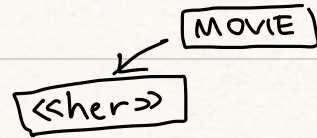
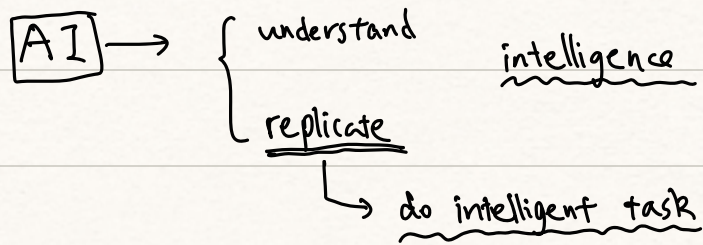
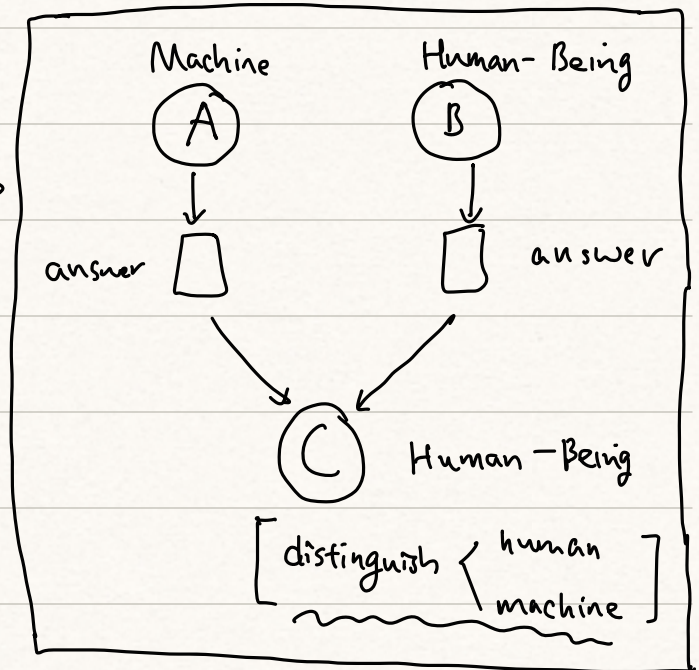
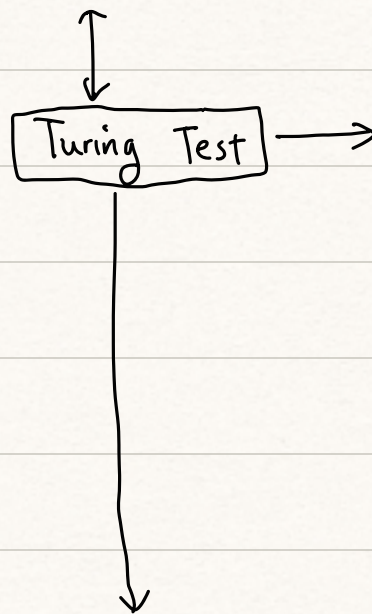


→ Low Yi Rui S17-05-19



Turing → Can Machine Thinks?



How machine learns to do things human-like?

(Machine Learning)

→ Adversarial Examples ⇒ try to analyze the change of response

when given different inputs

## Topics

1. architecture model

2. learning algo

→ learn how to learn

3. practical techniques

4. application

5. active areas of DL research

## Machine Learning Basis

1. Task :  $T \Rightarrow$  problem setting

Experience :  $E \Rightarrow$  data

{ prediction   
 transcription   
 translation   
 imputation   
 generation   
 }   
 regression   
 classification

Performance Measure :  $P$

(improve  $P$  through  $E$ )

{ MSE   
 Accuracy   
 image generation  $\Rightarrow$  hard to define }

$\Rightarrow$  Q: How to set learning objective to improve  $P$ ?

## 2. Linear Regression & Canonical Form



Oracle:  $f^*$

Learning:  $\hat{f} = \underset{f \in \mathcal{H}}{\operatorname{argmin}} \text{Distance}(f, f^*)$  (ERM Framework)

Empirical Risk Minimization

RV depends on  
training data

$$:= \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{R_{\text{emp}}(f, f^*)}{\text{empirical loss}}$$

$$= \underset{f \in \mathcal{H}}{\operatorname{argmin}} \frac{1}{2N} \sum_{i=1}^N (y_i - f(x_i))^2$$

$$\mathcal{H}_{\text{linear}} := \{ f: f(x) = \omega^T x, \omega \in \mathbb{R}^d \}$$

$$= \underset{\omega \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2N} \sum_{i=1}^N (y_i - \omega^T x_i)$$

$$= \underset{\omega \in \mathbb{R}^d}{\operatorname{argmin}} \|y - X\omega\|_2^2$$

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \in \mathbb{R}^{n \times d}$$

Benefit of Linearity

Design Matrix

exact solution

$\Rightarrow$  To solve this, we have  $\hat{\omega} = (X^T X)^{-1} X^T y$

For more complicated model (NN), we do not have closed-form sol<sup>n</sup>

iterative optimization algo  
(SGD, ADAM, ...)

3. Model Capacity &  $\begin{cases} \text{overfitting} \\ \text{underfitting} \end{cases}$

$$\mathcal{H}_{\text{linear}} \subseteq \mathcal{H}_{\text{quadratic}} \subseteq \dots$$

↓  
Small hypothesis space

↓  
big hypothesis space

#### 4. Linear Basis Model (extension of LR)

↓ → (improve the model capacity)  
Enlarge the hypothesis space

$$\mathcal{H}_{\text{LBM}}^{\phi} := \{ f: f(x) = \omega^T \phi(x), \omega \in \mathbb{R}^D \}$$

$$= \sum_{i=1}^D \omega_i \phi_i(x)$$

$$\phi(x) = \begin{pmatrix} \phi_1(x) \\ \vdots \\ \phi_D(x) \end{pmatrix} \in \mathbb{R}^D$$

( $\phi$  is fixed by design  $\Rightarrow \phi: x \in \mathbb{R}^d \mapsto \phi(x) \in \mathbb{R}^D$ )

↳ { polynomial  
radial basis (Gaussian)

→ universal approximator { Fourier Basis  
Polynomial Basis ( $D$  sufficiently large)

→ Easy to solve!  $\Rightarrow \hat{\omega} = (\Phi^T \Phi)^{-1} \Phi^T y$

$$\Phi = \begin{pmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{pmatrix} \in \mathbb{R}^{n \times D}$$

5. Large Model Capacity  $\Rightarrow$  over-fitting { large testing error  
small training error  
(behave bad on unseen data)  
use some Learning Theory to interpret this:



a) Generalization Error  $\leq$  Model Capacity<sup>①</sup> + Best we can do given the hypothesis space<sup>②</sup>

b) We are interested in

$$\min_f R_{\text{pop}}(f) = \mathbb{E}_{x \sim \mu} [L(f(x), f^*(x))]$$

↓ approximation (estimate the population risk)

$$\min_f R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^N L(f(x_i), f^*(x_i))$$

$$R_{\text{pop}}(f) \approx R_{\text{emp}}(f)$$

↓  
ground-truth value

↓  
RV estimator w.r.t training data

↓  
BIASED

★

unbiased estimator

→ testing set split

→ Cross Validation Estimator

6. K-class classification Algo

a) Binary Classification

$$\mathcal{H} = \{f: f(x) = \text{sigmoid}(\omega^\top \phi(x)), \omega \in \mathbb{R}^D\}$$

b) K-class Classification

$$\mathcal{H} = \{f: f(x) = \text{softmax}(W\phi(x)), W \in \mathbb{R}^{K \times D}\}$$

Q: How to train?

A: change the objective (KL-divergence for distribution)

$$(y_i - g(W\phi(x_i)))^2$$

$$g: \mathbb{R}^K \rightarrow \{1, 2, \dots, K\}$$

Unreasonable

Disadvantage for MSE Loss

- 1) Data is nominated, not ordinal
- 2) Hard to optimize (calculate gradient)

$y_i = 0, 1, 2, 3, \dots$   
for different class (ordinal)

### Strategy

- 1) one-hot encoding for  $y_i$
- 2) use softmax function to generate distribution

$$\mathcal{H} = \{f: f(x) = \text{softmax}(W\phi(x)), W \in \mathbb{R}^{K \times D}\}$$

### 3) Loss Function Design

$$\begin{cases} L(y, y') = 0 & \text{if } y' = y \\ L(y, y') \approx 0 & \text{if } y' \approx y \end{cases}$$

(distance measure for distribution)

Good choice is : KL-divergence

$$L(y, y') = - \sum_{k=1}^K y_k \log y'_k$$

$$\begin{cases} y_k \rightarrow \text{ground-truth} \\ y'_k \rightarrow \text{prediction} \end{cases}$$

$$= KL(y, y')$$