

Lecture 1: Basics of ML.

⇒ Theme: theoretical foundations of ML.

↓
feature learning

↓
AI safety.

⇒ Importance: We want to understand how & why ML works.

- Train a model \Rightarrow test on a new task
 - \rightarrow knowledge doesn't transfer
 - \rightarrow not enough data
 - \rightarrow computational issue
- } fundamental statistical problems.
- \rightarrow optimization issue.

(objective function is flawed)

- guide technical decision
- reduce trial and error
- forecast outcomes and risks.
- inspire new methods.

esp. useful
demand / supplies resource
is large

⇒ Structure and logic of this course.

ML (AI in general)

Chatbot.

text

many success stories.

Medical diagnosis

image

Alpha-Go.

RL environments'

specific structures.

We don't care about data modalities, have a general view high-dimensional samples, generated from a certain distribution.

$$X = \{x_1, \dots, x_n\}, \quad x_i \sim P_X.$$

$$Z = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad x_i, y_i \sim P_{X,Y}.$$

Basics of ML. consists of the following elements.

- data (MNIST, CIFAR, IMAGENET)
New era (AGI)
wiki, crawled from the Internet.
- loss function (measuring the difference of model's prediction vs y_i : l_2 , cross-entropy, etc.
next word prediction.
 $-P(x_{t+1} | x_1, \dots, x_t)$.
- model: (linear/affine, kernel, neural networks.
transformer, Bert.
- optimization alg. (SGD, Adam, adafactor, etc.)

• risks on AI safety

Security: adversarial attack
poisonous attack

privacy: reidentification.

clean data poisonous attack.

copyright

chatbot (finetuned w/

your dialog.

alignment.

⇒. Supervised learning.

Input data $\{(x_i, y_i)\}_{i=1}^n$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$.
input data, feature label response.
or \mathbb{R}^c .

Goal: Find prediction function $f_\theta: \mathbb{R}^d \rightarrow \mathbb{R}$ ($\mathbb{R}^d \rightarrow \mathbb{R}^c$ for multi class).
such that $f_\theta(x_i) \approx y_i$, $\forall i$.
↑
parametric model.

Empirical Risk Minimization:

$$\hat{R}_n := \min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i)$$

↓ concentrate

population.

$$R = \mathbb{E}_{(x,y) \sim P_{x,y}} \ell(f_\theta(x), y).$$

Hope:

f_θ s.t.
 $\hat{R}_n(f_\theta)$ small

↓

$R(f_\theta)$ small.

⇒ More types of ML tasks:

Traditional types of ML: } supervised learning
unsupervised learning, ↓
Reinforcement learning,

Spectrum between supervised \Rightarrow unsup. L :

More labeled
data involved

unsupervised L .
self-supervised.
meta-learning.
domain generalization.
semi-supervised learning
supervised learning.

feature learning.

L : labeled data. U : unlabeled data.

S : source data.

T : target data.

e : # environments or tasks.

Learning tasks	data that model trained on	data model is tested on.
self-supervised. (uns)	$\underline{U^S}$, L^T few-shot	U^T
meta-learning.	$\underline{L_1^S, L_2^S, \dots, L_e^S}, \underline{L^T}$	U^T
domain generalization.	$\underline{L_1^S, L_2^S, \dots, L_e^S}$	U^T
semi-supervised learning	L, U	U
supervised learning.	L	U
reinforcement learning	source environment	target environment.

⇒ Theoretical groundings of ML algorithms in general:

* objective function

supervised	unsupervised	privacy attack.
$\sum_{i=1} l(\underbrace{f_{\theta}(x_i)}_{\Delta}, y_i)$	$\min_x \ \underbrace{M}_{\Delta} - \underbrace{XX^T}_{\Delta} \ _F^2$	$\min_x l(\underline{Q}, \underline{F(x)}) + \text{prior}(x)$

LT requires studying the two aspects:

statistics



whether the optimal solution generalizes.

when is the optimal of (*) close to the ground truth.

$$\sum_{(i,j) \in \Omega} \| M_{ij} - (XY)_{ij}^* \|^2_F$$

optimization:



whether our algorithm can find a global minimum.

together they form Learning theory.