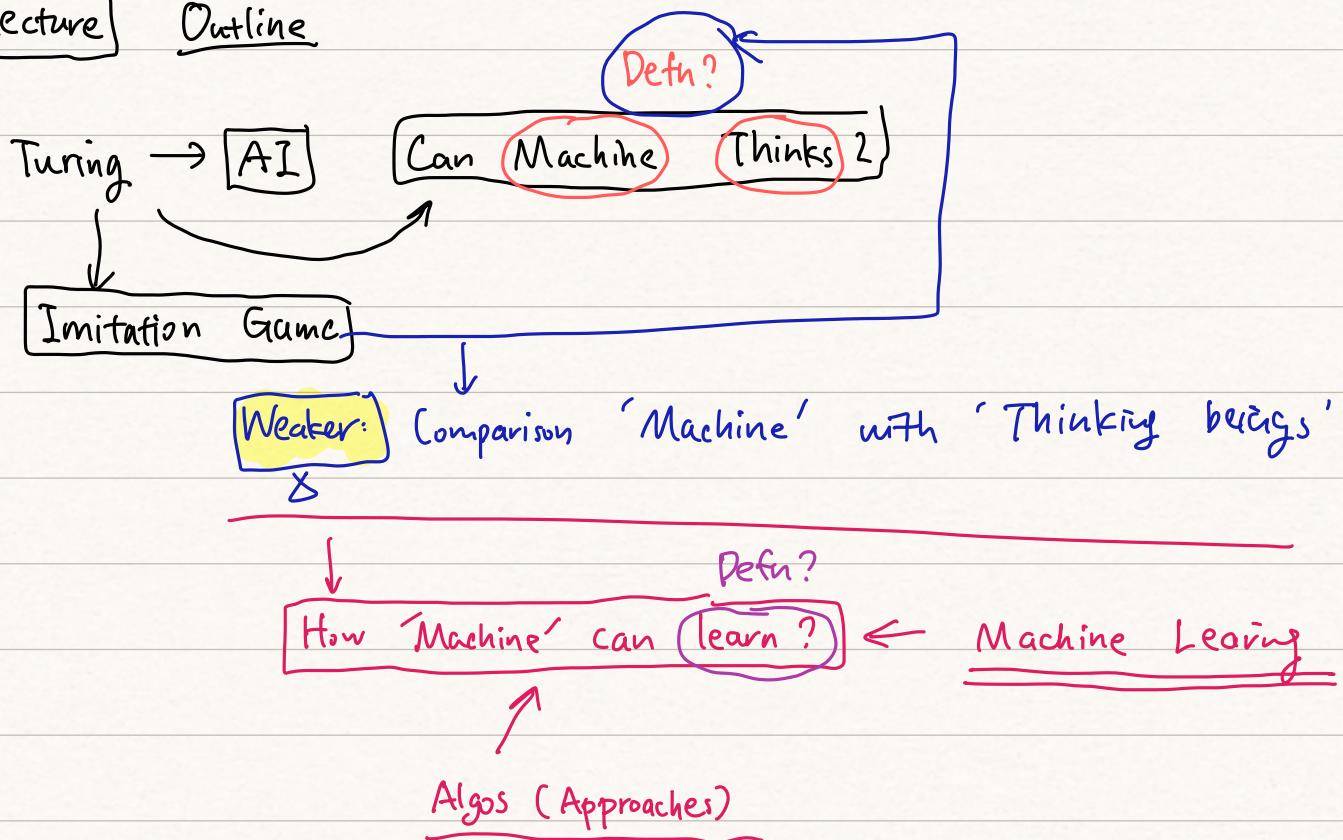


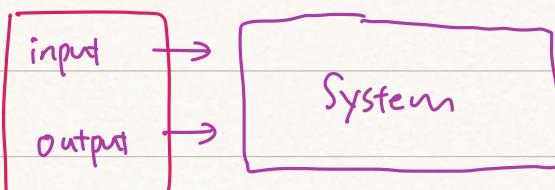
S17-05-08 ← officeLectureOutlineLearn Defn:

Experience
Task
Measure

E
T
P

AI

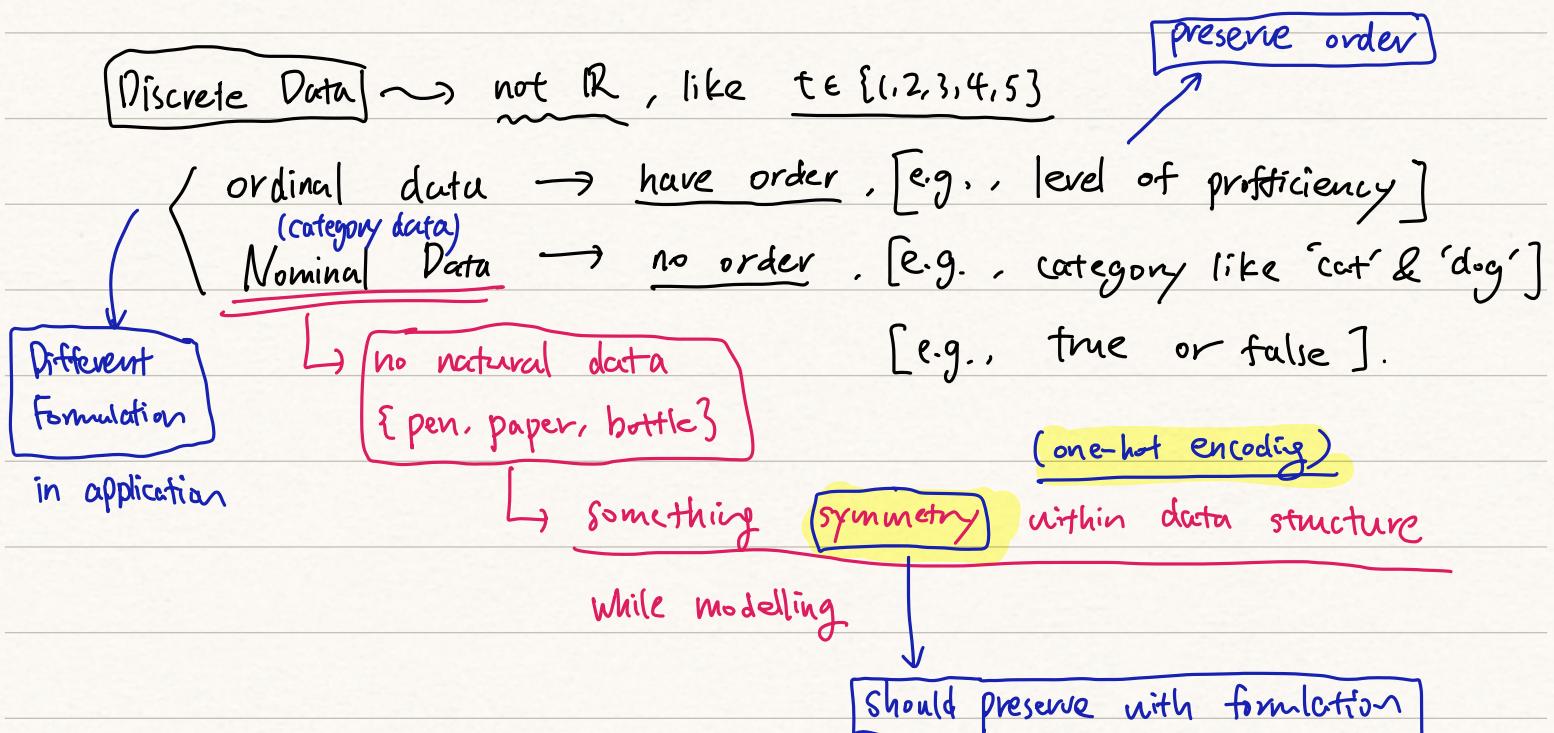
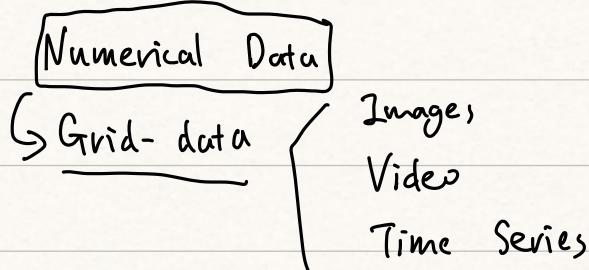
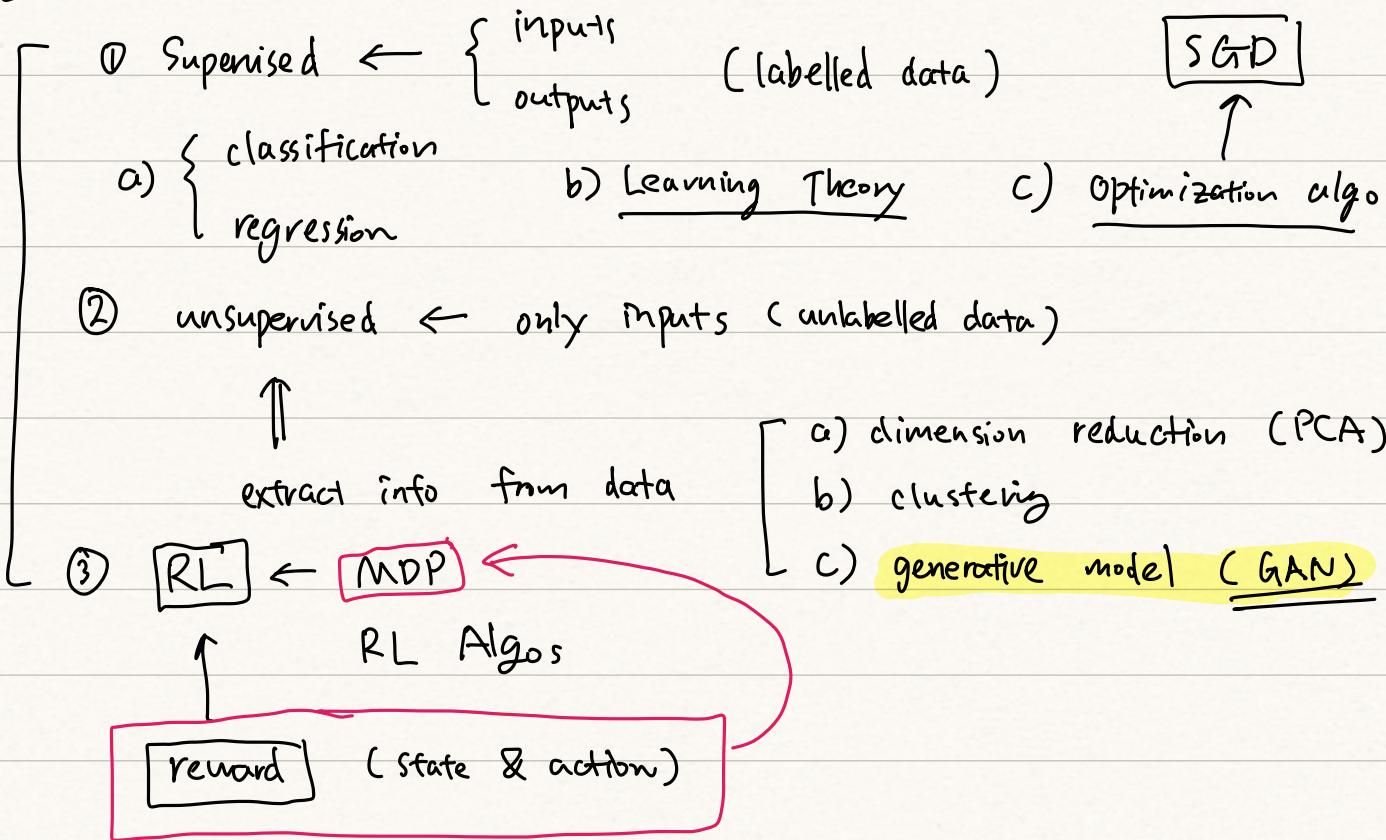
Learn Program from
data

Example (Learning Process)experience & Task

measure of how good the answer is

↑
Measure P

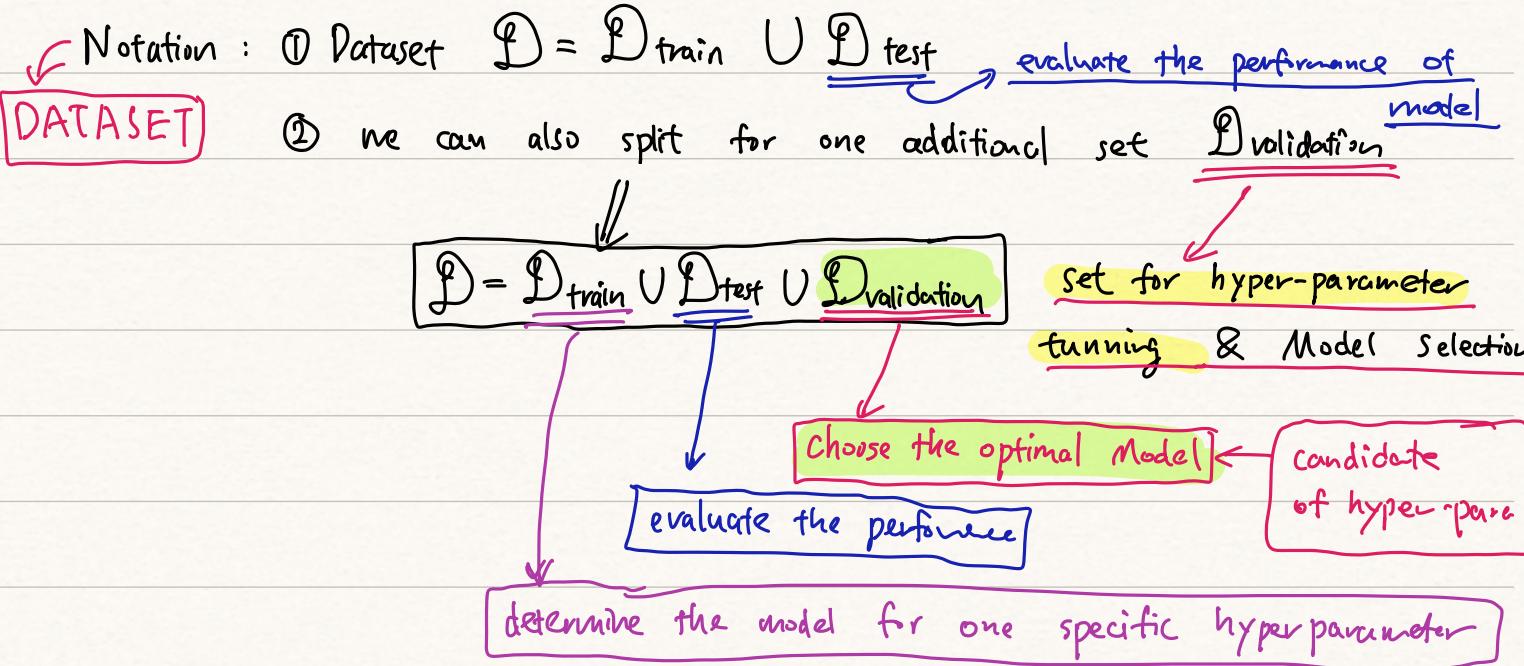
Outline



inverse problem / design

?

supervised learning



Framework:

parameters candidate

when training set is BIG, validation set can be eliminated. Do model selection on Training set Directly

train — validation

Find the optimal one

evaluate on test set

Supervised Learning Algo

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

Oracle $f^*(\cdot)$

goal: learn relationship :

$$x_i \rightarrow y_i$$



$$\text{learn } f^*(\cdot) \text{ s.t. } f^*(x_i) = y_i \quad i=1, 2, \dots, N$$

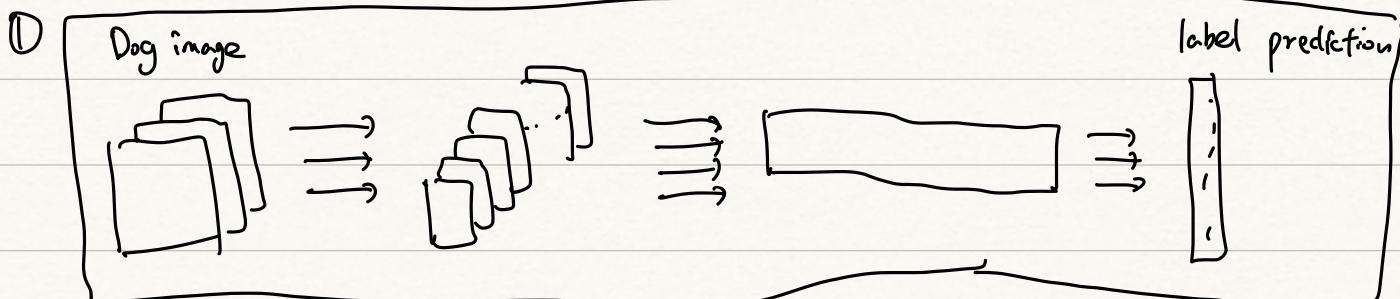


Image Recognition CNN architecture (basis)

②

$$\begin{cases} y_i = f^*(x_i) \rightarrow \text{choice 1: deterministic} \\ y_i = p(\cdot | x_i) \Leftrightarrow y_i = f^*(x_i) + \epsilon_i \rightarrow \text{choice 2: random} \end{cases}$$

↑
Different Model

→ Oracle f^* is unknown! we want to learn through ③

→ Approach

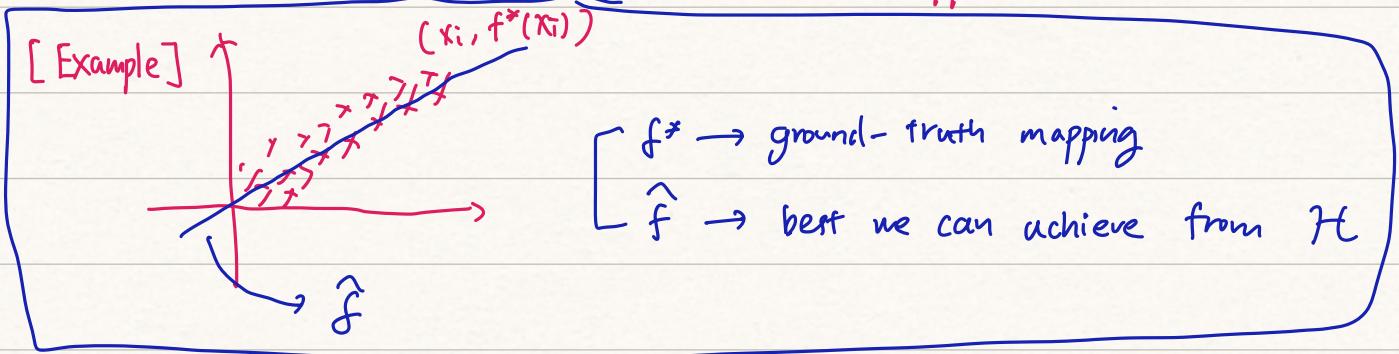
① Define Hypothesis class \mathcal{H}

$$\mathcal{H} = \{f: f = w_0 + w_1 x\}$$

specify ML model

define a function class

find a best f , namely \hat{f} to approximate f^*



Rank: from different \mathcal{H} (hypothesis), we achieve different \hat{f}
(i.e., $\hat{f} \rightarrow \hat{f}_H$)

Q: What is BEST APPROXIMATION means?

A: through LOSS FUNCTION!

[Example] $L(y', y) = \frac{1}{2}(y' - y)^2$

[IDEA]: Loss function: $L(y', y)$ $\begin{cases} \text{small if } y \approx y' \\ \text{large if } y \neq y' \end{cases}$

then BEST APPROXIMATION can be formulated as

One optimization problem:

$$\min_{f \in \mathcal{H}} \text{Remp}(f) = \min_{f \in \mathcal{H}}$$

$$\frac{1}{N} \sum_{i=1}^N$$

$$L(f(x_i), f^*(x_i))$$

(one way)

ERM Framework
to achieve \hat{f}

average loss on the
training set

y_i (ground-truth
label)

approximation of $\mathbb{E}_{x \sim D} [L(f(x), f^*(x))]$

⇒ convert learning approximation problem into optimization problem

⇒ use optimization framework to solve it

Recall: our goal is Generalization (on unseen data)

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \text{Remp}(f)$$

→ not generalization

Approximation

population risk minimization

Another choice

$$\tilde{f} = \arg \min_{f \in \mathcal{H}} R_{\text{pop}}(f)$$

$$= \arg \min_{f \in \mathcal{H}} \mathbb{E}_{x \sim \mu} [L(f(x), f^*(x))]$$

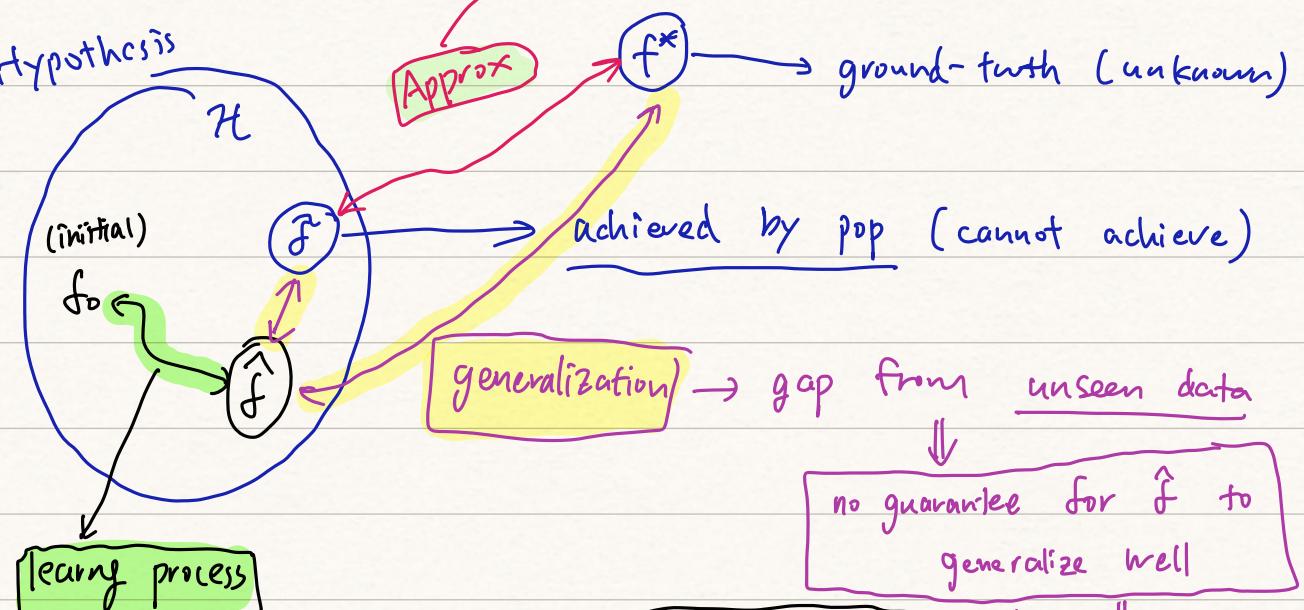
we want to solve

$$= \arg \min_{f \in \mathcal{H}} \mathbb{E}_{x \sim \mu} [L(f(x), y)]$$

but cannot achieve

Summary

depends on how large our Hypothesis Space is



through optimization and achieved with D_{train}

Statistical ML

generalization: $\hat{f} \sim \tilde{f}$ or f^*

optimization: $f_0 \sim \hat{f}$

approximation: $\tilde{f} \sim f^*$

Problem → How to find a \hat{f} to generalize well?

Linear Model

Linear Regression

$$\Rightarrow \mathcal{H} := \{ f: f = w_0 + w_1 x, w_0 \in \mathbb{R}, w_1 \in \mathbb{R} \}$$

$$\Rightarrow L(y', y) = \frac{1}{2} (y' - y)^2$$

$$\Rightarrow \min_{f \in \mathcal{H}} R_{emp}(f) = \min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

Necessary Condition

$$\begin{cases} \frac{\partial R_{emp}}{\partial w_0} = 0 \\ \frac{\partial R_{emp}}{\partial w_1} = 0 \end{cases} \rightarrow \begin{pmatrix} \hat{w}_0 \\ \hat{w}_1 \end{pmatrix}$$

$$\Rightarrow \hat{f}(x) = \hat{w}_0 + \hat{w}_1 x$$

→ Question: Is $\mathcal{H}_{\text{linear}}$ good enough? → Small Hypo → under-fit

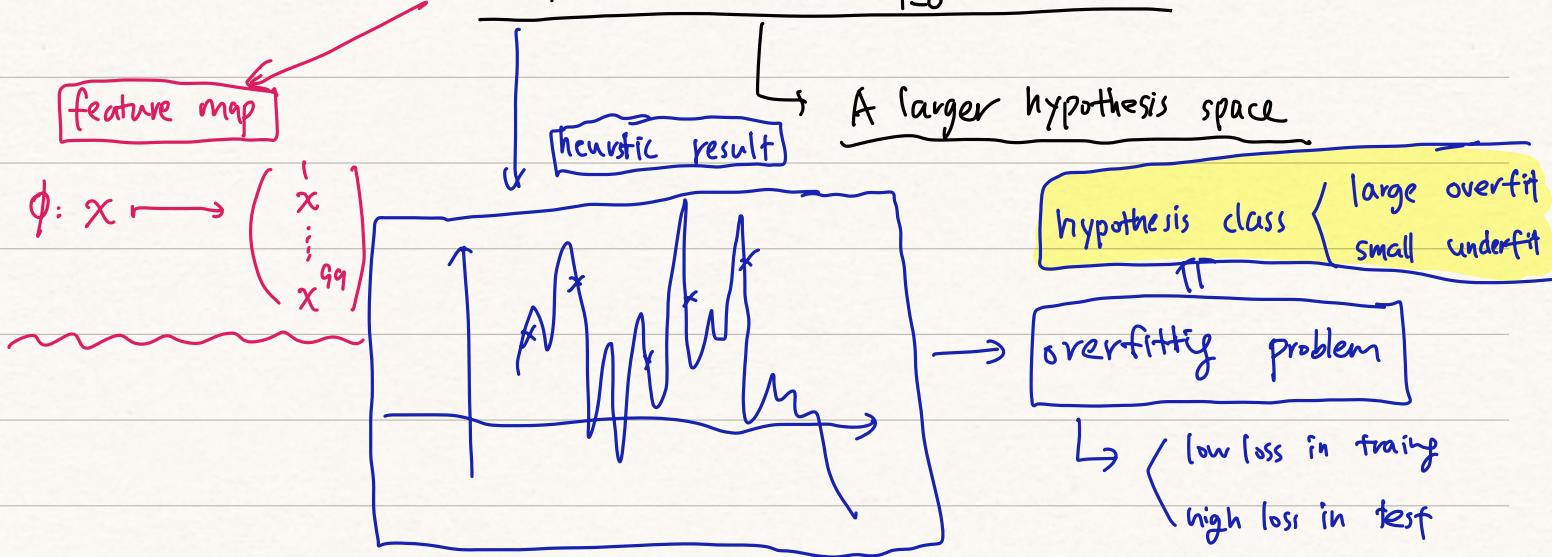
A: Depends on Task (Data), i.e. f^*

ground-truth mapping

→ Overfitting Problem

capacity increases

Consider $H_{99} = \{ f: f(x) = \sum_{i=0}^{99} w_i x^i \} \Rightarrow$ not linear model



Linear Basis Model

Generalization of Linear Model

$$\rightarrow H_M = \{ f: f(x) = \sum_{j=0}^{M-1} w_j \phi_j(x) \} = \{ f: f(x) = w^T \phi(x), w \in \mathbb{R}^M \}$$

in which $\phi_j: \mathbb{R}^d \rightarrow \mathbb{R}^M$ mapping (feature map)

\Leftrightarrow Linear Model (on feature space) $x \mapsto (\phi_0(x), \dots, \phi_{M-1}(x))$

[By the Way, another generalization for Linear Regression is Neural Network]

(use activation function)

Learn Process

Model

$$y_i = w^T \phi(x_i) = \phi^T(x_i) w$$

$$\phi(x_i) = \begin{pmatrix} \phi_0(x_i) \\ \vdots \\ \phi_{M-1}(x_i) \end{pmatrix}$$



$$y = \Phi \cdot w$$

$$\Phi = \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{pmatrix}$$



$$\boxed{\text{Loss}} \quad \min_{f \in \mathcal{H}_M} \text{Remp}(f)$$

$$= \min_{w \in \mathbb{R}^m} \frac{1}{2N} \|y - \Phi w\|^2$$

$$= \min_w \frac{1}{2N} (y - \Phi w)^\top (y - \Phi w)$$

$$= \min_w \frac{1}{2N} [y^\top y - 2 w^\top \Phi^\top y + w^\top \Phi^\top \Phi w]$$

$$\boxed{\frac{\partial L}{\partial w} = \frac{1}{2N} [-2y^\top \Phi + 2w^\top \Phi^\top \Phi]}$$

Necessary Condition

$$\nabla L = \frac{1}{2N} [-2\Phi^\top y + 2\Phi^\top \Phi w] \equiv 0$$

actually is sufficient

$$\nabla^2 L = \Phi^\top \Phi \geq 0$$

$$\Rightarrow \hat{w} = (\Phi^\top \Phi)^{-1} \Phi^\top y$$

(if $(\Phi^\top \Phi)$ is invertible!!!)

the feature map cannot be
so highly-dimensional

can be
formulated as
LM + Gaussian Noise

if not invertible, then Regularization

objective:

$$\min_w \frac{1}{2N} \|y - \Phi w\|^2 + \lambda C(w)$$

IDEA: Make w small

penalty term

$\lambda \rightarrow$ hyper-parameter

① $C(w) = \|w\|_2 \rightarrow$ Ridge regression

(l2 reg)

$$② C(w) = \|w\|_1 \rightarrow \boxed{\text{LASSO}} \xrightarrow{\text{advantage of sparse solution}}$$

$(\ell_1 \text{ reg})$

selection of features

Add $\boxed{\text{PRIORITY}}$ that, large weight will lead to over-fitting

Apply to K-class classification

Softmax

(Vector 2 probability distribution)

$\rightarrow L(y; y')$ change to Cross Entropy Loss

$$:= - \sum_{k=1}^K y_k \log y'_k$$

$$= - y_{k(c)} \log y'_{k(c)}$$

$$= - \log y'_{k(c)}$$

$\left\{ \begin{array}{l} y \rightarrow \text{ground-truth label} \\ y' \rightarrow \text{output of model} \end{array} \right.$

Rmk: $L_{\text{cross-entropy}}(y, y') = 0 \Leftrightarrow y = y'$