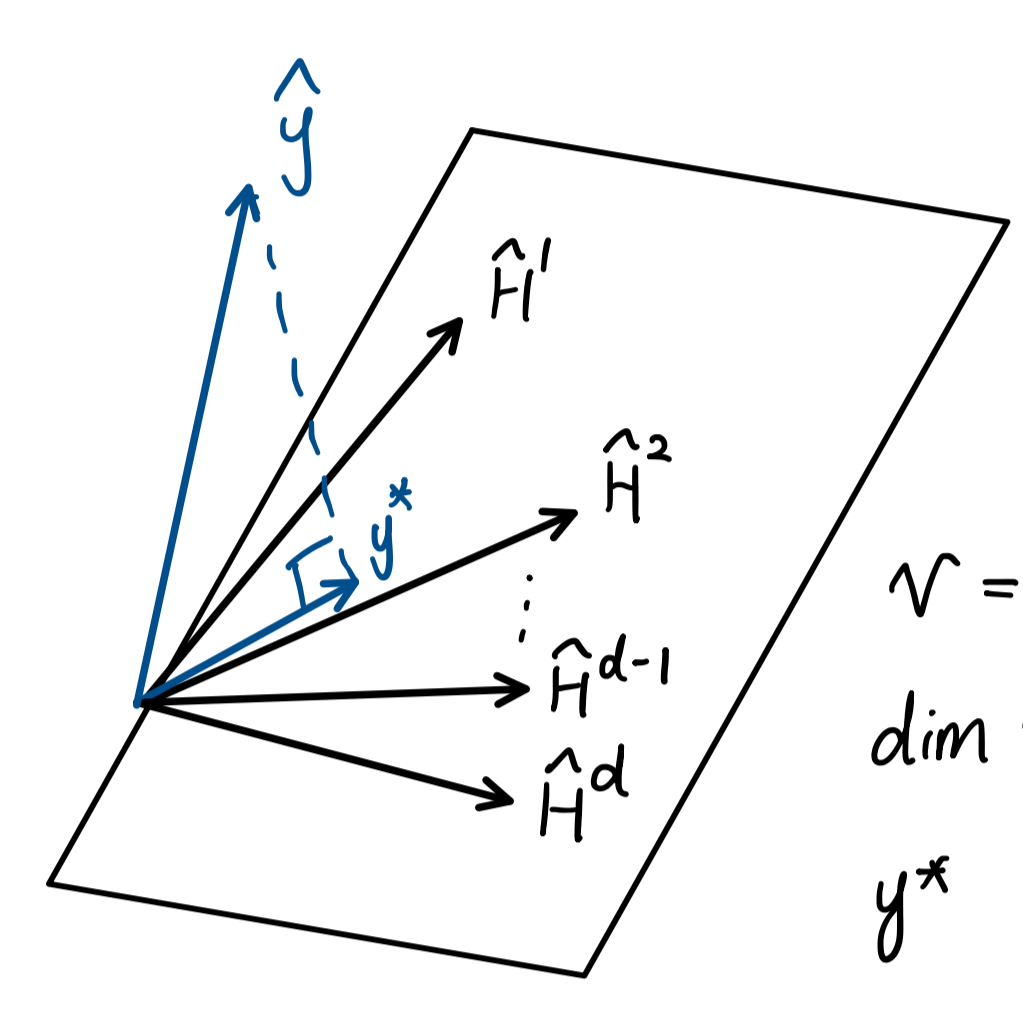
Best-Empirical Risk:  $\hat{R}(\hat{\theta}) = \frac{1}{n} \hat{y}^T [\mathbb{I}_n - \frac{1}{n} \hat{H} \hat{K}^{-1} \hat{H}^T] \hat{y}$   
 $\Pi = \hat{H} (\hat{H}^T \hat{H})^{-1} \hat{H}^T \in \mathbb{R}^{n \times n}$

Today: (1) Geometric Interpretation  
 (2) Statistical Analysis

Let  $\Pi = \hat{H} (\hat{H}^T \hat{H})^{-1} \hat{H}^T \in \mathbb{R}^{n \times n}$

Q: How to interpret the OLS solution?

$\hat{H} = \begin{bmatrix} | & & | \\ \hat{H}^1 & \dots & \hat{H}^d \\ | & & | \end{bmatrix}$  where  $\hat{H}^j \in \mathbb{R}^n, j=1, \dots, d$



Rmk.  $y^* = \hat{H} \hat{\theta}$  as  $\hat{\theta}$  is the orthogonal projection of  $\hat{y}$  onto  $\mathcal{V}$ . We will show  $\Pi \hat{y} = \hat{H} \hat{\theta}$ .  
 i.e.,  $\Pi$  is the orthogonal projector of  $\hat{y}$  onto  $\mathcal{V}$   
 $\mathcal{V} = \text{span}(H^1, \dots, H^d)$   
 $\dim \mathcal{V} = d$   
 $y^*$  is the orthogonal projection of  $\hat{y}$  onto  $\mathcal{V}$

A:  $\hat{H} \hat{\theta} = \Pi \hat{y}$ , as the orthogonal projection of  $\hat{y}$  onto  $\text{Col}(\hat{H})$

Pf.  $\min_{\theta \in \mathbb{R}^d} \|\hat{y} - \hat{H}\theta\|^2 = \min_{v \in \mathcal{V}} \|\hat{y} - v\|^2 = \text{Proj}_{\mathcal{V}}(\hat{y})$

- Equivalent properties of orthogonal projector  $P$ :
  - $x \in \mathcal{V}, Px = x$
  - $x \in \mathcal{V}^\perp, Px = 0$

We need to show that  $\Pi$  is the orthogonal projector onto  $\mathcal{V}$

Verify Property (1): Let  $x = \hat{H}\theta$  for some  $\theta$

$\Pi x = \hat{H} (\hat{H}^T \hat{H})^{-1} (\hat{H}^T \hat{H}) \theta = \hat{H} \theta = x \checkmark$

Property (2): for  $x \in \mathcal{V}^\perp \Leftrightarrow x \perp H^j$  for  $\forall j=1, \dots, d$   
 i.e.,  $\mathcal{V}^\perp = \text{null}(H^T)$   
 $\Rightarrow x \in \mathcal{V}^\perp, \hat{H}^T x = 0 \Rightarrow \Pi x = 0$

Hence  $\Pi$  is an orthogonal projector.

#

Statistical Analysis of Least Square

We distinguish two frameworks, to study generalization in linear regression

(i) "Random Design": We view  $(X_i, Y_i)$  as a random vector drawn from unknown distribution  $\mathcal{P}$

Test time  $R(\hat{\theta}) = \mathbb{E}_{\mathcal{P}} [(y - \hat{\theta}^T H(x))^2]$  (Random Setting)  
 Training data  $(X, Y) \sim \mathcal{P}$   
 Talk a bit more on this?  $\Rightarrow$

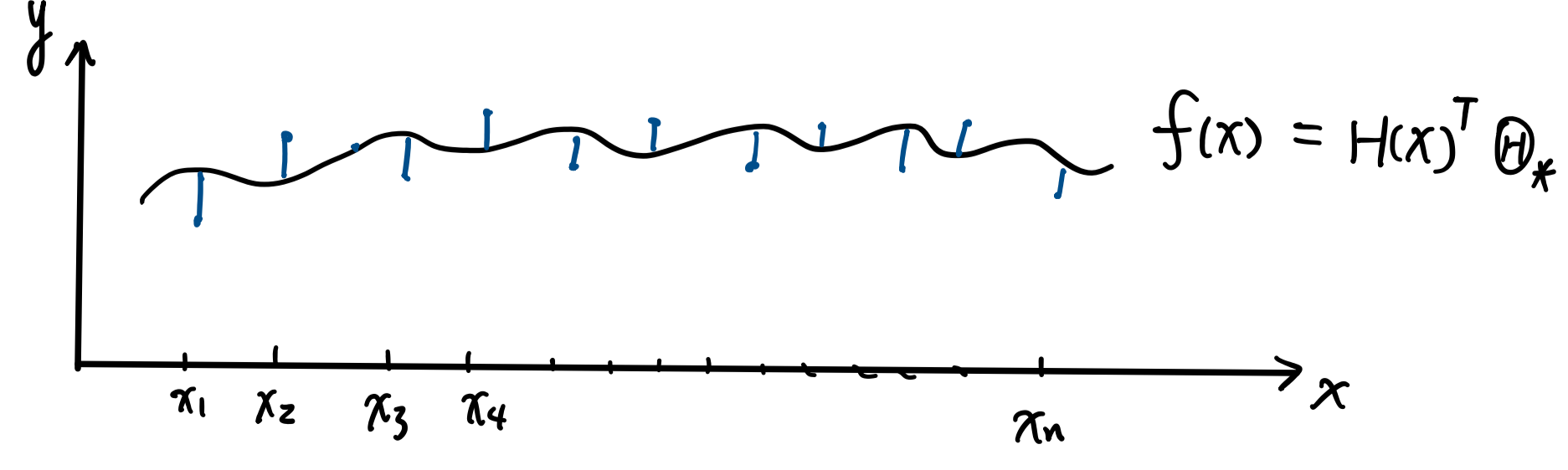
(ii) "Fixed Design": We view input features  $x_i$  as fixed, but outputs  $y_i$  still random

Eg. Coffee Experiment: same barista/machines but perhaps different observed conditions

Focus on fixed design setting:

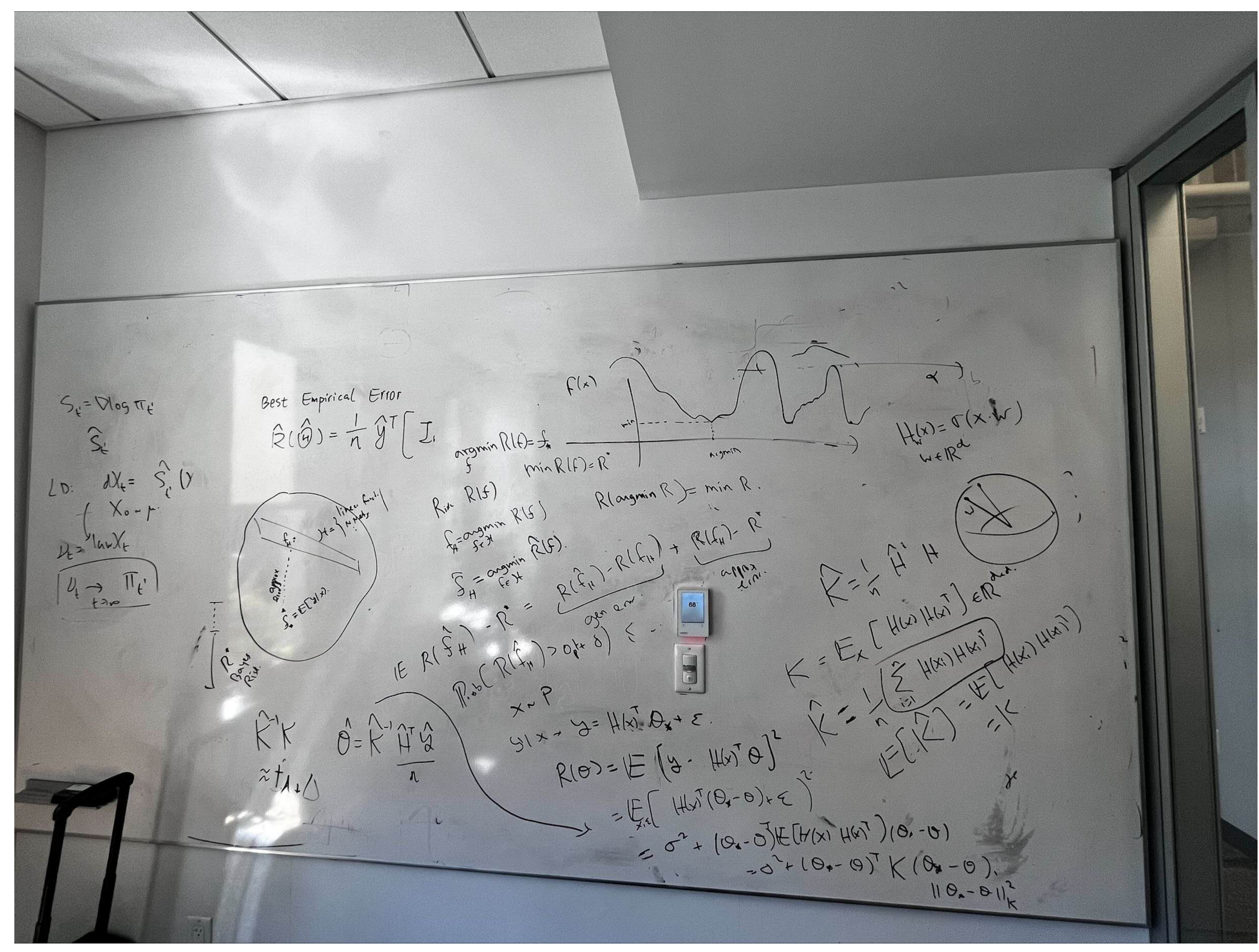
- As before, we assume that  $\hat{H} \in \mathbb{R}^{n \times d}$  has rank  $d$  (hence  $\hat{K}$  is invertible)
- We also suppose that outputs are generated using

$y_i = H(x_i)^T \theta_* + \epsilon_i$  such that  $\begin{cases} \mathbb{E}[\epsilon_i] = 0 \\ \text{Var}(\epsilon_i) = \sigma^2 \text{ for } \forall i=1, \dots, n \\ \epsilon_1, \dots, \epsilon_n \text{ are i.i.d} \end{cases}$   
 $\theta_*$ : deterministic "signal"  
 $\epsilon_i$ : random noise



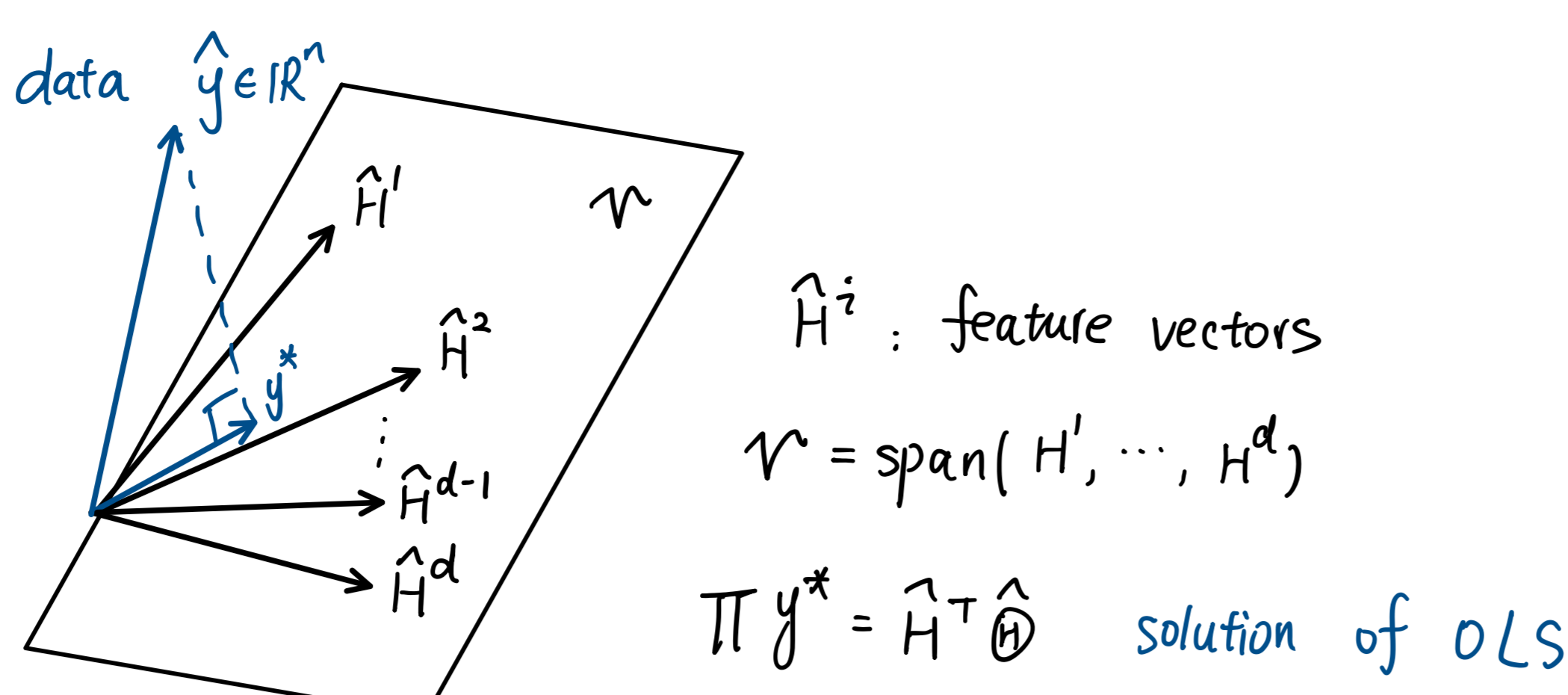
Stronger Assumption:  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$

Then  $y_i | x_i \sim \mathcal{N}(H(x_i)^T \theta_*, \sigma^2), i=1, \dots, n$



FoML Lecture 6: Linear Regression (cont'd): Statistical Analysis

Recap from L5:



Geometric View

Generalization

- \* "Fixed Design": input  $X_i$  fixed,  $y_i$  random
- \* Random Design:  $(X, Y) \sim p$

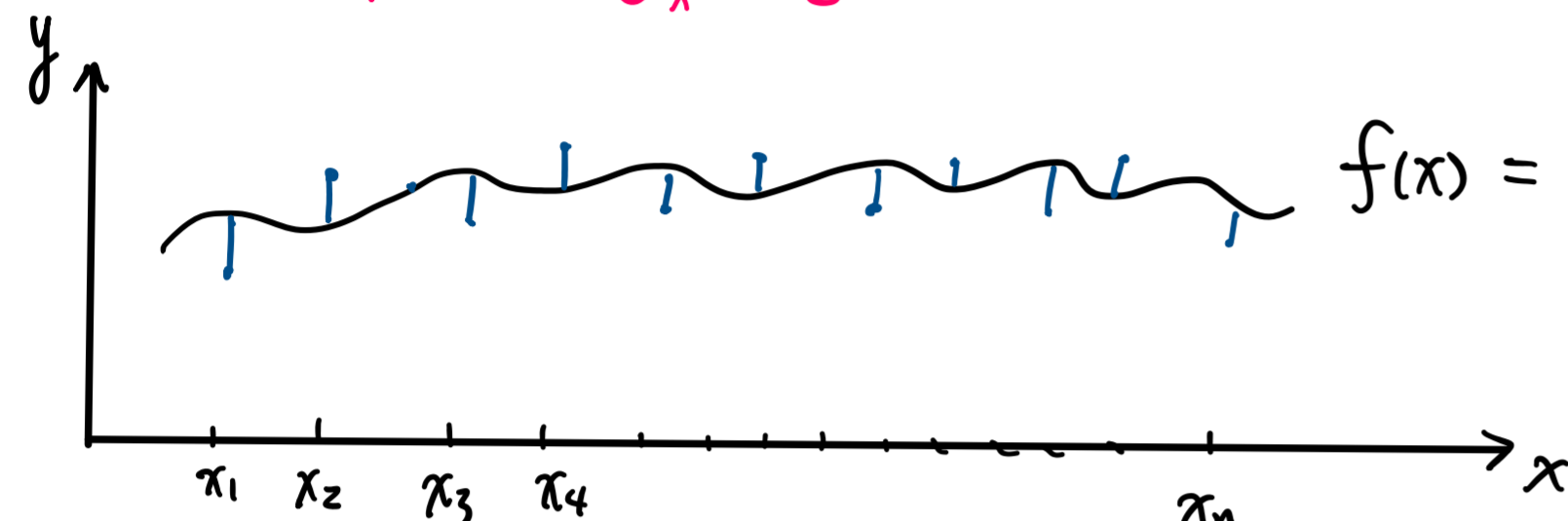
Today (1) Analyze LS on fixed design setting  
 (2) Regularization

Fixed Design:  $x_1, \dots, x_n$  fixed s.t.  $\hat{H} \in \mathbb{R}^{n \times d}$  is rank  $d$

Assumption: labels  $\{y_i\}$  are generated according to

for  $i=1, \dots, n$ ,  $y_i = \underbrace{\hat{H}(x_i)^T \theta_*}_{\text{"signal"}} + \underbrace{\varepsilon_i}_{\text{"noise"}}$ ,  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. R.V.'s with  $\mathbb{E}[\varepsilon] = 0$ ,  $\text{Var}(\varepsilon) = \sigma^2 > 0$

Vector Notation:  $Y = \hat{H}^T \theta_* + \varepsilon$ ,  $\theta_* \in \mathbb{R}^d$  This is what we need to figure out (this is the underlying structure)



One thing ambiguous: How do we determine whether the deviation is because of the lack of parameters or the random noises (Our Assumption)

Stronger Assumption: assume  $\varepsilon_i \sim \mathcal{N}(H(x_i)^T \theta_*, \sigma^2)$  (Gaussian Noise)

Goal: Estimate  $\theta_*$  (God's Parameters)

Risk Decomposition

Def. For any  $\theta \in \mathbb{R}^d$ , the generalization error is  $R(\theta) = \frac{1}{n} \mathbb{E}_Y \|Y - \hat{H}\theta\|^2$  (Note:  $R$  is NOT random)

Prop. Under the linear model assumption,  $R(\theta) = \sigma^2 + (\theta - \theta_*)^T \hat{K} (\theta - \theta_*)$

Rmk.  $R(\theta)$  takes minimum at  $\theta = \theta_*$  of  $\sigma^2$  (unique)

pf.  $R(\theta) = \frac{1}{n} \mathbb{E}_Y \|Y - \hat{H}\theta\|^2$   
 $\stackrel{\text{model assumption}}{=} \frac{1}{n} \mathbb{E}_\varepsilon \|\hat{H}(\theta_* - \theta) + \varepsilon\|^2$   
 $= \frac{1}{n} [\mathbb{E}_\varepsilon \|\varepsilon\|^2 + (\theta_* - \theta)^T \hat{H}^T \hat{H} (\theta_* - \theta) + 2 \mathbb{E}_\varepsilon \langle \varepsilon, \hat{H}(\theta_* - \theta) \rangle]$   
 $= \frac{1}{n} \mathbb{E}_\varepsilon [\sum_{i=1}^n \varepsilon_i^2] + (\theta_* - \theta)^T \hat{K} (\theta_* - \theta) + \underbrace{\langle \mathbb{E}_\varepsilon[\varepsilon], \hat{H}(\theta_* - \theta) \rangle}_0$

#

Next: Evaluate  $\mathbb{E}[R(\hat{\theta})]$  when  $\hat{\theta}$  is the OLS estimator (training data)

Bias-Variance Decomposition

Sps. that  $\hat{\theta}$  is a Random Vector, for example  $\hat{\theta} = \hat{\theta}_{OLS}$ , the OLS estimator

Then  $\mathbb{E}_{\hat{\theta}} [R(\hat{\theta})] = \sigma^2 + \underbrace{\| \mathbb{E}[\hat{\theta}] - \theta_* \|^2}_{\text{Bias}} + \underbrace{\mathbb{E}[\| \hat{\theta} - \mathbb{E}[\hat{\theta}] \|^2]}_{\text{Variance}}$   
 Bayes Risk: Smallest possible error that any method suffers

Def.  $\|x\|_{\hat{K}}^2 = x^T \hat{K} x$  (In particular,  $\|x\|_{I_d}^2 = \|x\|^2$ )

pf. By RD,  $R(\hat{\theta}) = \sigma^2 + (\hat{\theta} - \theta_*)^T \hat{K} (\hat{\theta} - \theta_*)$  (1)

Also,  $\mathbb{E}[(\hat{\theta} - \theta_*)^T \hat{K} (\hat{\theta} - \theta_*)] = \mathbb{E}[(\underbrace{\hat{\theta} - \mathbb{E}[\hat{\theta}]}_{\theta_A} + \underbrace{\mathbb{E}[\hat{\theta}] - \theta_*}_{\theta_B})^T \hat{K} (\underbrace{\hat{\theta} - \mathbb{E}[\hat{\theta}]}_{\theta_A} + \underbrace{\mathbb{E}[\hat{\theta}] - \theta_*}_{\theta_B})]$   
 $= \mathbb{E}[\theta_A^T \hat{K} \theta_A] + \mathbb{E}[\theta_A^T \hat{K} \theta_B] + \mathbb{E}[\theta_B^T \hat{K} \theta_A] + \mathbb{E}[\theta_B^T \hat{K} \theta_B]$   
 $= \mathbb{E}[\| \hat{\theta} - \mathbb{E}[\hat{\theta}] \|^2_{\hat{K}}] + \underbrace{\hat{K} \theta_B \mathbb{E}[\theta_A]}_0 + \underbrace{\theta_B^T \hat{K} \mathbb{E}[\theta_A]}_0 + \mathbb{E}[\theta_B^T \hat{K} \theta_B] = \| \mathbb{E}[\hat{\theta}] - \theta_* \|^2_{\hat{K}}$   
 $= \underbrace{\| \mathbb{E}[\hat{\theta}] - \theta_* \|^2_{\hat{K}}}_{\text{Bias}} + \underbrace{\mathbb{E}[\| \hat{\theta} - \mathbb{E}[\hat{\theta}] \|^2_{\hat{K}}]}_{\text{Variance}} \quad (2)$

Take  $\mathbb{E}(\cdot)$  on both sides of (1) and use (2), we complete the proof.

#

Now let's plug in our OLS estimator:  $\hat{\theta} = \frac{1}{n} \hat{K}^{-1} \hat{H}^T \hat{y}$  where  $\hat{y} = \hat{H} \theta_* + \varepsilon$   
 $\hat{\theta} = \hat{\theta}_{OLS}$   
 $= \hat{K}^{-1} \frac{\hat{H}^T}{n} (\hat{H} \theta_* + \varepsilon)$   
 $= \theta_* + \frac{1}{n} \hat{K}^{-1} \hat{H}^T \varepsilon$

Therefore,  $\mathbb{E}[\hat{\theta}] = \theta_*$  (3)

$\Rightarrow$  Bias Term:  $\| \mathbb{E}[\hat{\theta}] - \theta_* \|^2_{\hat{K}} = 0$ . [OLS is unbiased!]

$\Rightarrow$  Variance Term:  $\mathbb{E}[\| \hat{\theta} - \mathbb{E}[\hat{\theta}] \|^2_{\hat{K}}] = \mathbb{E}[\| \hat{\theta} - \theta_* \|^2_{\hat{K}}] = \sigma^2 \cdot \frac{d}{n}$  (to be cont'd next class)

Q1: Will regularization makes  $\hat{\theta}$  biased?

A1: Yes!

FML Lecture 7 : Linear Regression : Regularisation

Recap from last week :

1) Risk of any LS estimator  $\Theta$  in the fixed design

$$R(\Theta) = \sigma^2 + (\Theta - \Theta_*)^T \hat{K} (\Theta - \Theta_*)$$

2) When  $\hat{\Theta}$  is a random vector (eg.  $\hat{\Theta} = \hat{\Theta}_{OLS}$ ).

$$E[R(\hat{\Theta})] = \sigma^2 + \|E[\hat{\Theta}] - \Theta_*\|_{\hat{K}}^2 + E[\|\hat{\Theta} - E[\hat{\Theta}]\|_{\hat{K}}^2]$$

$$\hat{\Theta}_{OLS} = \hat{K}^{-1} \frac{\hat{H}^T \hat{Y}}{n} = \Theta_* + \hat{K}^{-1} \frac{\hat{H}^T \epsilon}{n} \Rightarrow \hat{\Theta}_{OLS} \text{ is unbiased}$$

Let's compute variance

$$E[\|\hat{\Theta}_{OLS} - \Theta_*\|_{\hat{K}}^2] = E\left[\left(\frac{\epsilon^T \hat{H}}{n}\right) \hat{K}^{-1} \hat{K} \hat{K}^{-1} \left(\frac{\hat{H}^T \epsilon}{n}\right)\right]$$

$$= \frac{1}{n^2} E\left[\epsilon^T \underbrace{\hat{H} \hat{K}^{-1} \hat{H}^T}_{n \Pi} \epsilon\right] = \frac{1}{n} E[\epsilon^T \Pi \epsilon] \text{ for } \epsilon \in \mathbb{R}^n, \Pi \in \mathbb{R}^{n \times n}$$

where  $\Pi = \hat{H}(\hat{H}^T \hat{H})^{-1} \hat{H}^T$  is an orthogonal projection onto  $V = \text{Col}(H)$

$$= \frac{1}{n} E\left[\sum_{i,j} \epsilon_i \epsilon_j \Pi_{ij}\right] = \frac{1}{n} \sum_{i,j} E[\epsilon_i \epsilon_j] \Pi_{ij} = \frac{1}{n} \sum_i \sigma^2 \Pi_{ii} = \frac{\sigma^2}{n} \cdot \text{Tr}(\Pi)$$

$$\begin{cases} \sigma^2, & i=j \\ 0, & i \neq j \end{cases}$$

$$= \frac{\sigma^2}{n} \cdot \text{Tr}(\hat{H}(\hat{H}^T \hat{H})^{-1} \hat{H}^T)$$

$$= \frac{\sigma^2}{n} \cdot \text{Tr}(\underbrace{(\hat{H}^T \hat{H})^{-1}}_{\text{rank } d} \hat{H}^T \hat{H})$$

$$= \frac{\sigma^2}{n} \cdot \text{Tr}(\text{Id}) = \sigma^2 \cdot \frac{d}{n}$$

A, B squared or some generalized condition  
 $\text{Tr}(AB) = \text{Tr}(BA)$

need stronger linear algebra tool

#

Conclusion : Variance  $E[\|\hat{\Theta} - E[\hat{\Theta}]\|_{\hat{K}}^2] = \sigma^2 \cdot \frac{d}{n}$

$$\Rightarrow E[R(\hat{\Theta})] = \sigma^2 + \sigma^2 \cdot \frac{d}{n} \rightarrow \text{excess risk, goes to 0 as } n \nearrow +\infty$$

↳ "Incompressible" error
[variance]  
[No bias]

Q1: What happens when  $n \approx d$ , even  $n < d$ ?

Q2: What happens beyond the fixed design setting?

Regularisation

Motivation :

When  $n = d$ , the normal equations  $\hat{K} \hat{\Theta} = \frac{\hat{H}^T y}{n}$  n equations  
n unknowns

Assuming  $\hat{K}$  invertible, unique solution  $\hat{\Theta}$  with error  $\hat{K}(\hat{\Theta}) = 0$

↳ We are memorizing data (which often includes noise)

When  $n < d$ ,  $\hat{K}$  is not invertible! (under-determined)

Very common regime (gene expression,  $n \ll d$ )

Solution: Add some "friction", "cost" to using features

↳ Explain the data using cheapest option

↳ Occam's Razor

Q<sub>1</sub>: Q: How to define a useful notion of cost?

Q<sub>2</sub>: Def. A **ridge regularisation** considers the cost as  $L^2$ -norm:  $\|\Theta\|_2^2 = \sum_{i=1}^n \Theta_i^2$

Q<sub>3</sub>: Note: An alternative is to consider the  $L^1$  norm  $\|\Theta\|_1 = \sum_i |\Theta_i|$

↳ this leads to **sparse regression**

Question: What is the essential difference of changing  $L^2$ -norm to  $L^1$ -norm

Ridge Regularisation

$$\hat{R}_\lambda(\Theta) = \frac{1}{n} \|\hat{Y} - \hat{H}\Theta\|^2 + \lambda \|\Theta\|_2^2, \lambda \geq 0$$

↑ regularisation strength

→ Minimize  $R_\lambda(\Theta) = \frac{1}{n} \|\hat{Y}\|^2 - 2 \frac{\hat{Y}^T \hat{H} \Theta}{n} + \Theta^T \hat{K} \Theta + \lambda \Theta^T \Theta$

$$\nabla_\Theta \hat{R}_\lambda(\Theta) = 0 \Leftrightarrow -\frac{2}{n} \hat{H}^T \hat{Y} + 2 \hat{K} \Theta + 2 \lambda \Theta = 0$$

$$\Leftrightarrow -\frac{1}{n} \hat{H}^T \hat{Y} + [\hat{K} + \lambda I_n] \Theta = 0$$

$$\Rightarrow \hat{\Theta}_\lambda = [\hat{K} + \lambda I_n]^{-1} \frac{\hat{H}^T \hat{Y}}{n}$$

→ Now we observe that  $[\hat{K} + \lambda I_n]$  is always invertible for  $\lambda > 0$

This is because  $\hat{K}$  is symmetric positive semi-definite, so  $y^T \hat{K} y \geq 0$  for all  $y \in \mathbb{R}^n$

hence  $y^T (\hat{K} + \lambda I_n) y = y^T \hat{K} y + \lambda \|y\|_2^2 > 0$

hence  $\hat{K} + \lambda I_n$  is symmetric positive definite.

→ Compare the generalisation error of  $\hat{\Theta}_\lambda$  in the fixed design setting, linear model  $y_i = H(x_i)^T \Theta_* + \epsilon_i$

$$E[R(\hat{\Theta}_\lambda)] = \sigma^2 + \lambda^2 \Theta_*^T (\hat{K} + \lambda I_n)^{-1} \hat{K} \Theta_* + \frac{\sigma^2}{n} \text{Tr}[\hat{K} (\hat{K} + \lambda I_n)^{-2}]$$

Bias Term ↑
Variance Term ↓

Q: Is the variance term always reduced? How should we choose  $\lambda$ ?

[TA] Solving Normal Equations

- Iterative Method
- SVD
- QR Decomposition
- Cholesky Decomposition

→ Linear Regression

$$X\beta = y$$

$$X \in \mathbb{R}^{n \times d}, y \in \mathbb{R}^n$$

$n$ : # of samples

$d$ : # of features

$n > d$  (otherwise no error)

Goal:  $\hat{\beta} = \arg \min_{\beta} \|X\beta - y\|^2$   
 $= (X^T X)^{-1} X^T y \in \mathbb{R}^d$

$\Leftrightarrow \underbrace{(X^T X) \beta = X^T y}_{\text{Normal Equation}} \in \mathbb{R}^d$

→ Iterative Method

$$\hat{R} = \|X\beta - y\|^2 \in \mathbb{R}$$

$$= (X\beta - y)^T (X\beta - y)$$

$$\nabla_{\beta} \hat{R} = 2X^T (X\beta - y) \in \mathbb{R}^d$$

- Gradient Descent.

for  $\forall t \geq 0$ ,  $\beta_{t+1} \leftarrow \beta_t - \eta \cdot \nabla_{\beta} \hat{R}(\beta_t)$  learning rate  
 $= \beta_t - \eta \cdot 2X^T (X\beta_t - y)$

Advantage:

- ① easy to understand and implement
- ② scalable (Stochastic Gradient Method)

Disadvantage:

- ① might be slow to converge
- Hessian:  $\nabla_{\beta} (\nabla_{\beta} \hat{R}) = 2X^T X \in \mathbb{R}^{d \times d}$   
 and min eigenvalue of  $X^T X \approx 0$  is possible  
 [make  $\eta$  as  $\eta(\lambda_i)$ , Adaptive Gradient Descent]

→ SVD: Singular Value Decomposition

Decompose any matrix  $A = U \Sigma V^T \in \mathbb{R}^{n \times d}$  where  $U \in \mathbb{R}^{n \times n}$ ,  $V \in \mathbb{R}^{d \times d}$  (orthogonal),  $\Sigma \in \mathbb{R}^{n \times d}$  (rectangular) diagonal matrix =  $\begin{bmatrix} * & * & \dots & 0 \\ & * & \dots & \\ & & \dots & \\ 0 & \dots & & 0 \end{bmatrix}$

Change  $X \leftarrow A$

$$(A^T A) x = A^T b$$

$$A = U \Sigma V^T$$

$$\Rightarrow V \Sigma^T \Sigma V^T x = V \Sigma^T U^T b$$

$$\Rightarrow x = V \Sigma^+ U^T b$$

where  $\Sigma^+$ : pseudo-inverse (Search)

Advantage:

- ① Handle ill-conditioned problem:  $\frac{\lambda_{\max}(X^T X)}{\lambda_{\min}(X^T X)}$

Disadvantage:

- ① Expensive computation

QR-decomposition:

Decompose any  $A = QR \in \mathbb{R}^{n \times d}$

$Q$ : orthogonal matrix  $\mathbb{R}^{n \times d}$

$R$ : upper triangle  $\mathbb{R}^{d \times d}$



FoML Lecture 9: Principles of Supervised Learning

Reminder: Office Hours (Joan) tmr. Wed. @ 2 p.m. (612 CDS)

Today: → From fixed to random design

→ Main elements of supervised learning

Limitations of fixed design: We need to predict outside training set

(Linear)

Random Design: Training  $\{\mathcal{X}_i\}_{i=1}^n$  i.i.d. from  $P$

$$y_i = H(\mathcal{X}_i)^T \theta_* + \varepsilon_i, \quad \varepsilon_i \text{ i.i.d.}, \quad \begin{cases} \mathbb{E}[\varepsilon_i] = 0 \\ \mathbb{E}[\varepsilon_i^2] = \sigma^2 \end{cases}$$

But now we evaluate it on a new point  $X \sim P$  (indep. of  $\{\mathcal{X}_i\}_{i=1}^n$ )

• Given  $\theta \in \mathbb{R}^d$ ,  $R(\theta) = \mathbb{E}_P[(H(\mathcal{X})^T \theta - y)^2]$

$$= \mathbb{E}_{\mathcal{X}, \varepsilon}[(H(\mathcal{X})^T (\theta - \theta_*) - \varepsilon)^2]$$

↳ main difference w.r.t. fixed design

$$= \sigma^2 + \mathbb{E}_{\mathcal{X}}[(\theta - \theta_*)^T H(\mathcal{X}) H(\mathcal{X})^T (\theta - \theta_*)]$$

$$= \sigma^2 + (\theta - \theta_*)^T \underbrace{\mathbb{E}_{\mathcal{X}}[H(\mathcal{X}) H(\mathcal{X})^T]}_{K \in \mathbb{R}^{d \times d}} (\theta - \theta_*) \quad , \quad \theta \in \mathbb{R}^d$$

$$= \sigma^2 + \|\theta - \theta_*\|_K^2$$

→ The only difference w.r.t. fixed design is that we have  $K$  instead of  $\hat{K} = \frac{1}{n} \hat{K} \hat{H} = \frac{1}{n} \sum_{i=1}^n H(\mathcal{X}_i) H(\mathcal{X}_i)^T$ ,

$$K = \int H(\mathcal{X}) H(\mathcal{X})^T P(\mathcal{X}) d\mathcal{X} \quad (\text{continuous version of } \hat{K})$$

→ Now we view  $\hat{K}$  as the sample version of  $K$

Let's now plug  $\theta = \hat{\theta}_{ols} = \hat{K}^{-1} \frac{\hat{H}^T \mathbf{y}}{n} = \theta_* + \hat{K}^{-1} \frac{\hat{H}^T \varepsilon}{n}$

$$\mathbb{E}_{\mathcal{X}, \varepsilon} R(\hat{\theta}_{ols}) = \sigma^2 + \mathbb{E} \left[ \frac{1}{n^2} \varepsilon^T \hat{H} \hat{K}^{-1} K \hat{K}^{-1} \hat{H}^T \varepsilon \right]$$

$$= \sigma^2 + \frac{1}{n^2} \mathbb{E} \left[ \text{Tr}(\varepsilon^T \hat{H} \hat{K}^{-1} K \hat{K}^{-1} \hat{H}^T \varepsilon) \right]$$

$$\text{Tr}(\varepsilon \varepsilon^T \hat{H} \hat{K}^{-1} K \hat{K}^{-1} \hat{H}^T)$$

Rmk. We can interchange  $\mathbb{E}[\cdot]$  &  $\text{Tr}(\cdot)$  as they are both linear.

$$= \sigma^2 + \frac{1}{n^2} \mathbb{E} \left[ \underbrace{\text{Tr}(\varepsilon \varepsilon^T)}_{\sigma^2} \cdot \mathbb{E} \left[ \underbrace{\text{Tr}(\hat{H} \hat{K}^{-1} K \hat{K}^{-1} \hat{H}^T)}_{\text{Tr}(\hat{H}^T \hat{H} \hat{K}^{-1} K \hat{K}^{-1})} \right] \right]$$

$$= \sigma^2 + \frac{\sigma^2}{n} \mathbb{E} \left[ \text{Tr}(K \hat{K}^{-1}) \right]$$

→ We need to understand the inverse of a random matrix  $\hat{K}$

This requires tools from Random Matrix Theory

Q: Does regularisation still make sense in random design setting

A: Yes, as  $\hat{K}$  might be poorly-conditioned and we should regularize on that.

Main Results we have seen:

- $\hat{\theta}_{ols}$  "best" model that fits data
- How to assess model outside training

Now let's describe general picture

→ Model for data: Training data  $n$  i.i.d. samples

$$\left\{ \begin{matrix} \mathcal{X}_i \\ \mathcal{Y}_i \end{matrix} \right\}_{i=1, \dots, n} \quad \begin{matrix} \uparrow & \uparrow \\ \mathcal{X} & \mathcal{Y} \end{matrix} \quad \text{drawn from unknown distribution } P \text{ in } \mathcal{X} \times \mathcal{Y}$$

Strong Assumption: Same distribution as training data

$$\text{Test data } (\mathcal{X}, \mathcal{Y}) \sim \hat{P} \text{ independent of training}$$

(when test distribution  $P_{test} \neq P_{train}$ , we have different problem: transfer learning)

→ Loss function: a function  $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

that measures agreement between true and predicted label

Ex:  $l(y, y') = (y - y')^2$  (LS regression)

$l(y, y') = \mathbb{1}_{y \neq y'}$  (Classification)

→ Now, given any mapping  $f: \mathcal{X} \rightarrow \mathcal{Y}$ ,

its (expected future) performance is:

$$R(f) = \mathbb{E}_{(\mathcal{X}, \mathcal{Y}) \sim P} [l(f(\mathcal{X}), \mathcal{Y})] \quad (\text{population risk / generalization error})$$

Ex.  $R(\theta) = \mathbb{E}[(H(\mathcal{X})^T \theta - y)^2]$  (least square in random setting)

Rmk. Now everything is on a random design

→ From population risk, we can define the optimum predictor:

$$f_* = \underset{f: \mathcal{X} \rightarrow \mathcal{Y}}{\text{argmin}} R(f)$$

↳ Recall in LS setting,  $f_*(x) = \mathbb{E}_P[Y|X=x]$  (lec. 3)

↳  $f_*$  is the **Bayes Predictor**,  $R^* = R(f_*)$  is called **Bayes Risk / Bayes Risk**

↳ We can have  $R^* > 0$  in general ( $R^* = \sigma^2$  in the linear model)

**Bayes Risk is unattainable in general**

→ At least two reasons:

(1) It requires knowledge of data distribution  $P$ !

(2)  $f_*$  might be arbitrarily crazy function → hard to even approximate!

FML Lecture 10: Elements of SL

Recall:  $f: \mathcal{X} \rightarrow \mathcal{Y}$   
input output

$R(f) = \mathbb{E}_{(X,Y) \sim P} [\ell(f(X), Y)]$  Population Risk

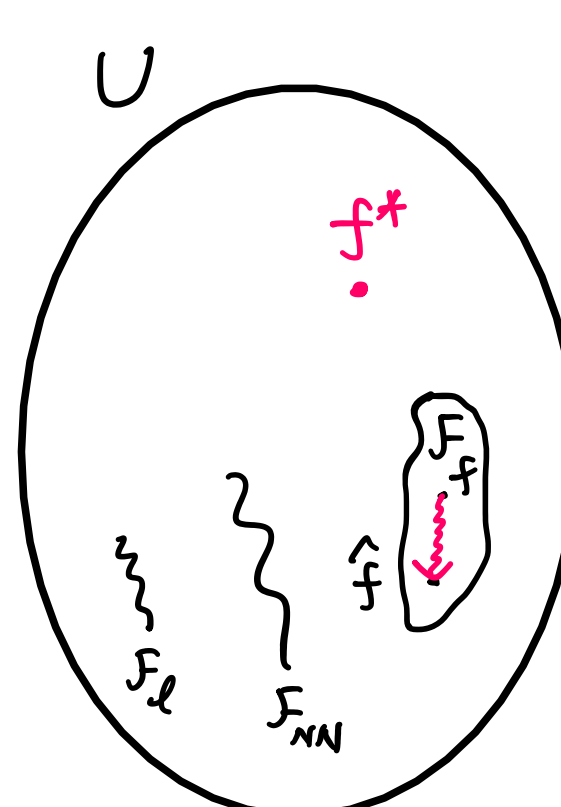
$\hookrightarrow (f^*, R^*)$  Bayes Predictor/Risk  
 $\text{argmin}_R \min R$

$\rightarrow$  Unpractical  $\begin{cases} \cdot \text{They depend on population} \\ \cdot \text{They can be arbitrarily complex} \end{cases}$

Instead in SL, we focus our attention on a hypothesis class  $\mathcal{F} = \{f_\theta: \mathcal{X} \rightarrow \mathcal{Y}, \theta \in \Theta\}$

Ex.  $\mathcal{F}_2 = \{f_\theta(x) = H(x)^T \theta, \theta \in \Theta = \mathbb{R}^d, x \in \mathcal{X} = \mathbb{R}^d\} \subset U = \{f: \mathbb{R}^d \rightarrow \mathbb{R}\}$  (linear hypothesis class),  $\dim \mathcal{F}_2 = d$

$\mathcal{F}_{NN} = \{f_\theta(x) = \underbrace{g_L(w_L, \dots, g_1(w_1, x))}_L, \theta \in \Theta = \{w_1, \dots, w_L\}\} \subset U$   $\dim \mathcal{F}_{NN} = L$   
L layers



$\rightarrow$  Now we can consider the best predictor in  $\mathcal{F}$

$\bar{f} = \text{argmin}_{f \in \mathcal{F}} R(f)$

$\rightarrow \inf_{f \in \mathcal{F}} R(f) - R^* \geq 0$  measures how accurate the hypothesis space is for our prediction task  
Approximation Error/Risk

$\rightarrow$  It is still impossible to find  $\bar{f}$  ( $P$  unknown)

$\rightarrow$  Instead, we can consider minimizing the Empirical Risk.

$\hat{R} = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$  where  $(x_i, y_i) \sim_{i.i.d.} P$

Since  $\{(x_i, y_i)\}_{i=1}^n$  is a Random Sample,  $\hat{R}$  is a Random Functional

Q1: What is the mean of  $\hat{R}(f)$  for any  $f$ ?

A:  $\mathbb{E}[\hat{R}(f)] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)\right] \stackrel{(x_i, y_i) \sim_{i.i.d.} P}{=} \mathbb{E}_P[\ell(f(x_1), y_1)] = R(f)$

i.e.,  $\hat{R}$  is an unbiased estimator of  $R$

$\rightarrow \{\mathcal{Z}_i = \ell(f(x_i), y_i)\}_{i=1}^n$  are i.i.d. R.V.'s

$\mathcal{Z} = \hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \mathcal{Z}_i$ , where  $\mathbb{E}\mathcal{Z} = R(f)$ ,  $\text{Var}(\mathcal{Z}) = \sigma_f^2$  (n large)

under mild moment assumptions, the sample mean is asymptotically normal

[CLT]:  $\sqrt{\frac{n}{\sigma_f^2}} (\hat{R}(f) - R(f)) \xrightarrow{d} \mathcal{N}(0, 1)$   
Prob. of being larger than  $\frac{\epsilon_f}{\sqrt{n}}$  is exponentially smaller

$\rightarrow$  for large  $n$  and fixed  $f$ ,  $|\hat{R}(f) - R(f)| \stackrel{\ominus}{\asymp} \frac{\sigma_f}{\sqrt{n}}$

$\hookrightarrow$  can be formalized in the non-asymptotic setting using Concentration Inequality (finite  $n$ )

We can define

$\hat{f} = \text{argmin}_{f \in \mathcal{F}} \hat{R}(f)$  Empirical Risk Minimization (ERM)

$\hookrightarrow$  Look for hypothesis in our class that best fits the training data

$\hookrightarrow$  Now, we have reduced learning to solving an optimization problem

Q: How to control the quality of ERM?

i.e., control generalization gap  $R(\hat{f}) - R^*$

For any  $f$ ,  $R(f) = \hat{R}(f) + (R(f) - \hat{R}(f))$  [tautology]

So, if we want LHS to be small, we can "hope" to have:

$\begin{cases} \hat{R}(f) \text{ small, and} \\ R(f) - \hat{R}(f) \text{ also small} \end{cases}$

ERM is designed to minimize  $\hat{R}(f)$ , then what about  $R(f) - \hat{R}(f)$

Key Observation: there is an inherent tension between the two terms

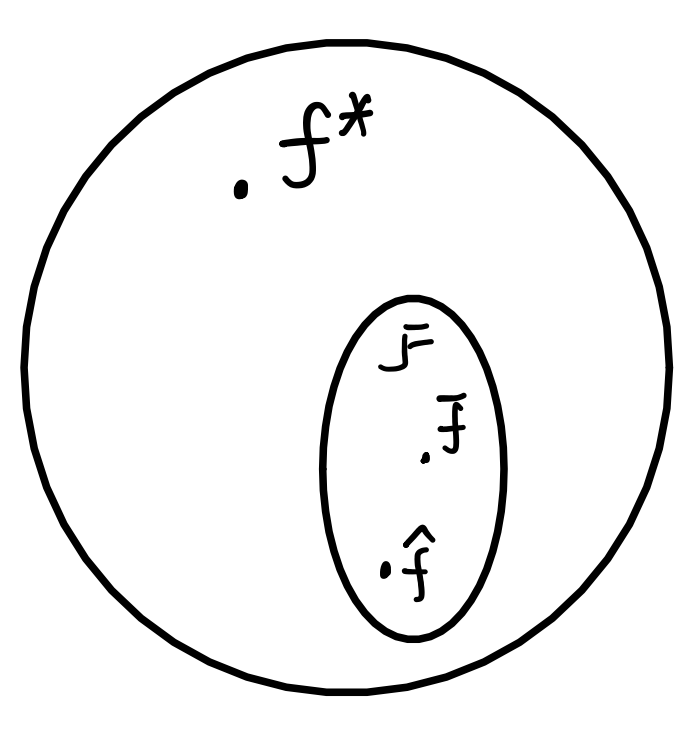
$\hat{f} = \text{ERM}$

$\hat{R}(\hat{f})$ : decreases as  $\mathcal{F}$  gets bigger  
but

$R(\hat{f}) - \hat{R}(\hat{f})$  might increase as  $\mathcal{F}$  gets bigger

Decomposition of Risk:

Consider  $\hat{f} = \text{argmin}_{f \in \mathcal{F}} \hat{R}(f)$  (ERM)



$R(\hat{f}) - R^* = \underbrace{R(\hat{f}) - \inf_{f \in \mathcal{F}} R(f)}_{\text{estimation error}} + \underbrace{\inf_{f \in \mathcal{F}} R(f) - R^*}_{\epsilon_A = \text{approximation error}}$

$= R(\hat{f}) - \hat{R}(\hat{f}) + \hat{R}(\hat{f}) - R(\bar{f}) + \epsilon_A$ , where  $\bar{f} = \text{argmin}_{f \in \mathcal{F}} R(f)$

$\leq R(\hat{f}) - \hat{R}(\hat{f}) + \hat{R}(\bar{f}) - R(\bar{f}) + \epsilon_A$ , as  $\hat{R}(\hat{f}) \leq \hat{R}(\bar{f})$  by def. of  $\hat{f}$

$\leq |R(\hat{f}) - \hat{R}(\hat{f})| + |R(\bar{f}) - \hat{R}(\bar{f})| + \epsilon_A$ , by triangular inequality

$\leq 2 \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| + \epsilon_A$

FoML Lecture 11: Decomposition of Risk

Recap:  $R(f)$ : Expected Risk :=  $\mathbb{E}_p[l(f(x), y)]$

$\hat{R}(f)$ : Empirical Risk :=  $\frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i)$

•  $\mathbb{E}[\hat{R}(f)] = R(f)$  ( $\hat{R}$  is an unbiased estimator of  $R$ )

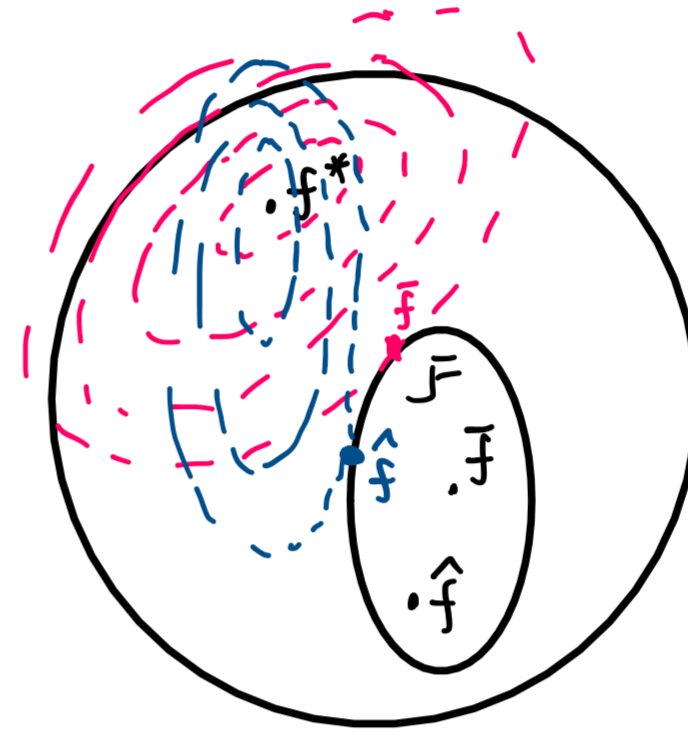
•  $|\hat{R}(f) - R(f)| \sim \frac{\sigma_f}{\sqrt{n}}$  |  $\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{R}(f)$   
 hypothesis class Empirical Risk Minimization

• ML "Tautology": for  $\forall f$ ,  $R(f) = \underbrace{\hat{R}(f)}_{\text{"under control"}} + \underbrace{(R(f) - \hat{R}(f))}_{\text{estimation}}$

§ Decomposition of Risk

Consider  $\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{R}(f)$

Goal: Control excess risk:  $R(\hat{f}) - R^*$  ↗ Bayes Risk



$$R(\hat{f}) - R^* = \underbrace{R(\hat{f}) - \inf_{f \in \mathcal{F}} R(f)}_{\text{estimation error}} + \underbrace{\inf_{f \in \mathcal{F}} R(f) - R^*}_{\varepsilon_A = \text{approximation error}}$$

$$= R(\hat{f}) - \hat{R}(\hat{f}) + \hat{R}(\hat{f}) - R(\bar{f}) + \varepsilon_A, \text{ where } \bar{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f)$$

$$\leq R(\hat{f}) - \hat{R}(\hat{f}) + \hat{R}(\bar{f}) - R(\bar{f}) + \varepsilon_A, \text{ as } \hat{R}(\hat{f}) \leq \hat{R}(\bar{f}) \text{ by def. of } \hat{f}$$

$$\leq |R(\hat{f}) - \hat{R}(\hat{f})| + |R(\bar{f}) - \hat{R}(\bar{f})| + \varepsilon_A, \text{ by triangular inequality}$$

$$\leq 2 \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| + \varepsilon_A \text{ where } R(f) = \mathbb{E}[\hat{R}(f)]$$

$\varepsilon_S = \text{Statistical error}$

"Rule of Thumb":

→ "Small" hypothesis space  $\mathcal{F}$ :  $\varepsilon_A$  dominates over  $\varepsilon_S$

→ "Large" hypothesis space  $\mathcal{F}$ :  $\varepsilon_S$  dominates over  $\varepsilon_A$

→ Instance of  $\frac{\text{bias}}{\varepsilon_A} - \frac{\text{variance}}{\varepsilon_S}$  decomposition of risk

Important Remark:

$\varepsilon_S$  is an upper bound of the estimation error

→ Upper Bound is pessimistic

Q: How does  $\varepsilon_S$  behave as a function of "size" of  $\mathcal{F}$  and size of training set  $n$ ?

$$\varepsilon_S = \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| \quad (\text{Uniform})$$

Recall that before, we measured fluctuations at an  $f \in \mathcal{F}$ :

$$|\hat{R}(f) - R(f)| \simeq \frac{\sigma_f}{\sqrt{n}} \quad (\text{Pointwise})$$

\* To get the main idea, consider idealized setting

(i)  $\mathcal{F} = \{f_1, \dots, f_M\}$  is a finite set of  $M$  hypothesis

(ii)  $\hat{R}(f_i)$  are indep. Gaussian R.V.'s with mean  $R(f_i)$  and variance  $\sigma^2$

$$\max_{i=1, \dots, M} \hat{R}(f_i) - R(f_i)$$

Then  $Z_i = \hat{R}(f_i) - R(f_i) \sim \mathcal{N}(0, \sigma^2)$ , i.i.d.,

$$\text{Now } \mathbb{E} \max_i Z_i \sim \sqrt{2\sigma^2 \log M}$$

FML Lecture 12 : Statistical Error in SL

Recall:  $\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f)$  **ERM**

$$R(\hat{f}) - R^* \leq 2\epsilon_s + \epsilon_A \quad \text{with} \quad \begin{aligned} \epsilon_A &= \inf_{f \in \mathcal{F}} R(f) - R^* \\ \epsilon_s &= \sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \end{aligned}$$

→ Natural Tension/Trade-off between approximation & statistical error

$\epsilon_A \downarrow$  as  $\mathcal{F}$  grows, while  $\epsilon_s \uparrow$  as  $\mathcal{F}$  grows

→ To understand  $\epsilon_s$ , we need to move from pointwise bound  $|\hat{R}(f) - R(f)| \sim \sqrt{\frac{\epsilon^2}{n}}$  to uniform bound

→ Simplified Settings:

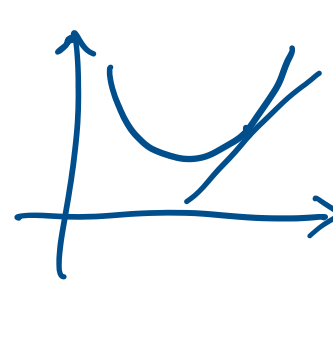
(1)  $\mathcal{F} = \{f_1, \dots, f_M\}$  finite discrete hypothesis class

(2)  $\hat{R}(f_i) \sim \mathcal{N}(R(f_i), \epsilon^2)$ ,  $i=1, \dots, M$

Q:  $\mathbb{E} \max_{i=1, \dots, M} (\hat{R}(f_i) - R(f_i))$  ?

Tools:

(a) Jensen's Inequality:  $f$  is convex, then  $f(\mathbb{E}X) \leq \mathbb{E}f(X)$

Convexity:  for  $\forall x, y$ ,  $f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle$  [Linear approximation of  $f$  at  $x$ ]

Then apply  $x$  by  $\mathbb{E}X$  to have:  $f(\mathbb{E}X) \leq f(y) - \langle \nabla f(\mathbb{E}X), y - \mathbb{E}X \rangle$

Taking  $\mathbb{E}[\cdot]$  on both sides,  $f(\mathbb{E}X) \leq \mathbb{E}f(y)$  then choose  $y$  to be  $X$

(b) Moment Generating Function

$$t \mapsto \mathbb{E}[e^{tX}] = e^{\frac{t^2 \epsilon^2}{2}} \text{ for } X \sim \mathcal{N}(0, \epsilon^2)$$

Let  $Z_i = \hat{R}(f_i) - R(f_i)$ , so  $Z_i \sim \mathcal{N}(0, \epsilon^2)$

$$\bar{Z} = \max(Z_1, \dots, Z_M)$$

We want  $\mathbb{E}\bar{Z}$ :

$$\begin{aligned} \text{Let } t > 0, \quad \mathbb{E}[\bar{Z}] &= \mathbb{E}\left[\frac{1}{t} \log(e^{t\bar{Z}})\right] \stackrel{\text{Jensen's (concave)}}{\leq} \frac{1}{t} \log(\mathbb{E}[e^{t\bar{Z}}]) = \frac{1}{t} \log\left(\mathbb{E}\left[\max_{i=1, \dots, M} e^{tZ_i}\right]\right) \\ &\leq \frac{1}{t} \log\left(\mathbb{E}\left[\sum_{i=1}^M e^{tZ_i}\right]\right) \\ &= \frac{1}{t} \log\left(\sum_{i=1}^M \mathbb{E}[e^{tZ_i}]\right) \\ &\stackrel{\text{MGF}}{=} \frac{1}{t} \log\left(M e^{\frac{t^2 \epsilon^2}{2}}\right) \\ &= \frac{\log M}{t} + \frac{t \epsilon^2}{2} \triangleq \phi(t) \end{aligned}$$

Recap:  $\mathbb{E}\bar{Z} \leq \phi(t)$  for all  $t > 0$

$$\text{So } \mathbb{E}\bar{Z} \leq \inf_{t > 0} \phi(t) = 2 \sqrt{\frac{\log M \epsilon^2}{2}} = \sqrt{2 \epsilon^2 \log M}$$

→ Price to pay (at most) for uniform deviations  $\sqrt{2 \log M}$

→ In fact, we can show (much harder) a **lower bound** of the form

$c \cdot \sqrt{\log n} \cdot \epsilon^2$  when  $Z_1, \dots, Z_M$  are i.i.d.

→  $\hat{R}(f) - R(f) \sim \mathcal{N}(0, \frac{\epsilon^2}{n})$  therefore  $\mathbb{E}[\max(\hat{R}(f) - R(f))] \leq \sqrt{\frac{2 \epsilon^2 \log M}{n}}$

→ Key Step: Replace the max by the sum **Union bound**

→ Q: How about infinite hypothesis Space?

e.g.  $\mathcal{F} = \{f(x) = H(x)^T \theta, \theta \in \mathbb{R}^d\}$

Intuition:



$$\text{Recap: } R(\hat{f}) - R^* \leq \underbrace{\inf_{f \in \mathcal{F}} R(f) - R^*}_{\text{approx.}} + 2 \underbrace{\sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)|}_{\text{Statistical. } \approx \sqrt{\frac{2 \log |\mathcal{F}|}{n}} \text{ (up to a constant)}}$$

Question:

- 1) When is the upper bound correctly capturing the tradeoff between approximation & estimation
- 2) Estimate the quantities in practice
- 3) How to efficiently adjust the "size" of  $\mathcal{F}$  to balance errors?

Answer:

(1) Upper bounds will be generally pessimistic

Exceptions: Sometimes we can directly analyze the generalization gap:  $R(\hat{f}) - R(\hat{f}) = \frac{\epsilon^2 d}{n}$  (in the fixed design of OLS)

(2) Cross-Validation

In practice, we split the available data into two buckets

$$\text{Training Set: } T = \{(x_i, y_i)\}_{i=1, \dots, n}$$

$$\text{Validation Set: } V = \{(x'_i, y'_i)\}_{i=1, \dots, m}$$

$$\text{ERM (using } T) \hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f)$$

Goal: Estimate  $R(\hat{f}) - R^*$

→ Recall that for each fixed  $f$ ,  $\hat{R}(f)$  is an unbiased estimator of  $R(f)$

Why isn't  $\hat{R}(\hat{f})$  a good estimator of  $R(\hat{f})$ ?

Because  $\hat{f}$  depends on randomness in  $T$ , can not treat  $\hat{f}$  as fixed

• Define another estimator  $\tilde{R}(f) = \frac{1}{m} \sum_{j=1}^m \ell(f(x'_j), y'_j)$

Still have that  $\mathbb{E}_V \tilde{R} = R$  and  $\tilde{R}(\hat{f})$  is an unbiased estimator of  $R(\hat{f})$

Next Class:  $|\tilde{R}(\hat{f}) - R(\hat{f})| \sim \sqrt{\frac{1}{m}}$  → the size of validation set.

FoML Lecture 13: Universal Approximation

Next Tuesday's Class: Florentin Gath Guest

Office Hours will be moved to Thursday.

In the past lectures, we saw that excess risk of ERM:

$$R(\hat{f}) - R^* \geq \mathcal{E}_A = \min_{f \in \mathcal{F}} R(f) - R^*$$

(approximation error)

Q: Design hypothesis class  $\mathcal{F}$  s.t.  $\mathcal{E}_A$  is as small as we want?

→ Let  $\mathcal{U} = \{f: \mathcal{X} \rightarrow \mathbb{R}, f \text{ is continuous}\}$

↳ Assume that Bayes estimator  $f^* \in \mathcal{U}$

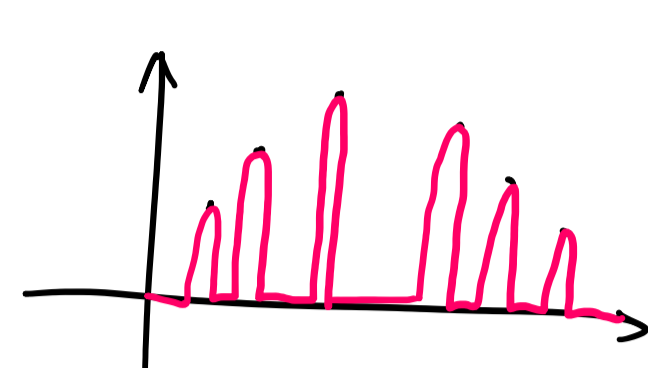
↳ We now consider a norm in  $\mathcal{U}$  given the supremum of  $f \in \mathcal{U}$ ,  $\|f\| = \sup_{x \in \mathcal{X}} |f(x)|$

↳ Can we do ERM on  $\mathcal{U}$  directly?

i.e., Given  $\{(x_i, y_i)\}_{i=1}^n$ ,  $\min_{f \in \mathcal{U}} \hat{R}(f) = \frac{1}{n} \sum_{i=1}^n |y_i - f(x_i)|^2$

No! Because there is no control of statistical error on  $\mathcal{U}$ !

$\sup_{f \in \mathcal{U}} |R(f) - \hat{R}(f)| = \|f^*\|$  for any  $n$ !



→ So we need to somehow "simplify" the universe.

Regularization Perspective

Consider a set  $A \subseteq \mathcal{U}$ , e.g.  $A = \{f: [0,1]^d \rightarrow \mathbb{R} \text{ polynomial}\}$

$A = \{f: [0,1]^d \rightarrow \mathbb{R}, f(x) = W_L \sigma(W_{L-1} \sigma(W_{L-2} \dots \sigma(W_1 x))\}$  Neural Nets of depth  $L$

$A = \{f(x) = H(x)^T \Theta, \Theta \in \mathbb{R}^d\}$  Linear Regression

→ Now we consider a "cost" measure over  $A$

$\gamma: A \rightarrow \mathbb{R}$ , measuring how expensive is it to use a given  $f \in A$

ex.  $A = \{\text{polynomials}\}$ ,  $\gamma(p) = \text{degree of } p$

$A = \{f \text{ Neural Networks}\}$ ,  $\gamma(f) = \# \text{ of parameters}$

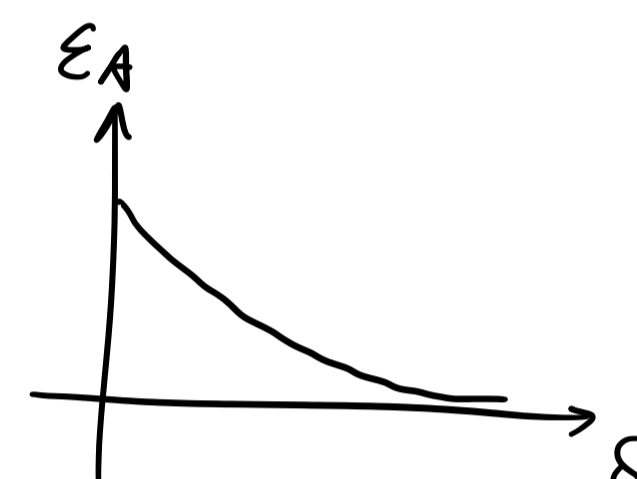
$A = \{f(x) = H(x)^T \Theta, \Theta \in \mathbb{R}^d\}$ ,  $\gamma(f) = \|\Theta\|$

→ We can now use  $\gamma$  to design a hypothesis class:

for each  $\delta > 0$ ,  $\mathcal{F}_\delta = \{f \in A: \gamma(f) \leq \delta\}$

→ Q: How does  $\mathcal{E}_A$  behave as we increase  $\delta$ ?

When can we have  $\mathcal{E}_A \rightarrow 0$  as  $\delta \rightarrow \infty$



→ a set  $A \subseteq \mathcal{U}$  is **dense** if for  $\forall f \in \mathcal{U}$  and any  $\epsilon > 0$ , there exists  $g \in A$  s.t.  $\|f - g\| \leq \epsilon$ .

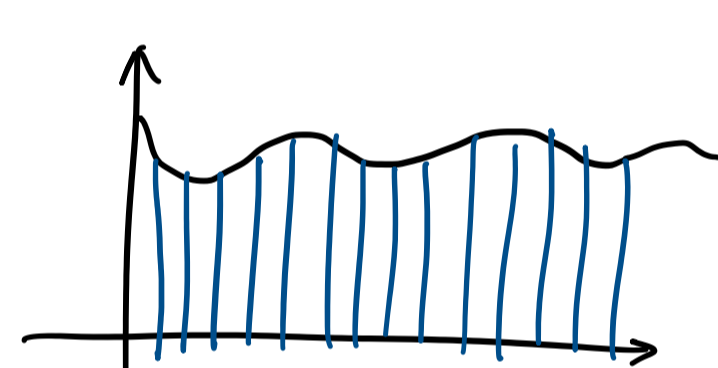
→ When a set  $A \subseteq \mathcal{U} = C(\mathcal{X})$  is dense, we say that it has the **universal approximation property**

§ Universal Approximation of Polynomials

Consider a continuous function  $f: [0,1] \rightarrow \mathbb{R}$

For each  $m$ , we consider the polynomial (deg.  $m$ )

$$p_m(x) = \sum_{j=0}^m f\left(\frac{j}{m}\right) \binom{m}{j} x^j (1-x)^{m-j}$$



Theorem (Weierstrass, early 20th century)

$$\lim_{m \rightarrow \infty} \|f - p_m\| = \lim_{m \rightarrow \infty} \sup_{x \in [0,1]} |f(x) - p_m(x)| = 0$$

So polynomials have the the universal approximation property

Pf. (By Bernstein)

Fix  $x \in [0,1]$ , Let  $Z_1, \dots, Z_m$  be i.i.d. Bernoulli R.V. of parameter  $x$

Let  $W = \frac{1}{m} \sum_{j=1}^m Z_j$  [ $mW \sim \text{Binomial}(m, x)$ ]

We have (i)  $1 = \sum_{j=0}^m \mathbb{P}(W = \frac{j}{m}) = \sum_{j=0}^m \binom{m}{j} x^j (1-x)^{m-j}$

(ii)  $\mathbb{E}W = \mathbb{E}Z_i = x = \sum_{j=0}^m \binom{m}{j} \frac{j}{m} x^j (1-x)^{m-j}$

(iii)  $\text{Var}(W) = \frac{\text{Var}(Z_1)}{m} = \frac{x(1-x)}{m} = \sum_{j=0}^m \binom{m}{j} \left(\frac{j}{m} - x\right)^2 x^j (1-x)^{m-j}$

$$\begin{aligned} p_m(x) - f(x) &= \sum_{j=0}^m \binom{m}{j} f\left(\frac{j}{m}\right) x^j (1-x)^{m-j} - f(x) \cdot 1 \\ &= \sum_{j=0}^m \binom{m}{j} \left[f\left(\frac{j}{m}\right) - f(x)\right] x^j (1-x)^{m-j} \quad (*) \end{aligned}$$

As  $f \in C[0,1]$ , we have

① for  $\forall \epsilon > 0$ ,  $\exists \delta > 0$  s.t.  $|f(x) - f(y)| \leq \epsilon$  whenever  $|x - y| \leq \delta$  ( $\delta$  is indep. of  $\epsilon$ , uniform continuity)

②  $f$  is bounded:  $\|f\| = \sup_{x \in [0,1]} |f(x)| = M < \infty$

We then break (\*) into two parts:

$$(*) = \sum_{j, |\frac{j}{m} - x| \leq \delta} \underbrace{|f\left(\frac{j}{m}\right) - f(x)|}_{b_{m,j}(x)} \binom{m}{j} x^j (1-x)^{m-j} + \sum_{j, |\frac{j}{m} - x| > \delta} |f\left(\frac{j}{m}\right) - f(x)| b_{m,j}(x)$$

$$\leq \epsilon \cdot \sum_{|\frac{j}{m} - x| \leq \delta} b_{m,j}(x) + 2M \cdot \sum_{|\frac{j}{m} - x| > \delta} b_{m,j}(x)$$

$$= \epsilon \cdot \mathbb{P}(|W - \mathbb{E}W| \leq \delta) + 2M \cdot \underbrace{\mathbb{P}(|W - \mathbb{E}W| > \delta)}_{\leq \frac{\text{Var}(W)}{\delta^2} = \frac{x(1-x)}{m\delta^2}} \quad \text{Chebyshev}$$

$$\leq \epsilon + 2M \cdot \frac{x(1-x)}{m\delta^2} \leq \epsilon + \frac{M}{2m\delta^2}$$

Setting  $m = \frac{M}{2\epsilon\delta^2}$  to get:

$$\sup_{x \in [0,1]} |p_m(x) - f(x)| \leq 2\epsilon \quad \text{for } \forall \epsilon > 0.$$

Therefore  $\lim_{m \rightarrow \infty} \|p_m - f\| = 0$ .

#

Remark:

→ The polynomial we have used here  $b_{m,j}(x) = \binom{m}{j} x^j (1-x)^{m-j}$

are called **Bernstein Polynomial**

→ They are not optimal, in the sense of having smallest degree

$m$  for a target error  $\epsilon$ . (Optimal approximation in the uniform

norm is obtained by **Chebyshev Polynomials**)

FML Lecture 15: The Curse of Dimensionality

Recap: Excess Risk:  $R(\hat{f}) - R^* = \overset{\text{(approx.)}}{\mathcal{E}_A} + \overset{\text{(estimation)}}{\mathcal{E}_S}$   
 $\underset{f \in \mathcal{F}}{\inf} R(f) - R^* \quad R(\hat{f}) - \underset{f \in \mathcal{F}}{\inf} R(f)$

↳ Two parameters guiding this error

→  $n$ : # of training points

→  $\delta$ : "size" of hypothesis space  $\mathcal{F}_\delta = \{f \in \mathcal{A}, r(f) \leq \delta\}$   
 $f: \mathcal{X} \rightarrow \mathbb{R}, d \triangleq \dim(\mathcal{X})$

We saw:

(i)  $\mathcal{E}_A \rightarrow 0$  as  $\delta \rightarrow \infty$  (Universal Approximation)

(ii)  $\mathcal{E}_S \rightarrow 0$  as  $n \rightarrow \infty$  (recall  $\mathcal{E}_S \leq \sqrt{\frac{\log |\mathcal{F}_\delta|}{n}}$ )

→ Supervised Learning works "asymptotically"

Today: Practical aspect (i.e., finite  $\delta, n$ )?

Key extra parameter: dimension  $d$  of input space

↳ Generic Phenomena:  $n, \delta$  need to grow exponentially in  $d$

• Curse of Dimensionality [Bellman 1950s]

Two vignettes of CoD:

(1) Approximation with polynomials

Last week we saw that polynomials have UAP

In  $d=1$ ,  $f: [0,1] \rightarrow \mathbb{R}$  conti., then  $\lim_{k \rightarrow \infty} \inf_{p \in \mathcal{P}_k} \|f-p\| = 0$

$\mathcal{P}_k = \{p: [0,1] \rightarrow \mathbb{R}, \text{polynomials of degree } k\}$

$\exists f \in C^1$ , then  $\inf_{p \in \mathcal{P}_k} \|f-p\| \leq \frac{1}{k}$

In other words, if we want  $\epsilon$ , we set degree of poly. to be  $k = \frac{1}{\epsilon}$

In  $d=1$ ,  $\mathcal{P}_k$  contains  $p(x) = x^k + a_{k-1}x^{k-1} + a_{k-2}x^{k-2} + \dots + a_0, a_i \in \mathbb{R}^k$

Q: What happens as  $d$  decreases?

$f: [0,1]^d \rightarrow \mathbb{R}, f \in C^1$

$\mathcal{P}_k = \{p: [0,1]^d \rightarrow \mathbb{R}, p \text{ is a multivariate poly. of deg. } k\}$

e.g.  $d=2, k=3: x_1^3, x_1^2x_2, x_1x_2^2, x_2^3$

→ It is not hard to check that  $\mathcal{P}_{k,d}$  has UAP in the class of smooth functions (using e.g. Stone-Weierstrass)

→ We also preserve the rate of approximation:  $\inf_{p \in \mathcal{P}_{k,d}} \|f-p\| \leq \frac{1}{k} \Rightarrow$  we need at least  $\frac{1}{\epsilon}$  degree to reach approx. error  $\epsilon$

↳ How many parameters do we need to express  $\mathcal{P}_{k,d}$ ?

$x_1^{s_1} x_2^{s_2} \dots x_d^{s_d}$  where  $s_i \in \mathbb{N}, s_i \geq 0$  &  $k = \sum_{i=1}^d s_i$

# of possible choices:  $\binom{d+k-1}{k} = \binom{d+k-1}{d-1} \stackrel{k \gg d}{\approx} \binom{k}{d} \approx k^d = \epsilon^{-d}$

→ Same is true if we replace polynomials by Neural Nets

→  $\delta = \epsilon^{-d}$  is a "signature" of Curse of Dimensionality

• Estimation of Continuous/Lipschitz Functions

Say we want to learn a target function  $f^*: [-1,1]^d \rightarrow \mathbb{R}$

from examples  $\{(X_i, y_i = f(X_i))\}_{i=1, \dots, n}$ , under the assumption

that  $f^*$  is 1-Lipschitz:  $|f^*(x) - f^*(x')| \leq \|x - x'\|$  for  $\forall x, x'$

→ A natural estimator in this setting is the Nearest Neighbor estimator

$\hat{f}(x) = f^*(X_{i(x)})$  where  $i(x) = \underset{i=1, \dots, n}{\operatorname{argmin}} \|x - X_i\|$  (Fundamental nonparametric estimator)

↳ Existence of memorization, and exploit smoothness prior

Q: How well does Nearest Neighbor do?

$\mathbb{E}_x | \hat{f}_{NN}(x) - f^*(x) | = \mathbb{E}_x | f^*(X_{i(x)}) - f^*(x) | \leq \mathbb{E}_x \|X_{i(x)} - x\|$

Uniform Distribution is optimal for the lower bound.

In this case, the expected error is  $\epsilon \sim n^{-\frac{1}{d}}$

⇔ To reach error  $\epsilon$ , we need  $n \sim \epsilon^{-d}$  points

⇔ high dimensional spaces are very lonely places!

FoML Lecture 16: Optimization in ML

Recap so far: → Focus on statistical & approximation in SL

→ We have viewed ERM as a black-box

$$\text{ERM: } \min_{f \in \mathcal{F}_S} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) = \hat{R}(\hat{f})$$

Beyond OLS, this problem does not admit a closed-form solution

→ We need to resort to iterative, optimization methods

→ We will focus on the two most important methods

- (i) Gradient Descent
- (ii) Stochastic Gradient Descent

Optimization Basics:

Consider a generic optimization  $\min_{\theta \in \mathbb{R}^d} F(\theta)$

- ① When can we solve this problem efficiently?
- ② How expensive?

Def. (Global Minimizer)

A point  $\theta^* \in \mathbb{R}^d$  is a global minimizer of  $F$  if  $F(\theta^*) \leq F(\theta)$  for  $\forall \theta \in \mathbb{R}^d$

(Local Minimizer)

A point  $\theta^* \in \mathbb{R}^d$  is a local minimizer if  $\exists \epsilon > 0$  s.t.  $F(\theta^*) \leq F(\theta)$  for all  $\theta \in B_\epsilon(\theta^*)$

Remark. Global minimizer is a (much) stronger property than local minimizer

Q: How hard is to solve a (generic) optimization problem (in high dimension)?

→ We only access the function via local queries

In general, we need to grid/explore all the domain to find the global minimum.

→ We need an exponential number of queries (Curse of dimensionality)

→ Contrary to the worst case, many typical global optimization problems can be solved by breaking them into a sequence of local optimisation problems

↳ Eg. Navigation

Given some point  $\theta_0$ , we aim to find a nearby point  $\theta_1$  s.t.  $F(\theta_1) < F(\theta_0)$

How to find such update?

(d=1 setting):

$\text{sign}(F'(\theta_0))$  indicates whether to go left/right

$F(\theta_0 + t) = F(\theta_0) + t \cdot F'(\theta_0) + o(t^2)$   
 (Taylor Expansion)

$f(t) := F(\theta_0 + t \cdot v)$   
 $f'|_{t=0} = \langle \nabla F(\theta_0), v \rangle$

So  $\theta_1 = \theta_0 - \epsilon \cdot \text{sign}(f'(0)) \cdot v$

In particular, choosing  $v = -\nabla F(\theta_0)$  is the steepest descent direction whenever  $\nabla F(\theta_0) \neq 0$

In particular, at local minimizers  $\theta$ , we must have  $\nabla F(\theta) = 0$  (necessary)

The set of points  $C = \{\theta; \nabla F(\theta) = 0\}$  are called the first-order stationary/critical points

GM =  $\{\theta; \theta \text{ is a global min}\}$

LM =  $\{\theta; \theta \text{ is a local min}\}$

$GM \subset LM \subset C \supset \text{Max}$   
 Saddle

• Gradient Descent (Cauchy, 1847)

(i) Picking an initial point  $\theta_0 \in \mathbb{R}^d$ .

(ii) For each  $t = 0, 1, 2, 3, \dots$

Pick a step-size  $\eta_t > 0$  (learning rate) and set

$\theta_{t+1} = \theta_t - \eta_t \nabla F(\theta_t)$

Remark:  $\{\theta_t\}$  is a random sequence if either  $F$  and/or  $\theta_0$  are random

Key Questions:

→ When can we guarantee that GD finds the global optimum?

How long do we need to run it? How to adjust LR?

→ How to apply it to solve ERM.

→ How to scale it to large problems?

• GD succeeds whenever  $GM = C$

In particular,  $F$  convex satisfies this property.

Pf. We only need to show that  $C \subseteq GM$  if  $F$  convex

Recall  $F$  convex: for  $\forall \theta, \forall \alpha \in [0, 1], F((1-\alpha)\theta^* + \alpha\theta) \leq (1-\alpha)F(\theta^*) + \alpha F(\theta)$   
 $\Leftrightarrow F(\theta) \geq F(\theta^*) + \frac{1}{\alpha}(F((1-\alpha)\theta^* + \alpha\theta) - F(\theta^*))$

Let  $g(\alpha) = F(\theta^* + \alpha(\theta - \theta^*))$

$F(\theta) \geq F(\theta^*) + \frac{g(\alpha) - g(0)}{\alpha}$

Mean-Value Thm.  $\exists \tilde{\alpha} \in (0, \alpha)$  s.t.  $\frac{g(\alpha) - g(0)}{\alpha} = g'(\tilde{\alpha}) = \langle \nabla F(\theta^* + \tilde{\alpha}(\theta - \theta^*)), \theta - \theta^* \rangle$

By sending  $\alpha \rightarrow 0$ ,

$F(\theta) \geq F(\theta^*) + \langle \nabla F(\theta^*), \theta - \theta^* \rangle$

So,  $F(\theta) \geq F(\theta^*)$  for all  $\theta$  if  $\theta^* \in C \Rightarrow \theta^* \in GM$

#

FML Lecture 17: Optimization II

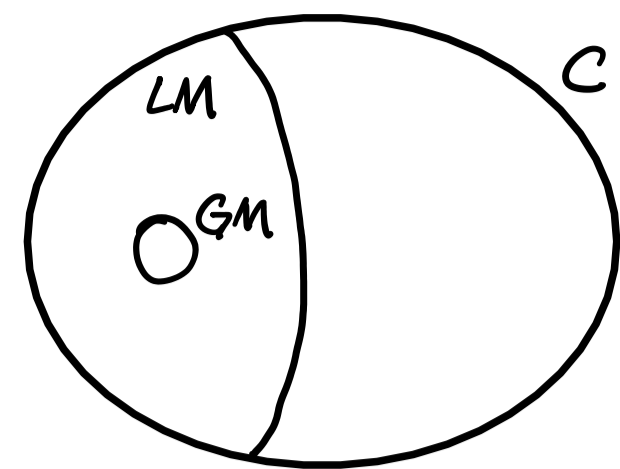
Recap: • Optimization in worst case: too hard (in high-dim)

• Approach: Use local descent method iteratively

$$C = \{\theta; \nabla F(\theta) = 0\}$$

$$LM = \{\theta; \theta \text{ is a local minimum of } F\}$$

$$GM = \{\theta; F(\theta) \leq F(\theta') \text{ for } \forall \theta'\}$$



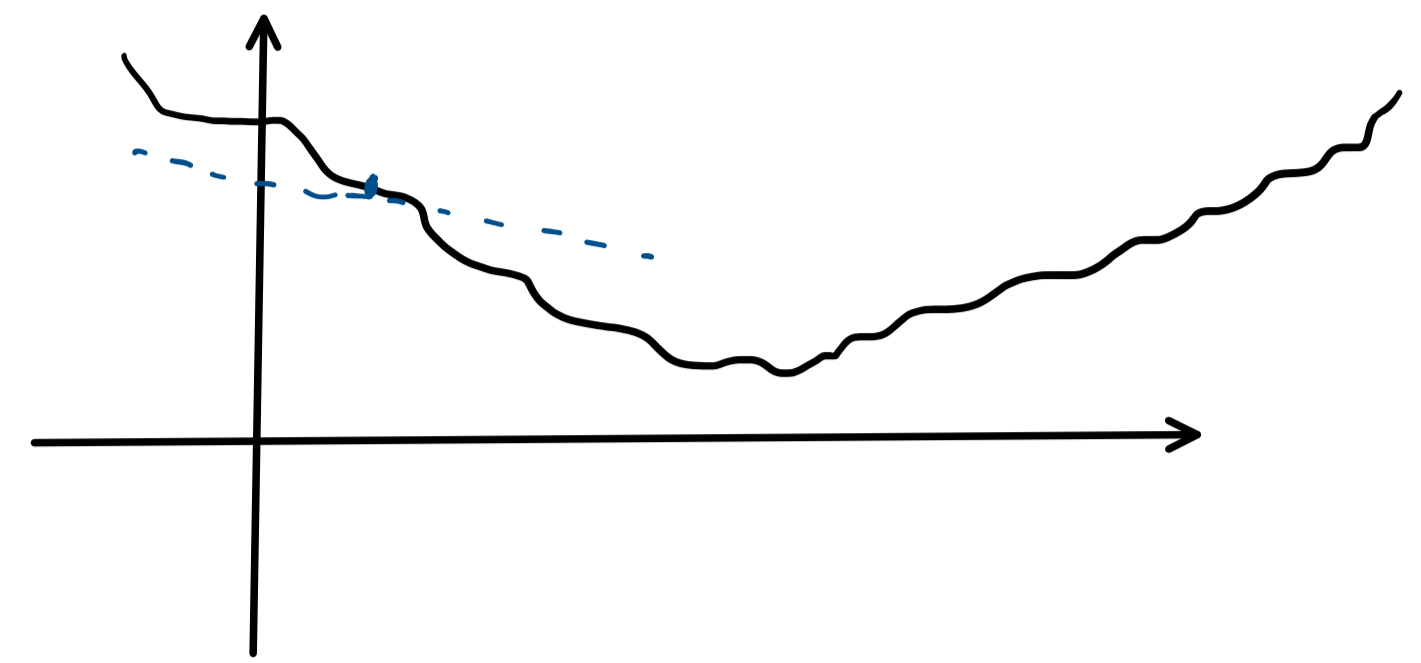
Generically, C are the equilibrium points of gradient descent

LM are the stable equilibrium points

→ There is a class of functions where  $C = GM$ : Convex Functions

Remark: There are other functions F s.t.  $C = GM$

A good example: (quasi-convex but not convex)



(\*) quasi-convex functions

$$F \text{ s.t. its level sets } S_\lambda = \{\theta; F(\theta) \leq \lambda\} \text{ are convex for } \forall \lambda$$

• If F is quasi-convex, then GD will find a global optima



$$\text{If } \lambda \leq \tilde{\lambda}, \text{ then } S_\lambda \subseteq S_{\tilde{\lambda}}$$

(\*\*) F with a Polyak - Lovajcivice (PL) inequality:

$$\|\nabla F(\theta)\| \geq a |F(\theta) - F(\theta^*)|^b$$

GD will find global optima

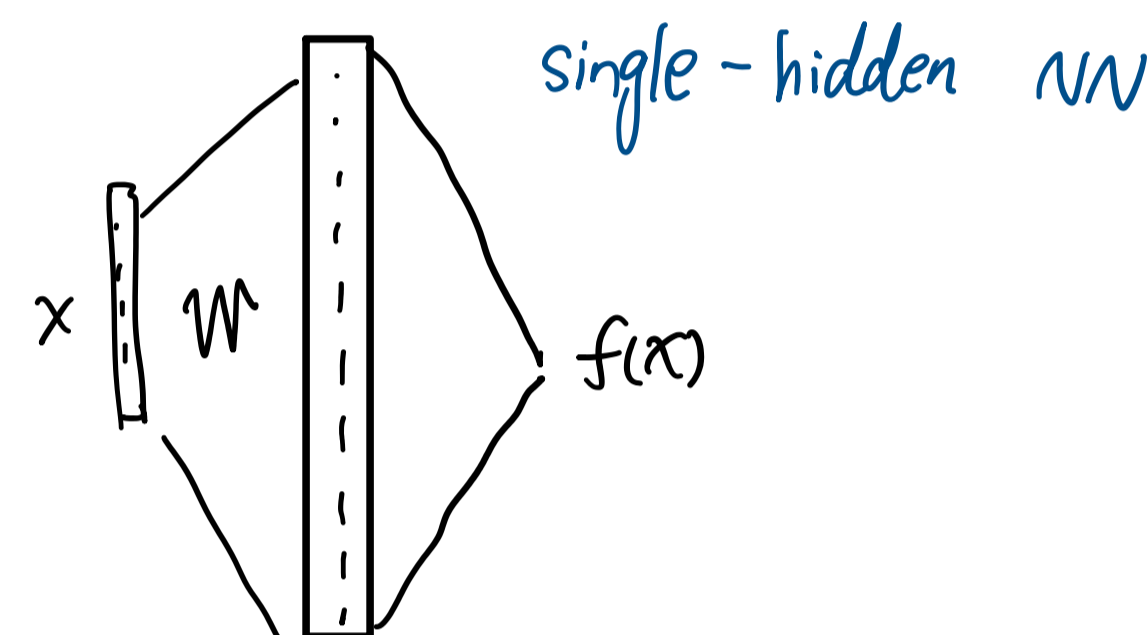
(\*\*\*) F with discrete symmetries,

$$F(T_k \theta) = F(\theta), \forall \theta, \{T_1, \dots, T_k\} \text{ is a family of transform}$$

Eg.  $\theta$  - parameters of a Neural Network

$$\theta = \{W, a\}; f(x; \theta) = a^T \cdot \sigma(Wx)$$

$$W \in \mathbb{R}^{m \times d}; a \in \mathbb{R}^m$$



Insights from quadratic functions

$$F(\theta) = \frac{1}{2n} \|H\theta - y\|^2 \text{ (ordinary least square)}$$

$$\theta \in \mathbb{R}^d \quad \cdot \nabla F(\theta) = \frac{1}{n} H^T (H\theta - y) = K\theta - \frac{1}{n} H^T y, \quad \nabla^2 F(\theta) = K$$

$$H \in \mathbb{R}^{n \times d} \quad \cdot \text{Recall that } \theta^* \text{ is a GM iff } \nabla F(\theta^*) = 0 \text{ (Normal Equations)}$$

$$y \in \mathbb{R}^n \quad K\theta^* = \frac{1}{n} H^T y$$

• F equals its 2nd-order Taylor Approximation (since F quadratic)

$$F(\theta) = F(\theta^*) + \langle \nabla F(\theta^*), \theta - \theta^* \rangle + \frac{1}{2} (\theta - \theta^*)^T \nabla^2 F(\theta^*) (\theta - \theta^*)$$

$$= F(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T K (\theta - \theta^*)$$

→ Recall that K is symmetric psd, i.e., K has eigenvalues  $\lambda_1, \dots, \lambda_d$  where  $\lambda_i \geq 0$

→ Define  $u = \min(\lambda_i), L = \max(\lambda_i), \rho = \frac{L}{u} \geq 1$  Condition Number of K

→ Gradient Descent with fixed step-size  $\eta > 0$  & initial point  $\theta_0$ .

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \cdot \nabla F(\theta_t) \\ &= \theta_t - \eta (K\theta_t - \frac{1}{n} H^T y) \\ &= \theta_t - \eta K (\theta_t - \theta^*) \end{aligned} \quad \text{as } \frac{1}{n} H^T y = K\theta^*$$

$$\Rightarrow \theta_{t+1} - \theta^* = [I - \eta K] (\theta_t - \theta^*)$$

$$= [I - \eta K]^{t+1} (\theta_0 - \theta^*)$$

→ Now we can track progress of GD: (not well defined for non-unique  $\theta^*$ )

$$\|\theta_t - \theta^*\|^2 = (\theta_0 - \theta^*)^T A^{2t} (\theta_0 - \theta^*) \text{ (Iterate Convergence)}$$

$$\rightarrow F(\theta_t) - F(\theta^*) = \frac{1}{2} \|\theta_t - \theta^*\|_K^2 \quad (= (\theta_t - \theta^*)^T K (\theta_t - \theta^*))$$

$$= \frac{1}{2} (\theta_0 - \theta^*)^T A^t K A^t (\theta_0 - \theta^*)$$

$$= \frac{1}{2} (\theta_0 - \theta^*)^T A^{2t} K (\theta_0 - \theta^*)$$

Recall that K has eigenvalues in  $[u, L]$

↳ Let us first assume that  $\theta^*$  is unique  $\Leftrightarrow u > 0$

Q: What are the eigenvalues of A?

$$\text{eigenvalues of } K: \{\lambda_i\}_{i=1}^d \longleftrightarrow \{1 - \eta \lambda_i\}_{i=1}^d: \text{ eigenvalues of } A = I - \eta K$$

$$\text{(matrix calculus)} \quad \longleftrightarrow \{(1 - \eta \lambda_i)^{2t}\}_{i=1}^d: \text{ eigenvalues of } A^{2t}$$

→ To guarantee that  $\|\theta_t - \theta^*\|^2 \xrightarrow{t \rightarrow \infty} 0$ , we want  $|1 - \eta \lambda_i| < 1$  for  $\forall i = 1, \dots, d$

Eg. Pick  $\eta = \frac{1}{L}$  where  $L = \max \lambda_i$

$$\lambda_i \in [u, L] \Rightarrow \eta \lambda_i = \frac{\lambda_i}{L} \in \left[\frac{u}{L}, 1\right]$$

$$\Rightarrow 1 - \eta \lambda_i \in \left[0, 1 - \frac{u}{L}\right] \subset [0, 1)$$

$$\parallel$$

$$[0, 1 - \rho^{-1}]$$

$$\Rightarrow \|A\|^{2t} \leq (1 - \rho^{-1})^{2t}$$



Lecture 18: Optimization (cont'd)

Recap: Analysis of GD on quadratic functions  $F(\theta) = \|\mathbf{H}\theta - \mathbf{y}\|^2 = F(\theta^*) + (\theta - \theta^*)^T \mathbf{K} (\theta - \theta^*)$

$\mathbf{K} \in \mathbb{R}^{d \times d}$  GD:  $\theta_{t+1} = \theta_t - \eta \nabla F(\theta_t)$

→ We saw (when  $\mathbf{K} \succ 0$ ):  $\|\theta_t - \theta^*\|^2 \leq (1 - \frac{1}{\rho})^t \|\theta_0 - \theta^*\|^2$

$\rho$ : Condition Number of  $\mathbf{K}$ :  $\frac{\lambda_{\max}(\mathbf{K})}{\lambda_{\min}(\mathbf{K})} \triangleq \rho$   $\xrightarrow{\text{by setting } \eta = \frac{1}{\lambda_{\max}(\mathbf{K})}}$

- Remark:
- (1) This is what we call a "linear" convergence (error decays exponentially fast)
  - (2) The bound  $1 - \frac{1}{\rho}$  comes from the operator norm of  $\mathbf{A} = \mathbf{I} - \eta \mathbf{K}$

⇒ Any choice of  $\eta \in (0, \frac{2}{\lambda_{\max}(\mathbf{K})})$  guarantees exponential convergence

- Questions:
- ① Optimality of GD?
  - ② What happens when  $\mathbf{u} = 0$  (in particular when  $d > n$ )

Answers:

- ① GD is NOT optimum amongst algorithms that only rely on gradients (first-order method)  
 Newton at 90s: Using "Momentum", one can replace  $\rho$  by  $\sqrt{\rho}$  on convergence

②:  $\mathbf{u} = 0 \Rightarrow \rho = +\infty \Rightarrow$  Previous bound says  $\|\theta_t - \theta^*\| \leq \|\theta_0 - \theta^*\|$

→ Rather than tracking  $\|\theta_t - \theta^*\|$ , now we can track  $|F(\theta_t) - F(\theta^*)|$   
 → Using again  $\eta = \frac{1}{L}$ , and recall  $F(\theta_t) - F(\theta^*) = (\theta_0 - \theta^*)^T (\mathbf{I} - \eta \mathbf{K})^{2t} \mathbf{K} (\theta_0 - \theta^*)$   
 → Let's again bound the eigenvalues of  $(\mathbf{I} - \eta \mathbf{K})^{2t} \mathbf{K}$

$$\|[\mathbf{I} - \mathbf{K}/L]^{2t} \mathbf{K}\|_{op} \leq \sup_{\lambda \in [0, L]} \left| \lambda (1 - \lambda/L)^{2t} \right| = \frac{L}{2t+1} \cdot \underbrace{\left(1 - \frac{1}{2t+1}\right)^{2t}}_{\downarrow t \uparrow \infty \rightarrow e^{-1}} \leq \frac{L}{2t+1}$$

Therefore,  $F(\theta_t) - F(\theta^*) \leq \frac{L}{2t+1} \|\theta_0 - \theta^*\|^2$   
 Convergence but much slower than  $\mathbf{u} > 0$

Recap: Till now, we have:

When  $\mathbf{K} \succ 0$  ( $\mathbf{u} > 0$ ),  $\|\theta_t - \theta^*\|^2 \leq (1 - \rho^{-1})^t \|\theta_0 - \theta^*\|^2$   
 $F(\theta_t) - F(\theta^*) \leq L (1 - \rho^{-1})^t \|\theta_0 - \theta^*\|^2$

When  $\mathbf{u} = 0$ ,  $F(\theta_t) - F(\theta^*) \leq \frac{L}{2t} \|\theta_0 - \theta^*\|^2$

In other words, to reach error  $\Sigma$ , we need  $\begin{cases} t \approx \rho \cdot \log(\frac{1}{\epsilon}) \text{ iterations } (\mathbf{u} > 0) \\ t \approx L/\epsilon \text{ iterations } (\mathbf{u} = 0) \end{cases}$

Remark:

→ We have shown that upper bounds for the loss convergence at a certain rate

↳ This may be pessimistic in practice!

↳ Contrast with "scaling laws", which looks at typical case

Q: What happens when  $\eta \rightarrow 0$ ?  $[\theta_{t+1} = \theta_t - \eta \nabla F(\theta_t)]$

The sequences  $\{\theta_t^{(j)}\}_t$  accumulates to a conti. curve  $\{\theta(t)\}_{t \in \mathbb{R}_+}$

$\frac{\theta_{t+1} - \theta_t}{\eta} = -\nabla F(\theta_t)$  Say now  $\theta_t = \theta(\eta t)$

Then  $\dot{\theta}(t_0) = \frac{\theta(t_0 + \eta) - \theta(t_0)}{\eta} = -\nabla F(\theta(t_0))$

→ having more theoretical significance

Therefore  $\{\theta(t)\}_{t \in \mathbb{R}_+}$  satisfies:  $\dot{\theta}(t) = -\nabla F(\theta(t))$  which is called Gradient Flow

→ At small  $\eta$ , GD is a time discretisation of the gradient flow ODE

FoML Lecture 19: Optimization cont'd: Convex Functions

Recap: Analysis of GD on quadratic functions

$$F(\theta) = F(\theta^*) + \frac{1}{2}(\theta - \theta^*)^T K (\theta - \theta^*) \quad \text{eigenvalues of } K: [u, L] \text{ \& condition number } \rho = \frac{L}{u}$$

→ Choosing  $\eta = 1/L$ , iteration complexity to reach error  $\epsilon$

$$t \approx \begin{cases} \rho \cdot \log(1/\epsilon) & \text{when } u > 0 \\ L \cdot \frac{1}{\epsilon} & \text{when } u = 0 \end{cases}$$

→ When  $\eta \ll 1/L \rightarrow$  GD can be analyzed using differential calculus (Gradient Flow ODE:  $\dot{\theta}(t) = -\nabla F(\theta(t))$ )

→ When  $\eta > 2/L$ , then GD diverges

Today: Beyond Quadratic Functions

• The Convex case: (Recall  $F$  convex if  $\forall x, y \in \mathbb{R}^d, \forall \alpha \in [0, 1], F(\alpha x + (1-\alpha)y) \leq \alpha F(x) + (1-\alpha)F(y)$ )

In quadratic case,  $\nabla^2 F(\theta) = K$  as a constant (matrix)

→ In the convex case,  $\nabla^2 F(\theta)$  is no longer constant, but it satisfies  $\nabla^2 F(\theta) \geq 0$

Def. A function  $F \in C^2$  is  **$\mu$ -strongly convex** if  $\nabla^2 F(\theta) - \mu I \geq 0$  for  $\forall \theta$  &  $\mu > 0$

In other words, for all  $\theta$ , all eigenvalues of  $\nabla^2 F(\theta)$  are  $\geq \mu > 0$

Def.  $F \in C^2$  is  **$L$ -smooth** if  $\nabla F$  is  $L$ -Lipschitz. ( $L > 0$ )

i.e., all eigenvalues of the Hessians  $\nabla^2 F(\theta)$  are bounded by  $L > 0$

→ **"Sandwich" Property**: If  $F$  is  $\mu$ -strongly convex &  $L$ -smooth, then:

$$\forall x, y, F(x) + \langle \nabla F(x), y-x \rangle + \frac{\mu}{2} \|x-y\|^2 \leq F(y) \leq F(x) + \langle \nabla F(x), y-x \rangle + \frac{L}{2} \|x-y\|^2$$

→ Strongly Convex Setting ( $\mu$ -sc)

→ Do we have a unique minimizer in this case?

Assume (towards contradiction) that  $\theta_1^*, \theta_2^*$  where  $\theta_1^* \neq \theta_2^*$  and both minimize.

Plug  $x = \theta_1^*$  &  $y = \theta_2^*$  in the sandwich lower bound

$$F(\theta_1^*) + \frac{\mu}{2} \|\theta_1^* - \theta_2^*\|^2 \leq F(\theta_2^*) \quad \text{!! Contradiction}$$

So minimizer is unique.

→ Strongly convex functions satisfy a P-L inequality:

$$\|\nabla F(\theta)\|^2 \geq 2\mu(F(\theta) - F(\theta^*)) \quad \text{where } \theta^* \text{ is the unique minimizer of } F$$

Pf. Recall Sandwich Property:

$$G_x(y) \stackrel{\text{variable}}{\triangleq} F(x) + \langle \nabla F(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \leq F(y) \quad (*)$$

$$\nabla_{\theta} G_x(y) = \nabla F(x) + \mu(y-x)$$

$$\Rightarrow y_x^* = x - \frac{1}{\mu} \nabla F(x)$$

So  $G_x(y_x^*) \leq F(\theta^*)$  [taking min(.) on both sides of (\*)]

$$\| \nabla F(x) + \mu(y_x^* - x) \|^2 \leq \mu \|y_x^* - x\|^2$$

$$= F(x) - \frac{1}{2\mu} \|\nabla F(x)\|^2$$

$$\Rightarrow \|\nabla F(x)\|^2 \geq 2\mu(F(x) - F(\theta^*))$$

#

→ This PL inequality allows us to establish linear convergence

Prop. Choose  $\eta = 1/L$ . The iterates of GD satisfy  $F(\theta_t) - F(\theta^*) \leq (1 - \rho)^t (F(\theta_0) - F(\theta^*))$  where  $\rho = L/\mu$

Pf. From  $\theta_t = \theta_{t-1} - \nabla F(\theta_{t-1})/L$

$$F(\theta_t) = F(\theta_{t-1} - \nabla F(\theta_{t-1})/L) \leq F(\theta_{t-1}) + \langle \nabla F(\theta_{t-1}), -\nabla F(\theta_{t-1})/L \rangle + \frac{L}{2} \|\nabla F(\theta_{t-1})/L\|^2$$

$$\Rightarrow F(\theta_t) \leq F(\theta_{t-1}) - \frac{1}{2L} \|\nabla F(\theta_{t-1})\|^2 \quad \text{[Descent Lemma]}$$

$$\text{So } F(\theta_t) - F(\theta^*) \leq F(\theta_{t-1}) - F(\theta^*) - \frac{1}{2L} \|\nabla F(\theta_{t-1})\|^2$$

$$\stackrel{\text{(previous Lem.)}}{\leq} F(\theta_{t-1}) - F(\theta^*) - \frac{\mu}{L} (F(\theta_{t-1}) - F(\theta^*))$$

$$= (1 - \rho)(F(\theta_{t-1}) - F(\theta^*))$$

$$\leq (1 - \rho)^t (F(\theta_0) - F(\theta^*))$$

#

→ As in the quadratic setting, condition number of Hessians  $\rho = L/\mu$  determines speed of convergence

Q: Continuous-time Analysis?

Recall: PL Inequality:  $\|\nabla F(\theta)\|^2 \geq 2\mu(F(\theta) - F(\theta^*))$

Gradient Flow:  $\dot{\theta}(t) = -\nabla F(\theta(t))$

Track  $F(\theta(t)) - F(\theta^*) \triangleq f(t) \geq 0$

$$f'(t) = \langle \nabla F(\theta(t)), \dot{\theta}(t) \rangle = -\|\nabla F(\theta(t))\|^2 \leq -2\mu \cdot f(t)$$

$$\Rightarrow f(t) \leq f(0) e^{-2\mu t} \quad \swarrow \text{Gronwall's lemma}$$

So the loss decreases exponentially.

Rmk. In the continuous setting, we don't see  $L$  appears.

FML Lecture 20: Optimization: Surrogate Loss, SGD

Recap: Analysis of GD on convex functions

Strongly convex setting:  $F(\theta_t) - F(\theta^*) \leq (1 - \frac{\mu}{L})^t (F(\theta_0) - F(\theta^*))$

Today: \* Analysis of GD in (vanilla) convex setting.

\* Discuss examples where convex functions appear in ML: Linear classification

\* Stochastic Gradient Descent

Reminder: from quadratic setting

When we lost strong convexity ( $\mu=0$ ), we went from a  $O((1-\mu/L)^t)$  rate to a  $O(1/t)$  rate

Q: Same thing in the general convex setting?

A: Yes!

Focus on the continuous time:  $\dot{\theta}(t) = -\nabla F(\theta(t))$

Consider the function:  $L(t) = t \cdot (F(\theta(t)) - F(\theta^*)) + \frac{1}{2} \|\theta(t) - \theta^*\|^2$ : Lyapunov Function

where  $\theta^* \in \text{argmin}_{\theta} F(\theta)$  (global minimiser)

(dynamical system)

(represents stability)

(in general, it always decreases, i.e., the system converges to stationary)

Let's compute  $L'(t) = (F(\theta(t)) - F(\theta^*)) + t \langle \nabla F(\theta(t)), -\nabla F(\theta(t)) \rangle + \langle \dot{\theta}(t), \theta(t) - \theta^* \rangle$

$= F(\theta(t)) - F(\theta^*) - t \|\nabla F(\theta(t))\|^2 - \langle \nabla F(\theta(t)), \theta(t) - \theta^* \rangle$

$= F(\theta(t)) - F(\theta^*) + \underbrace{\langle \nabla F(\theta(t)), \theta^* - \theta(t) \rangle}_{\leq 0 \text{ as } F \text{ is convex}} - t \|\nabla F(\theta(t))\|^2$

$\leq 0$

Therefore  $L(t) \leq L(0)$

$\Rightarrow t \cdot (F(\theta(t)) - F(\theta^*)) \leq L(t) \leq L(0) = \frac{1}{2} \|\theta(0) - \theta^*\|^2$

$\Rightarrow F(\theta(t)) - F(\theta^*) \leq \frac{1}{2t} \|\theta(0) - \theta^*\|^2$

Beyond Gradient Descent

↳ Momentum and acceleration: Use memory to improve convergence.  $O(1/t) \rightarrow O(1/t^2)$

$O((1-\mu/L)^t) \rightarrow O((1-\mu/L)^{t^2})$

↳ Normalisation / Adaptive Learning rates (Adam, Adagrad, ...)

↳ Second-Order Methods: use gradient  $\nabla F$  but also Hessian information  $\nabla^2 F(\theta)$

(eg. Gauss-Newton:  $\theta_{t+1} = \theta_t - \nabla^2 F(\theta_t)^{-1} \nabla F(\theta_t)$ )

(very fast but very expensive)

↳ Stochastic Gradient Descent: See next!

Examples of Convex ERM:

→ Ex 0: Linear Regression:  $\hat{R}(\theta) = \|\hat{H}^T \theta - \hat{y}\|^2$ ,  $\hat{R}$  is convex & quadratic

→ Linear Classification:  $\{(x, y)\}$  where  $y \in \{1, \dots, K\}$  ( $y$  is discrete label)

→ Simplest instance:  $K=2 \rightarrow$  Binary classification

eg. Spam filter / Fraud detection / text is AI-generated

→ Natural Loss  $\ell(y, \hat{y}) = \mathbb{1}_{\{y \neq \hat{y}\}}$ , or,  $\begin{matrix} y & -1 & +1 \\ \hat{y} & -1 & 0 & 1 \\ & +1 & 1 & 0 \end{matrix}$ ,  $\mathbb{1}_{\{y \hat{y} < 0\}}$

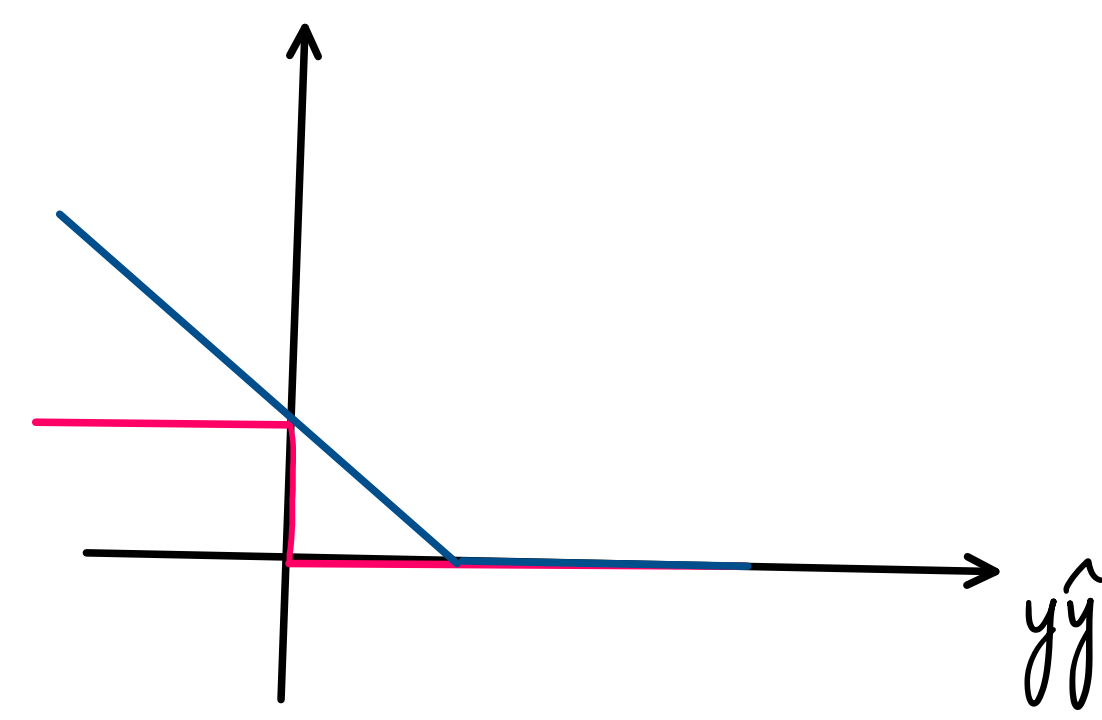
→ Associated ERM:  $\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) \rightarrow$  counts average # of mistakes

→ Q: Can we use gradient descent methods to solve this ERM?

A: No! Gradients are zero a.e.!

Sol. Replace this loss by a smoother one (with nonzero gradient): Surrogate Loss

$z = y \hat{y}$ ,  $\ell(z) = \max(0, 1-z)$  Hinge Loss

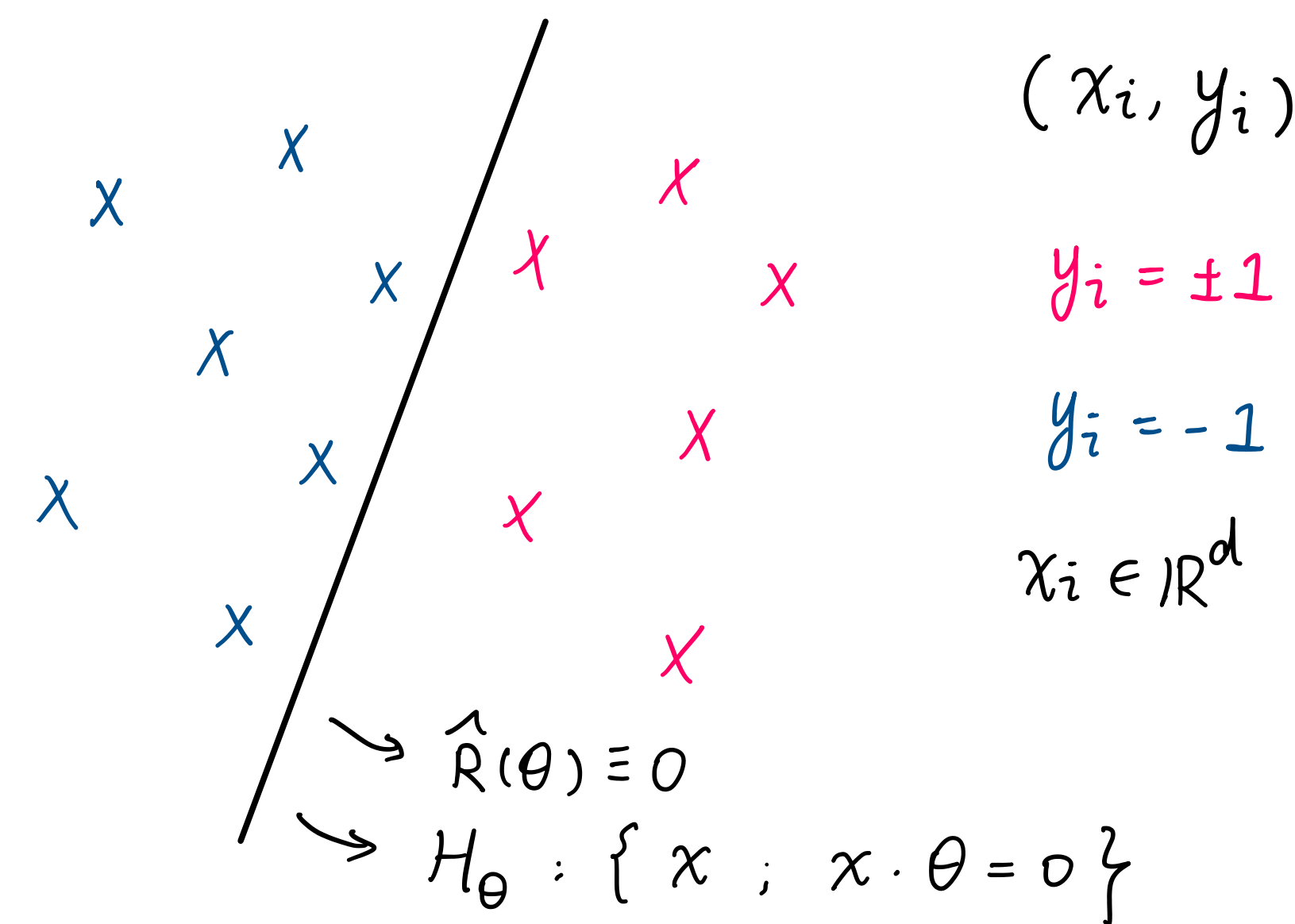


Geometric Interpretation of Hinge Loss

We want to find  $\theta \in \mathbb{R}^d$  s.t.  $\begin{cases} \theta \cdot x_i > 0 & \text{if } y_i = +1 \\ \theta \cdot x_i < 0 & \text{if } y_i = -1 \end{cases}$

Assume that data is linearly separable:  $\exists$  such hyperplane

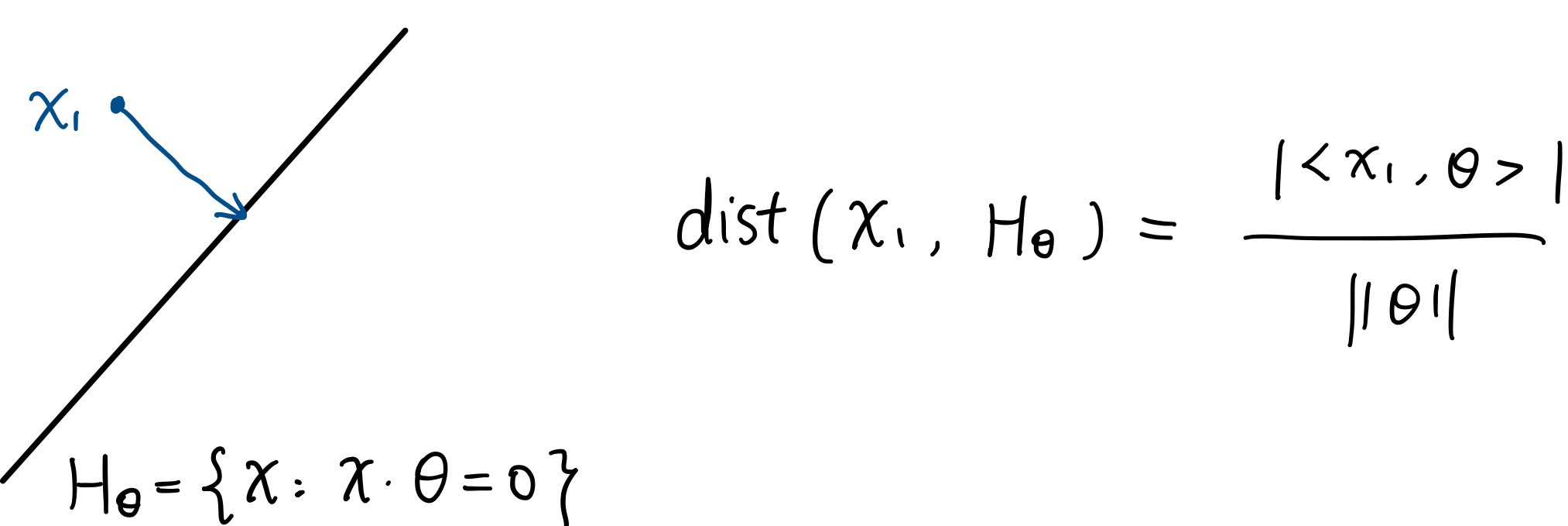
$\begin{matrix} \times & \times \\ \times & \times \end{matrix}$  (X or problem) [not linearly separable]



Q: Which hyperplane to pick amongst those that separate data?

↳ We may want to pick a hyperplane as far as possible from the data

→ Maximise the margin



Several Options:

(i) Find hyperplane with largest margin:  $\max_{\theta} \min_i \frac{y_i \langle x_i, \theta \rangle}{\|\theta\|} \Leftrightarrow \min_{\theta} \|\theta\|$  subject to  $y_i \langle x_i, \theta \rangle \geq 1$  for  $\forall i=1, \dots, n$

Support Vector Machine (Vapnik)

FoML Lecture 21: Stochastic Gradient Descent

- Recap:
- Binary Classification: Error measure  $\ell(y, \hat{y}) = 1_{\{y\hat{y} < 0\}}$  defines a loss with no gradients!
  - Introduce a surrogate loss  $\tilde{\ell}(y, \hat{y})$  with "useful" gradients
  - Margin: first example of surrogate loss

for SVM,  $\text{Margin}(\theta) = \min_i \frac{y_i \langle x_i, \theta \rangle}{\|\theta\|}$

$\rightarrow \max_{\theta} \text{Margin}(\theta)$

$\rightarrow$  penalize small margins:  $\tilde{\ell}(y, \hat{y}) = \max(1 - y\hat{y}, 0)$

$\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n \max(1 - y_i \langle x_i, \theta \rangle, 0) + \frac{\lambda}{2} \|\theta\|^2$

$\rightarrow$  This ERM is associated with the perceptron  
 [McCulloch & Pitts, 1943, Rosenblatt' 50s]

•  $\hat{R}$  is convex w.r.t.  $\theta$  (in fact it is  $\lambda$ -strongly convex)

$\Leftrightarrow$  Logistic Loss:  $\tilde{\ell}(y, \hat{y}) = \log(1 + e^{-y\hat{y}})$

$\rightarrow$  Probability Interpretation

Model:  $y|x \sim \text{Bern}\left(\frac{e^{\frac{1}{2}\langle x, \theta \rangle}}{e^{\frac{1}{2}\langle x, \theta \rangle} + e^{-\frac{1}{2}\langle x, \theta \rangle}}\right)$

$P_{\theta}(y = +1|x) = \frac{e^{\frac{1}{2}\langle x, \theta \rangle}}{e^{\frac{1}{2}\langle x, \theta \rangle} + e^{-\frac{1}{2}\langle x, \theta \rangle}} = \frac{1}{1 + e^{-\langle x, \theta \rangle}} = \frac{1}{1 + e^{-y\langle x, \theta \rangle}}$

$P_{\theta}(y = -1|x) = 1 - P_{\theta}(y = +1|x) = \frac{1}{1 + e^{\langle x, \theta \rangle}} = \frac{1}{1 + e^{-y\langle x, \theta \rangle}}$

$\rightarrow$  Consider the MLE:

$\max_{\theta} \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y_i|x_i) \Leftrightarrow \min_{\theta} \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \langle x_i, \theta \rangle}) = \hat{R}(\theta)$  using logistic loss

$\rightarrow \hat{R}$  is also convex (as  $t \mapsto \log(1 + e^{-t})$  is convex)

$\rightarrow$  All these surrogate losses can be optimized by GD (thanks to convexity)

$\rightarrow$  Big caveat: any ERM of the form  $\hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i, \theta)$  has a gradient of the form:  $\nabla_{\theta} \hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(y_i, x_i, \theta) \rightarrow$  need to use all the data all the time  
 unfeasible in large scale ML!

• Stochastic Gradient Descent [Robbins & Munro 50s]

We can view the training loss  $\hat{R}(\theta)$  as an expectation over the data:

$\hat{R}(\theta) = \mathbb{E}_{(x_i, y_i) \sim T} [\ell(x_i, y_i, \theta)] (+ \lambda H(\theta))$  optional regularization

$\nabla_{\theta} \hat{R}(\theta) = \mathbb{E}_{(x_i, y_i) \sim T} [\nabla_{\theta} \ell(x_i, y_i, \theta)]$

$i_t \sim \text{Unif}\{1, \dots, n\}$

$\rightarrow$  At each iteration  $t$ , we draw a point  $i_t \sim T$  and define

$\theta_t = \theta_{t-1} - \eta_t \underbrace{\nabla \ell(x_{i_t}, y_{i_t}, \theta_{t-1})}_{g_t(\theta_{t-1})}$   $\hookrightarrow$  (or also a "minibatch" of  $k \ll n$  points)  
 Stochastic approximation of the gradient

Q: Is SGD a descent method?

A: No! Updates might increase the loss, but should decrease "on average"

Some key questions:

(\*) Underlying assumptions that make SGD valid?

(\*) Role of the learning rate  $\eta_t$ ?

(\*) Performance in convex functions?

Key Assumptions:

(i) Unbiased Gradient Descent:  $\mathbb{E}[g_t(\theta_{t-1}) | \theta_{t-1}] = \nabla F(\theta_{t-1})$

eg.  $g_t(\theta_{t-1}) = \nabla \ell(x_{i_t}, y_{i_t}, \theta_{t-1})$   $F$ : objective function

$\mathbb{E}[g_t(\theta_{t-1}) | \theta_{t-1}] = \frac{1}{n} \sum_{i=1}^n \nabla \ell(x_i, y_i, \theta_{t-1}) = \nabla F(\theta_{t-1}) \checkmark$

(ii) Variance Control:  $\|g_t(\theta_{t-1})\|^2 \leq B^2$  a.s.

FoML Lecture 22: SGD

Recap:  $\rightarrow$  Stochastic Gradient Descent

$\theta_t = \theta_{t-1} - \eta_t g(\theta_{t-1})$ ,  $g(\theta_{t-1})$ : estimator of gradient of  $F(\theta_{t-1})$  at  $\theta_{t-1}$

$\rightarrow$  Main Example:  $F(\theta) = \mathbb{E}_x[l(\theta, X)]$

and  $g(\theta_t) = \nabla_{\theta} l(\theta_t, X_t)$  (gradient w.r.t. a single sample)

$\rightarrow$  Today: understand the role learning rate  $\eta_t$

**Problem Set-up:**

$\rightarrow X \sim P$  in  $\mathbb{R}^d$  s.t.  $\mathbb{E}_P(X) = \theta^*$ ,  $\mathbb{E}_P[\|X - \theta^*\|^2] = \sigma^2 < +\infty$

$\rightarrow$  Define  $F(\theta) = \frac{1}{2} \mathbb{E}_P[\|X - \theta\|^2]$  Global Min  $\theta^* = \mathbb{E}_P[X]$

$\rightarrow$  Goal: Minimize  $F$  using SGD:

At iteration  $t$ , we draw  $X_t \sim P$  (indep. from all previous data)

$$\theta_t = \theta_{t-1} - \eta_t \nabla_{\theta} \left[ \frac{1}{2} \|X_t - \theta\|^2 \right]_{\theta = \theta_{t-1}}$$

$$= (1 - \eta_t) \theta_{t-1} + \eta_t X_t$$

Q: How to pick  $\eta_t$ ?

Idea 1: If  $\eta_t = 1/t$ , then  $\theta_t = \frac{1}{t} \sum_{j=1}^t X_j$

We can see this by induction:

$$\theta_t = (1 - \eta_t) \theta_{t-1} + \eta_t X_t = (1 - \frac{1}{t}) \frac{1}{t-1} \sum_{j=1}^{t-1} X_j + \frac{1}{t} X_t = \frac{1}{t} \sum_{j=1}^t X_j$$

Idea 2: If  $\eta_t = \frac{2}{t+1}$ , then  $\theta_t = \frac{1}{t(t+1)} \sum_{j=1}^t j \cdot X_j$

$\hookrightarrow$  ex: Check recurrence

$\rightarrow$  Q: Principled way to select learning rate?

$\rightarrow$  From  $\theta_t = (1 - \eta_t) \theta_{t-1} + \eta_t X_t$

We have a recurrence error:  $\theta_t - \theta^* = (1 - \eta_t)(\theta_{t-1} - \theta^*) + \eta_t (X_t - \theta^*)$

$$\Rightarrow \theta_t - \theta^* = (1 - \eta_t) \left[ (1 - \eta_{t-1})(\theta_{t-2} - \theta^*) + \eta_{t-1}(X_{t-1} - \theta^*) \right] + \eta_t (X_t - \theta^*)$$

$$= \dots$$

where  $\prod_{k=t+1}^t (1 - \eta_k) \cong 1$

$$= \prod_{j=1}^t (1 - \eta_j) (\theta_0 - \theta^*) + \sum_{j=1}^t \left( \prod_{k=j+1}^t (1 - \eta_k) \right) \eta_j (X_j - \theta^*)$$
 [random]

Using that  $X_1, \dots, X_t$  are i.i.d., Rmk. We can see that as long as we have  $\prod_{j=1}^t (1 - \eta_j)^2 (\theta_0 - \theta^*)^2 \xrightarrow{t \rightarrow \infty} 0$

$$\mathbb{E}[\|\theta_t - \theta^*\|^2] = \prod_{j=1}^t (1 - \eta_j)^2 (\theta_0 - \theta^*)^2 + \sum_{j=1}^t \underbrace{\left( \prod_{k=j+1}^t (1 - \eta_k) \right)^2 \eta_j^2 \text{Var}(X_j - \theta^*)}_{\sigma^2}$$

$$= \prod_{j=1}^t (1 - \eta_j)^2 (\theta_0 - \theta^*)^2 + \sigma^2 \sum_{j=1}^t \eta_j^2 \prod_{k=j+1}^t (1 - \eta_k)^2$$
 [deterministic]

we have  $\prod_{j=1}^t (1 - \eta_j) (\theta_0 - \theta^*) \xrightarrow{t \rightarrow \infty} 0$ , so  $\theta_t \xrightarrow{a.s.} \theta^*$

So it is asymptotically unbiased

Q: Do we have any idea on controlling it to be unbiased non-asymptotically?

$\rightarrow$  To get smaller error as  $t$  increases, we need:

- (i) forget initial conditions: we need  $\prod_{j=1}^t (1 - \eta_j)^2 \rightarrow 0$  as  $t \uparrow \infty$
- (ii) Control of the variance:  $\sum_{j=1}^t \eta_j^2 \prod_{k=j+1}^t (1 - \eta_k)^2 \rightarrow 0$  as  $t \uparrow \infty$

A: (i) Assume  $\eta_t \rightarrow 0$  as  $t \uparrow \infty$

$$\log \prod_{j=1}^t (1 - \eta_j)^2 = 2 \sum_{j=1}^t \log(1 - \eta_j) \leq -2 \sum_{j=1}^t \eta_j$$
 by  $\log(1 - \eta_j) \approx -\eta_j$

$\Rightarrow$  We need  $\sum_{j=1}^t \eta_j \uparrow \infty$  as  $t \uparrow \infty$

(ii) Decomposition of variance term: assume  $\eta_t \geq 0$  & is non-increasing &  $\eta_1 \leq 1$ .

Let  $m \in [t]$ .

$$\sum_{j=1}^t \eta_j^2 \prod_{k=j}^t (1 - \eta_k)^2 \leq \sum_{j=1}^m \eta_j^2 \prod_{k=j}^t (1 - \eta_k) + \sum_{j=m+1}^t \eta_j^2 \prod_{k=j}^t (1 - \eta_k)$$

$$\leq \prod_{k=m}^t (1 - \eta_k) \sum_{j=1}^m \eta_j^2 + \eta_m \sum_{j=m+1}^t \eta_j \prod_{k=j}^t (1 - \eta_k) = (*)$$

$\hookrightarrow \prod_{k=j}^t (1 - \eta_k) \leq \prod_{k=m}^t (1 - \eta_k)$  for  $m \geq j$  ( $(1 - \eta_k) \in [0, 1]$ )

&  $\eta_j^2 \leq \eta_m \cdot \eta_j$  for  $m < j$ . (non-decreasing)

Since  $\prod_{k=j}^t (1 - \eta_k) = \exp(\log \prod_{k=j}^t (1 - \eta_k)) = \exp(\sum_{k=j}^t \log(1 - \eta_k)) \leq \exp(-\sum_{k=j}^t \eta_k)$ .

Then  $(*) \leq \exp(-\sum_{k=m}^t \eta_k) \sum_{j=1}^m \eta_j^2 + \eta_m \sum_{j=m+1}^t (1 - \prod_{k=j}^t (1 - \eta_k))$

$\underbrace{\prod_{k=j}^t (1 - \eta_k) - \prod_{k=j+1}^t (1 - \eta_k)}_{a_j}$

By observation,  $\sum_{j=m+1}^t (a_j - a_{j-1}) = a_t - a_m$

$$= \exp(-\sum_{k=m}^t \eta_k) \sum_{j=1}^m \eta_j^2 + \eta_m \left( 1 - \prod_{k=m}^t (1 - \eta_k) \right)$$

$$\leq \exp(-\sum_{k=m}^t \eta_k) \sum_{j=1}^m \eta_j^2 + \eta_m$$
 for  $m \leq t$

Recall that  $\sum_{j=1}^{+\infty} \eta_j = +\infty$  from the bias term,

$\Rightarrow$  In particular, if we consider  $\sum_{j=1}^{+\infty} \eta_j^2 < +\infty$  &  $\sum_{j=1}^{+\infty} \eta_j = +\infty$

By picking  $m = t/2$ , we have variance  $\searrow 0$  as  $t \uparrow \infty$

$\rightarrow$  Remark: ① It's a careful balance between forgetting I.C. & controlling overall variance ( $\sum_j \eta_j = +\infty$ ) & ( $\sum_j \eta_j^2 < +\infty$ )

② Our previous examples of choosing  $\eta_t = 1/t$  or  $2/(t+1)$  make sense

③  $\sum_{j=1}^{+\infty} \eta_j^2 < +\infty$  is sufficient but not necessary.

Even constant learning rate is valid (if we perform averaging of iterates)?

§ SGD in action: the Perceptron

Consider a dataset  $\{(x_i, y_i)\}_{i \in [n]}$  with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \{\pm 1\}$

We want to train a linear classifier:  $\hat{y}(\pi) = \text{sign}(\langle \pi, x \rangle)$

Perceptron Algorithm: [Rosenblatt 1950s]

(.) Start from  $\theta_0 = 0$

(.) At each step  $t = 0, 1, 2, \dots$

$\hookrightarrow$  Select a random example  $i \in [n]$

$\rightarrow$  If  $y_i \langle x_i, \theta_t \rangle < 1 \rightarrow$  mistake

$$\theta_{t+1} = \theta_t + y_i x_i$$

Otherwise  $\theta_{t+1} = \theta_t$

Rmk. If we made a mistake (wrong side / too small margin), we push it to the right side by adding  $y_i \langle x_i, y_i x_i \rangle = y_i^2 \|x_i\|^2 = \|x_i\|^2$

FoML Lecture 23: The perceptron and SGD

Recap: Given a dataset:  $S = \{(x_i, y_i)\}_{i \in [n]}$  with  $x_i \in \mathbb{R}^d$ ,  $y_i \in \{\pm 1\}$

We train a linear classifier:  $x \mapsto \text{sign}(\langle x, \theta \rangle)$ ,  $\theta \in \mathbb{R}^d$

Using a perceptron:

(i) Initialize  $\theta_0 = 0$

(ii) At each iteration  $t$ , select a sample  $i_t$

if  $y_{i_t} \cdot \langle x_{i_t}, \theta_t \rangle < 1$ , then  $\theta_{t+1} = \theta_t + y_{i_t} x_{i_t}$

Otherwise  $\theta_{t+1} = \theta_t$

Q: Link between perceptron & SGD?

Recall: The hinge loss:  $\ell(y\hat{y}) = \max(1 - y\hat{y}, 0)$

→ This defines the empirical loss:  $L(\theta) = \sum_{i=1}^n \ell(y_i \langle x_i, \theta \rangle)$

→ Consider SGD on this empirical loss:  $\theta_{t+1} = \theta_t - \eta_t \nabla_{\theta} \ell(y_{i_t} \langle x_{i_t}, \theta \rangle)$

$$= \begin{cases} -y_{i_t} x_{i_t} & \text{if } 1 - y_{i_t} \langle x_{i_t}, \theta_t \rangle > 0 \\ 0 & \text{otherwise} \end{cases}$$

→ If pick  $\eta_t = 1$  for  $\forall t$ , then we get perceptron !!!

Q: Does the perceptron learn?

Q1: Can it fit the training set?

Q2: Will it correctly classify a test data point?

Assumption: Dataset is linearly separable

→ Recall: notion of margin of a separating hyperplane

$H_{\theta} = \{x; \langle x, \theta \rangle = 0\}$  and a dataset  $S = \{(x_i, y_i)\}_{i \in [n]}$

Define the margin:  $\sigma(S, \theta) = \min_{i \in [n]} \frac{y_i \langle x_i, \theta \rangle}{\|\theta\|} > 0$

$$\sigma(S) = \max_{\theta} \sigma(S, \theta) \quad \& \quad \theta^* = \text{argmax}_{\theta} \sigma(S, \theta)$$

→ Define  $D(S) = \max_{i \in [n]} \|x_i\|$

(For Q1)

→ Thm. The perceptron algorithm makes at most  $\frac{2 + D(S)^2}{\sigma(S)^2}$  margin mistakes on any linearly-separable dataset  $S$

pf. Main Idea: Controlling # of mistakes ~ Controlling how much can  $\theta_t$  change, in fact,  $\|\theta_t\|$  suffices

Upper bound: Sp. we made a mistake at iteration  $t$ :

$$\|\theta_{t+1}\|^2 = \|\theta_t + y_{i_t} x_{i_t}\|^2 = \|\theta_t\|^2 + \underbrace{\|x_{i_t}\|^2}_{\leq D(S)^2} + \underbrace{2y_{i_t} \langle x_{i_t}, \theta_t \rangle}_{< -2}$$

So,  $m_t = \#$  of margin mistakes after  $t$  iterations

$$\|\theta_t\|^2 \leq m_t (D(S)^2 + 2)$$

Lower bound: Let  $\theta$  be any unit vector s.t.  $H_{\theta}$  is a separating hyperplane

If we make a mistake at step  $t$ :

$$\langle \theta, \theta_{t+1} - \theta_t \rangle = \langle \theta, y_{i_t} x_{i_t} \rangle \geq \sigma(S, \theta)$$

In particular,  $\langle \theta^*, \theta_{t+1} - \theta_t \rangle \geq \sigma(S)$  &  $\|\theta^*\| = 1$

$$\|\theta_t\| \geq \langle \theta_t, \theta^* \rangle = \sum_{j=1}^t \langle \theta_j - \theta_{j-1}, \theta^* \rangle \geq m_t \sigma(S)$$

Therefore,  $m_t^2 \cdot \sigma(S)^2 \leq \|\theta_t\|^2 \leq m_t \cdot (2 + D(S)^2)$

$$\Rightarrow m_t \leq \frac{2 + D(S)^2}{\sigma(S)^2}$$

→ Thus, perceptron eventually correctly classifies all training points

For (Q2):

Assume datapoints  $z_i = (x_i, y_i)$  are drawn i.i.d. from  $D$  and test point  $z \sim D$  (i.i.d.)

Thm. [Vapnik, Chvornokis]

(run until convergence)

Assume dataset  $S = \{z_1, \dots, z_n\}$  is linearly separable. Let  $\theta(S)$  be the output of the perceptron on  $S$ .

Then the prob. of making a margin mistake on  $z = (x, y)$  satisfies

$$P(y \langle \theta(S), x \rangle < 1) \leq \frac{1}{n+1} \mathbb{E} \left[ \frac{2 + D(\bar{S})^2}{\sigma(\bar{S})^2} \right] \text{ where } \bar{S} = S \cup \{z\}$$

pf. We exploit the exchangeability of the data  $\{z_i\} = \{(x_i, y_i)\}_{i \in [n]}$  and  $z = (x, y)$

Joint distribution of  $x_1, \dots, x_n$  does not depend on the order

$$(1) P[y \langle \theta(S), x \rangle < 1] = \mathbb{E} [1_{\{y \langle \theta(S), x \rangle < 1\}}]$$

• Define  $S^{-k} \triangleq \{z_1, \dots, z_{k-1}, z_{k+1}, \dots, z_n, z\}$

Exchangeability: the order of these R.V.'s does not affect the test error

i.e., Running perceptron on  $S^{-k}$  and testing on  $z_k$  gives the same prob. error for each  $k$

$$\text{So } P[y \langle \theta(S), x \rangle < 1] = \frac{1}{n+1} \sum_{k=1}^{n+1} \mathbb{E} [1_{\{y_k \langle \theta(S^{-k}), x_k \rangle < 1\}}]$$

→ Recall that running perceptron on  $\bar{S}$  makes at most  $m = \frac{2 + D(\bar{S})^2}{\sigma(\bar{S})^2}$  mistakes.

There are at most  $m$  indices  $i_2, \dots, i_m \in [n]$  where we have made mistakes ( $m \leq n$ )

If  $k \notin \{i_2, \dots, i_m\}$ , then  $\theta(\bar{S}) = \theta(S^{-k})$

$$\Rightarrow y^k \langle \theta(S^{-k}), x_k \rangle > 1$$

Other terms contribute at most 1. So

$$P[y \langle \theta(S), x \rangle < 1] \leq \frac{1}{n+1} \cdot \mathbb{E}[m] = \frac{1}{n+1} \cdot \mathbb{E} \left[ \frac{2 + D(\bar{S})^2}{\sigma(\bar{S})^2} \right]$$

Remark: Unlike SVM, perceptron does not necessarily converge to a unique hyperplane

it stops as soon as it makes no mistake

From the perspective of SGD, the reason we can converge by just

choosing a constant convergence rate can be boiled down to the simplicity

of our hinge loss

# FML Lecture 25: Geometric Deep Learning

→ Basic Supervised Learning Set-up:

Input Space  $\mathcal{X}$  (high-dimensional)

Output Space  $\mathcal{Y}$  (low-dimensional, e.g.,  $\mathcal{Y} = \mathbb{R}$ )

Hypothesis Class:  $\mathcal{F} = \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$ , often indexed by a complexity parameter  $\mathcal{F}_\delta = \{f \in \mathcal{F}, \mathcal{V}(f) \leq \delta\}$

Goal: Approximate unknown target  $f^*$  via ERM

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{F}_\delta} \frac{1}{n} \sum \ell(f(x_i), y_i) \text{ where we assume } y_i = f^*(x_i) + \varepsilon$$

Recall Decomposition of error:

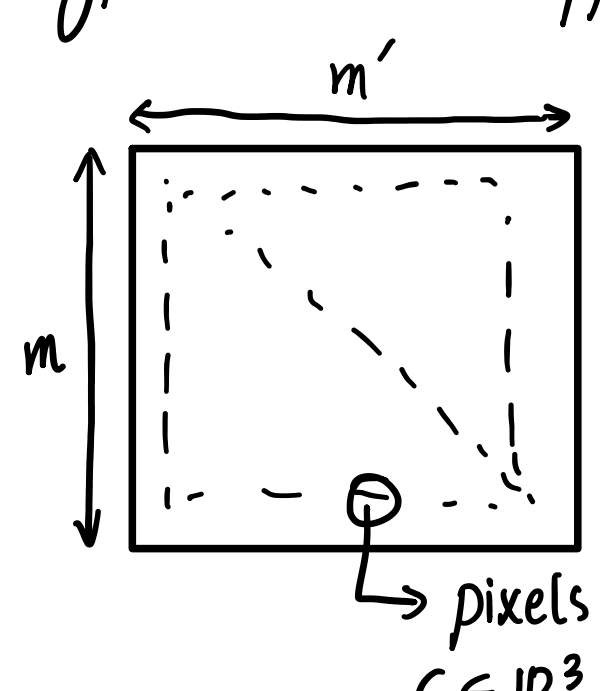
$$R(\hat{f}) \leq \varepsilon_{\text{approx}}(\delta) + \varepsilon_{\text{stat}}(\varepsilon)$$

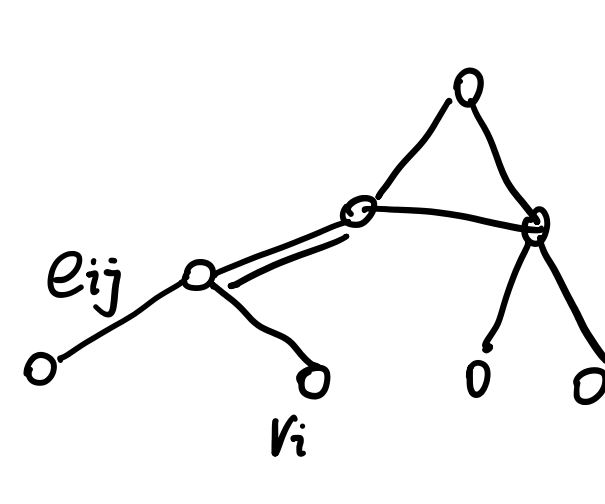
Conclusion: To efficiently learn, we need accurate ( $\varepsilon_{\text{approx}}$  small) yet "small" hypothesis  $\mathcal{F}_\delta$  ( $\varepsilon_{\text{stat}}$  small)

⇒ Need to exploit any prior information on target  $f^*$

## Learning in the Physical World

High-dimensional input space  $\mathcal{X}$  in typical ML applications?

→  $\mathcal{X} = \{\text{images}\}$  represented as   $x \in \mathbb{R}^{3 \times m \times m'}$   
 2D-grid encoding RGB

→  $\mathcal{X} = \{\text{molecules}\}$  represented as   $v_i \in \{O, C, H, N, \dots\}$ ,  $e_{ij} \in \mathbb{R}^5$   
 graphs, encoding the atoms and their chemical bounds

→  $\mathcal{X} = \{\text{text/language}\}$ , represented as a sequence  $\{w_1, w_2, \dots, w_n\} \in \text{Dictionary}$

$\mathcal{X}$  is in fact a space of signals that live on a physical domain  $\Omega$ :  $\mathcal{X} = \{x: \Omega \rightarrow \mathbb{C}\}$   
 $u \mapsto x(u)$  (2D-grid, graphs, sequence) channels

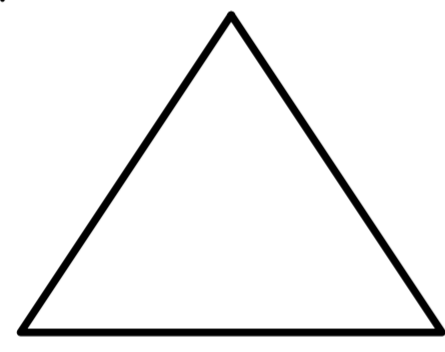
→ We can add signals or scale them,  $(\alpha x + \beta y)(u) = \alpha x(u) + \beta y(u) \Rightarrow \mathcal{X}$  is a vector space

→ Inner product structure in  $\mathbb{C}$  "upgrades" to inner product structure in  $\mathcal{X}$ :  $\langle x, y \rangle_{\mathcal{X}} = \int_{\Omega} \langle x(u), y(u) \rangle_{\mathbb{C}} du$

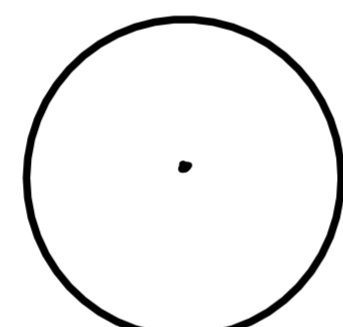
Q: Why is this physical domain useful?

Symmetry: A symmetry of an object is a transformation that leaves the object unchanged

Ex. 1:



finite # of symmetries



infinite # of symmetries

Ex 2:  $f(x; W_1, W_2) = W_2 \rho(W_1 x)$

A symmetry of this architecture is a transformation of parameters  $W = \{W_1, W_2\}$

s.t.  $f(x; g(W)) = f(x; W)$  for  $\forall W$

$\sigma: \{1, m\} \rightarrow \{1, m\}$  a permutation

Given  $\sigma$ , we define perm. matrix  $\Pi_\sigma \in \{0, 1\}^{m \times m}$  where  $(\Pi_\sigma)_{ij} = \begin{cases} 1 & \text{if } \sigma(i) = j \\ 0 & \text{otherwise} \end{cases}$

So  $g_\sigma(\{W_1, W_2\}) = \{\Pi_\sigma W_1, W_2 \Pi_\sigma^T\}$  is a symmetry of  $f$  s.t.  $f(x; g_\sigma(W)) = f(x; W)$  for  $\forall \sigma, \forall x, \forall W$

Rmk. ① Permutation Symmetry is indep. of the form  $\sigma$  !! (So we at least have  $m!$  symmetries for one-layer NN)

② e.g.  $\rho(t) = t$ , then  $f(x; g_\sigma(W)) = W_2 \Pi_\sigma^T \Pi_\sigma W_1 x = W_2 W_1 x = f(x; W)$

so, it's natural to think: if we have some assumptions on  $\rho$ , we can explore more symmetries (like orthogonal)  
 (consider homogeneous functions like ReLU)

→ Most importantly, we are interested in symmetries of the target  $f^*: \mathcal{X} \rightarrow \mathcal{Y}$

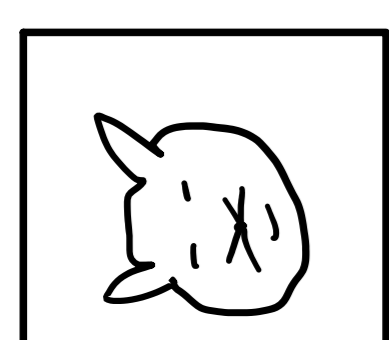
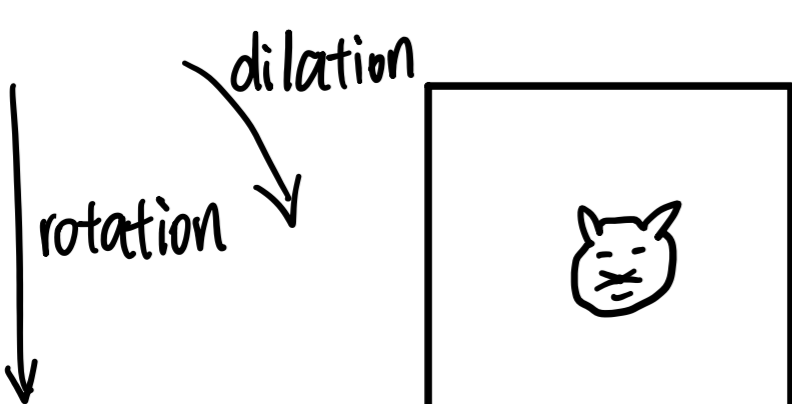
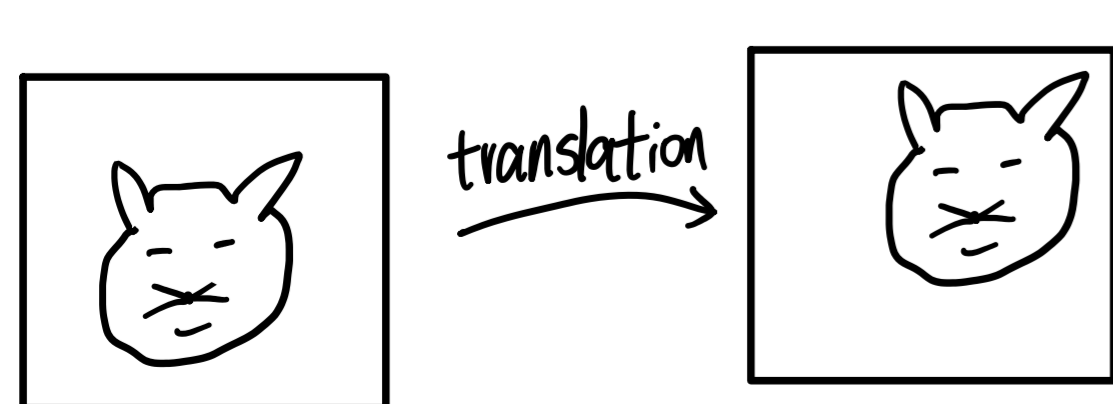
↳ Transformations  $g: \mathcal{X} \rightarrow \mathcal{X}$  s.t.  $f^*(g(x)) = f^*(x)$  for  $\forall x \in \mathcal{X}$

↳ Challenging in generic high-dimension!

↳ Instead, use physical domain  $\Omega$  to describe symmetries!

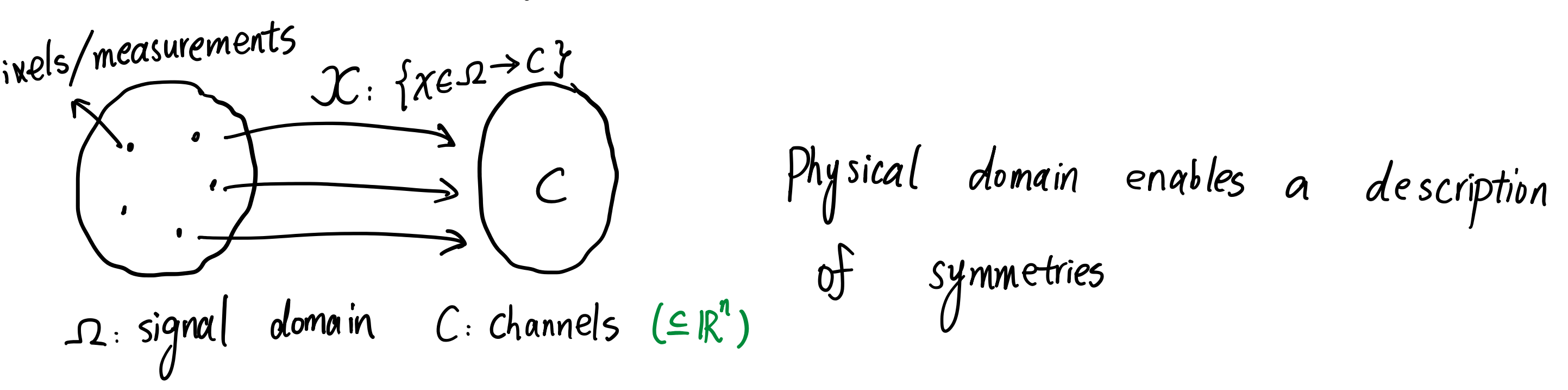
## Symmetries of $\Omega$

→ Images  $f^*(x) = \text{Is there a cat in } x$ ?

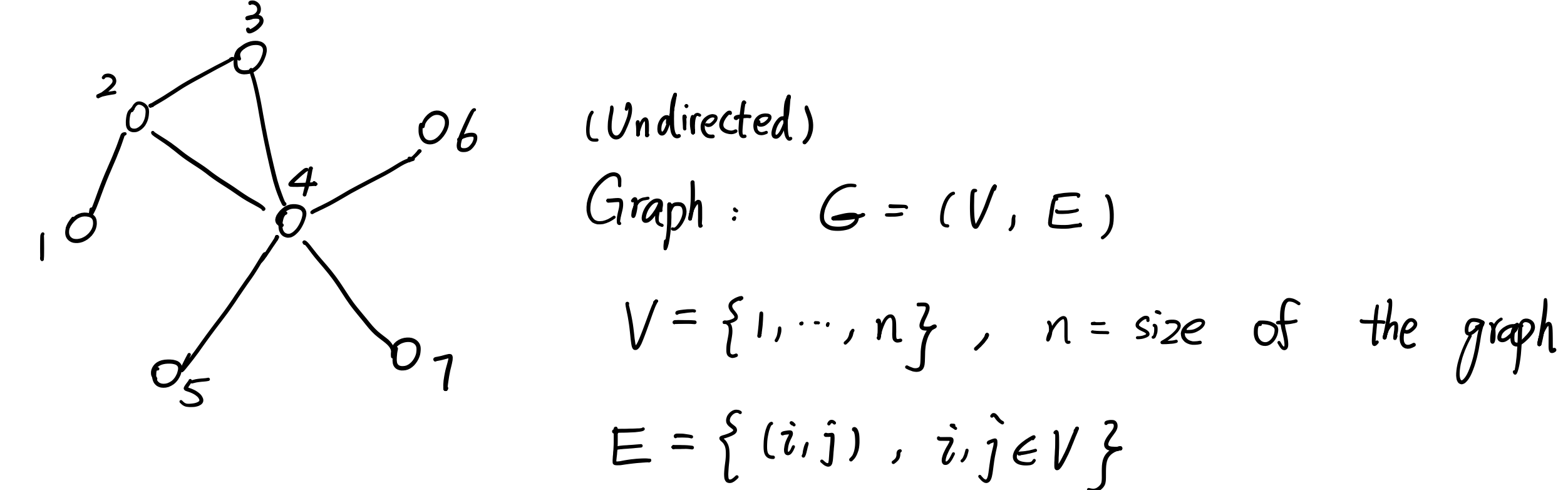


FoML Lecture 26: Learning with Symmetries

Recap: Learning in Physical World



Symmetries arising on graphs (motivation: molecules, traffic network)



Represented with Adjacency Matrix:

$A \in \{0, 1\}^{n \times n}$ , where  $A_{ij} = \begin{cases} 1, & \text{if } (i, j) \in E \\ 0 & \text{otherwise} \end{cases}$

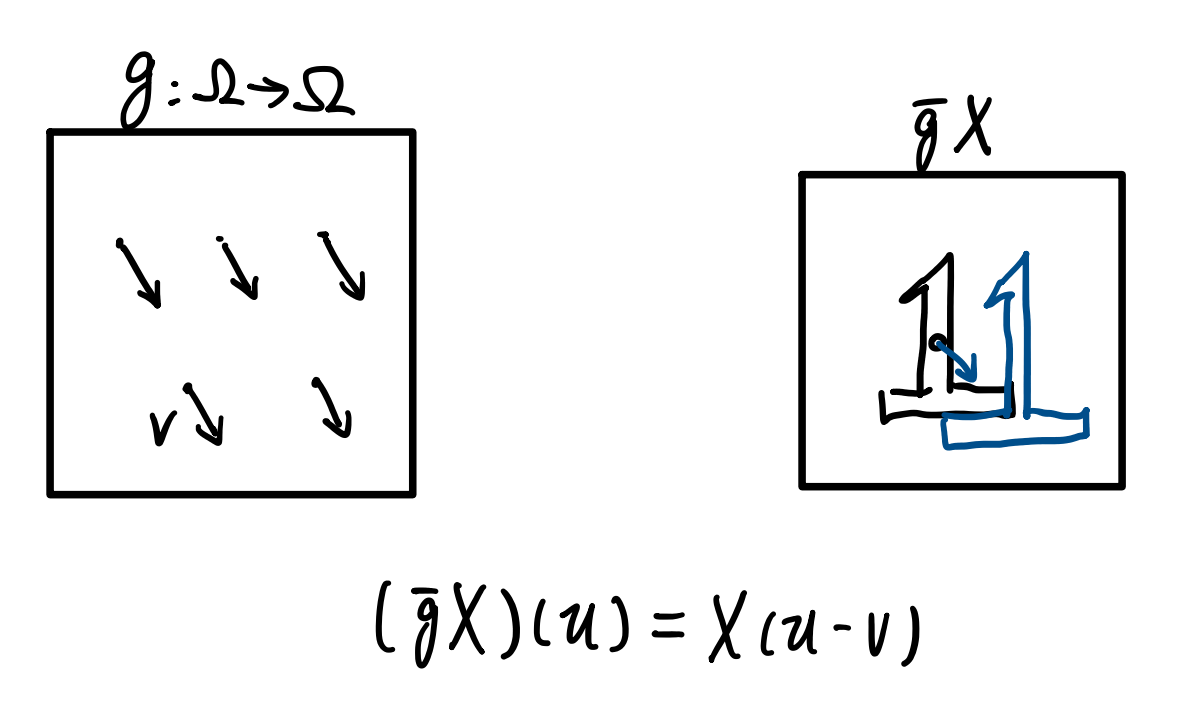
$G$  of size  $n$   $\xrightarrow{\text{encoding}}$   $A$

$\rightarrow$  We can relabel the vertices in  $n!$  ways, while all the adjacency matrices are related by:  $\bar{A} = \Pi A \Pi^T \Leftrightarrow \bar{A}$  obtained by permuting rows & columns of  $A$

From Symmetries of  $\Omega$  to the symmetries of  $\mathcal{X}$

A domain transformation  $g: \Omega \rightarrow \Omega$  defines a transformation  $\bar{g}: \mathcal{X} \rightarrow \mathcal{X}$  by:  $(\bar{g}X)(u) = X(g^{-1}u)$  for  $u \in \Omega$

$\rightarrow \bar{g}$  defines a linear transformation  $\bar{g}(\alpha X + \beta Y) = \alpha \bar{g}(X) + \beta \bar{g}(Y)$



Symmetries & Groups

- We observe that (i)  $g = \text{Id}$  is a symmetry
- (ii)  $g$  and  $h$  are symmetries, then  $g \circ h$  and  $h \circ g$  are also symmetries
- (iii) If  $g$  is a symmetry, then its inverse  $g^{-1}$  is also a symmetry

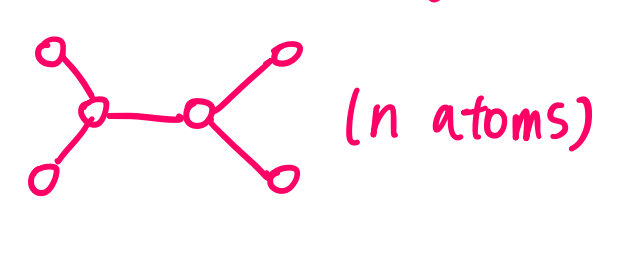
Symmetries form a group using composition

$\rightarrow$  Groups can either be discrete (finite elements) or continuous

- Ex.  $\rightarrow \mathbb{Z}_2$ : cyclic group of integers modulo 2
- $\rightarrow$  Rubik's Cube
- $\rightarrow (\mathbb{R}, +), (\mathbb{R}^n, \cdot), (\mathbb{R}^n, \times)$

Summary so far: (i) We use physical domain  $\Omega$  to define a group  $G$  of transformations  
(ii) This group defines a symmetry group of the target function  $f^*$ :  $f^*(g \cdot x) = f^*(x)$  for  $\forall g \in G, x \in \mathcal{X}$

$\rightarrow$  In many ML applications, we have prior knowledge of (some) symmetries of the target

Examples:	Arithmetic	Symmetry
	$x_1 + x_2 = ?$	Commutative Structure
	Image Classification	Orthogonal Group $O_n$
	Proteins / Biology  (n atoms)	Permutation Group $S_n$ & Orthogonal Group $O_n$

- Q1: Why symmetries are useful for learning?
- Q2: How to leverage them in practice?

Invariant Learning

$\rightarrow$  Let  $f^*: \mathcal{X} \rightarrow \mathcal{Y}$  be the target function  
 $\rightarrow$  we consider symmetry group  $G$  acting on  $\mathcal{X}$   
 $\rightarrow$  We say that  $f^*$  is invariant to  $G$  (or  $G$ -invariant) if  $f^*(g \cdot x) = f^*(x)$  for  $\forall g \in G, x \in \mathcal{X}$

$\rightarrow$  assume w.l.o.g.,  $G$  is discrete

$\rightarrow$  Given any  $f: \mathcal{X} \rightarrow \mathcal{Y}$ , we define the average w.r.t.  $G$  as  $Sf: \mathcal{X} \rightarrow \mathcal{Y}$  s.t.  $Sf(x) = \frac{1}{|G|} \sum_{g \in G} f(g \cdot x)$  for  $\forall x \in \mathcal{X}$   
 $\{g \cdot x, g \in G\} = O(x)$ : Orbit of  $G$  passing through  $\mathcal{X}$

$\rightarrow$  So  $S$  is thus averaging over group orbits

$\rightarrow Sf^* = f^*$  as  $f^*(g \cdot x) = f^*(x)$  for  $\forall x \in \mathcal{X}, \forall g \in G$

$\rightarrow$  Given a hypothesis class  $\mathcal{F}$ , we can make it  $G$ -invariant:  $S\mathcal{F} = \{Sf, f \in \mathcal{F}\}$  (as long as  $g$  &  $g^2$  acts transitively on  $\mathcal{F}$ )

Eg.  $\Omega = \{1, \dots, q\}$  for  $q > 2$  &  $q$  is prime.  $G = \mathbb{Z}_q$   
 $(\chi_i: i \rightarrow C)$   
 $\mathcal{F} = \{\text{polynomials } P(\chi_1, \dots, \chi_2) \text{ of degree } k\}$

Eg. for  $k=2$ ,  $\mathcal{F}$  could contain  $\chi_1^2, \chi_1 + \chi_2 \chi_3, \dots$   
 $S\mathcal{F} = \{Sp(\chi_1, \dots, \chi_2) = \frac{1}{q} \sum_{j=1}^q P(\chi_{1+j}, \chi_{2+j}, \dots, \chi_{2+j})\}$ ,  $P$  polynomials

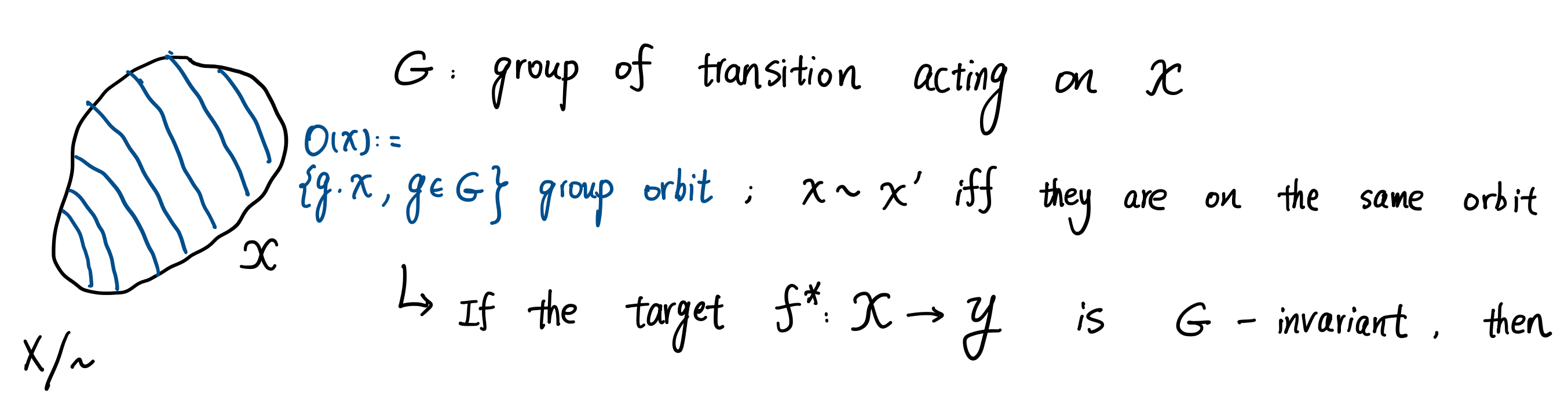
Essential Idea:  $\mathcal{X} = G \otimes \underline{\underline{\mathcal{X}/G}}$

$\rightarrow$  and we only need this information  
we pick  $Sf$  as the representative of each equivalence class  $[\bar{f}]$



FoML Lecture 27: The Geometric DL Blueprint

Recap: Decomposition of input space into orbits



→ Averaging operator  $Sf(x) = \frac{1}{|G|} \sum_{g \in G} f(g \cdot x)$  maps arbitrary hypothesis space  $\mathcal{F} = \{f: X \rightarrow Y\}$  into an invariant hypothesis class:  $S\mathcal{F} = \{Sf, f \in \mathcal{F}\}$

→ When  $f^*$  is  $G$ -invariant, should we use  $\mathcal{F}$  or  $S\mathcal{F}$ ?

↳ Verify that  $Sf$  is  $G$ -invariant:  $Sf(g \cdot x) = Sf(x)$  for  $\forall x \in X, \forall g \in G$

• Approximation Error:  $\inf_{f \in \mathcal{F}} \|f^* - f\|^2$  v.s.  $\inf_{\tilde{f} \in S\mathcal{F}} \|f^* - \tilde{f}\|^2$

$$Sf(g \cdot x) = \frac{1}{|G|} \sum_{g' \in G} f(g' \cdot (g \cdot x)) = \frac{1}{|G|} \sum_{\tilde{g} \in G} f(\tilde{g} \cdot x) = Sf(x) \quad \#$$

Fact: The averaging operator  $S$  is an orthogonal projection

pf. for  $\forall h: X \rightarrow Y$ :  $\|h\|^2 = \|Sh + (I-S)h\|^2$   
 $= \|Sh\|^2 + \|(I-S)h\|^2 + 2\langle Sh, (I-S)h \rangle$

As  $\langle Sh, (I-S)h \rangle = \int_{X/\sim} \left( \int_{O(x)} Sh(\bar{x}) \cdot (I-S)h(\bar{x}) d\bar{x} \right) dx = \int_{X/\sim} h(x) \left[ \underbrace{\int_{O(x)} h(\bar{x}) d\bar{x}}_0 - h(x) \right] dx = 0 \quad (*)$

So  $S$  is an orthogonal projection.

By this fact, we have

$$\|f^* - f\|^2 \stackrel{\text{Fact}}{=} \|Sf^* - Sf\|^2 + \|(I-S)f^* - (I-S)f\|^2$$

$$= \|f^* - Sf\|^2 + \|(I-S)f\|^2$$

$$\geq \|f^* - Sf\|^2$$

$$\Rightarrow \inf_{\tilde{f} \in S\mathcal{F}} \|f^* - \tilde{f}\| \leq \inf_{f \in \mathcal{F}} \|f^* - f\|^2$$

→ So Approximation error is not degraded (if  $S\mathcal{F} \subseteq \mathcal{F}$ , then they're equal)

→ Statistical Error?

$S\mathcal{F}$  is defined over smaller space  $X/\sim$ , so stat. error is not degraded either

→ Using Symmetries helps the learning task

→ The larger the symmetry group, the smaller the quotient space  $X/\sim$

• Big caveat so far: Computing  $Sf$  is expensive, especially as  $|G|$  is large, even  $|G|$  is  $\infty$ !

↳ Q: Efficient Algorithm?

The Geometric DL Blueprint

Consider a linear hypothesis  $f$ :

$$Sf(x) = \frac{1}{|G|} \sum_{g \in G} f(g \cdot x) = f\left(\left[\frac{1}{|G|} \sum_{g \in G} g\right] \cdot x\right) = f(\bar{x})$$

where  $\bar{x}$ : group average of  $x$

→ the group average can be computed efficiently in our cases of interest

Eg 1.  $\mathcal{X} = \{x: \Omega \rightarrow \mathbb{R}\}, \Omega = \{1, \dots, m\}, G = S_m$   
 $\cong \mathbb{R}^m$

$$\bar{x} \in \mathbb{R}^m \text{ and } \bar{x}_j = \frac{1}{m!} \sum_{\sigma \in S_m} (\sigma \cdot x)_j = \frac{1}{m} \sum_{i=1}^m x_i \rightarrow \text{simple average over all coordinates!}$$

easy ~

Eg 2.  $\mathcal{X} = \{x: \mathbb{R}^2 \rightarrow \mathbb{R}\}, G = \text{Translation Group}$

$$\bar{x}(u) = \int_G (g \cdot x)(u) = \int_G x(u-v) dv = \int_G x(l) dl \rightarrow \text{average of image}$$

Problem: The averaging loses too much information!

How to complement the linear invariant?

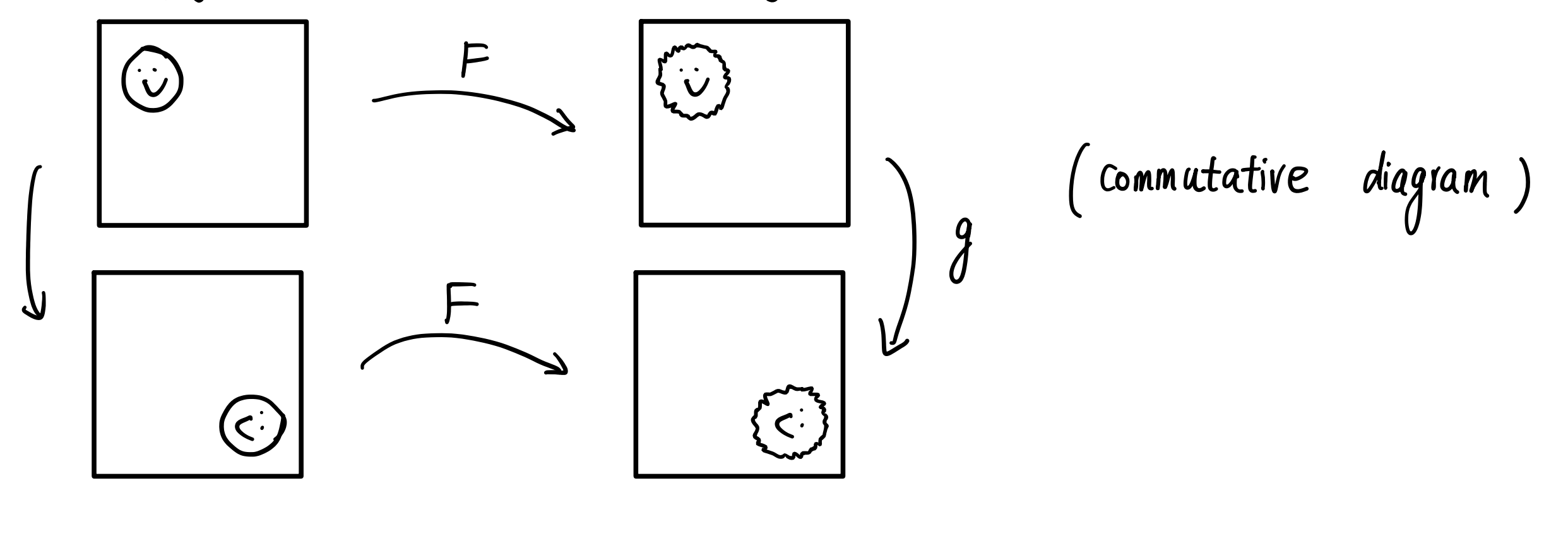
From invariance to equivariance

→ Sps. we have  $F: \mathcal{X} \rightarrow \mathcal{X}'$  s.t.  $G$  acts on both  $\mathcal{X}$  and  $\mathcal{X}'$

Eg.  $\mathcal{X}' = \mathcal{X} = \text{image}, F(x) = x$  with Van Gogh style

Rmk.  $\mathcal{X}, \mathcal{X}'$  can differ, Eg.  $\mathcal{X}' = \Omega, F(x) = \text{location of certain object}$

→ We say that  $F$  is  $G$ -equivariant if  $F(g \cdot x) = g \cdot F(x)$  for  $\forall x \in \mathcal{X}$  and  $\forall g \in G$



Q: How to compute linear equivariants?

We start with  $G$ : Translation Group in  $\Omega = \mathbb{R}^2, \mathcal{X} = \{x: \Omega \rightarrow \mathbb{R}\}, F: \mathcal{X} \rightarrow \mathcal{X}$  and  $F$  is linear and commutes with translations

⇒  $F$  is a convolution:  $(Fx)(u) = \int_{\Omega} x(v) h(u-v) dv$  where  $h: \Omega \rightarrow \mathbb{R}$  is a filter

Verify linearity & commutativity w/ translation

$$F(\alpha x + \beta x') = \alpha F(x) + \beta F(x') \text{ for } \forall x, x' \in \mathcal{X} \text{ \& } \alpha, \beta \in \mathbb{R}$$

$$F(g \cdot x)(u) = \int_{\Omega} (g \cdot x)(v) h(u-v) dv$$

↳ translation by  $u_0$

$$= \int_{\Omega} x(v-u_0) h(u-v) dv$$

change of variables

$$= \int_{\Omega} x(v') h(u-u_0-v') dv' = g \cdot (Fx)(u)$$

Eg 2.  $F(x)(u) = \phi(x(u))$  pointwise transformation for  $\forall x, \forall u$

$$F(g \cdot x)(u) = \phi(x(g^{-1}u))$$

$$[g \cdot F(x)](u) = g \cdot \phi(x(u)) = \phi(x(g^{-1}u))$$

$F_1: \mathcal{X} \rightarrow \mathcal{X}'$   
 $F_2: \mathcal{X}' \rightarrow \mathcal{X}''$   
 $F_3$  is invariant

are both equivariant } ⇒  $F_2 \circ F_1$  is also equivariant  
 $F_3 \circ F_1$  is invariant