

(Deep) Learning Theory.

Instructor: Andrea Agazzi

andrea.agazzi@unipi.it

Motivation: neural networks have emerged as a powerful tool to solve complex problems, and many of the successes of the ongoing AI revolution would be unthinkable without such tool.

Examples include :

- Deepfakes
- Alpha GO
- Snapchat filters
- Ticks / Chatbots

But also :

- Medical diagnosis (data)
- Self-driving cars / robotics
- Advertising (data)
- Translation
- Autopilots

List of links to videos / tutorials available on class website.

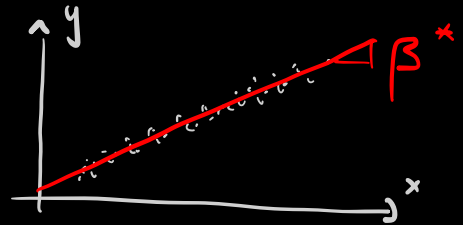
⇒ aim to develop a mathematical theory explaining this success.

Mathematical preliminaries: statistical learning theory

Example (linear regression) assume we have measured the price of texas crude oil and the price of diesel, giving rise to a database

$$D_n = \{(x_i, y_i)\}_{i=1}^n \quad \begin{array}{l} \xrightarrow{\text{crude \$}} \\ \xrightarrow{\text{diesel \$}} \end{array}$$

We can plot the dataset and plot the linear regression



Then, for a new day, measuring the price of crude x_{new} we can predict

$$y_{\text{new}} \approx \beta^* x_{\text{new}}$$

Question: what can I say about my prediction?
is it going to be accurate?

We now introduce a setting where this question can be formulated

Classical Supervised Learning problem:

Given: a dataset/training set $D_n(\underline{P})$ that we aim to "learn" (the world).
 $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} \quad i \in \{1, \dots, n\}$

where:

- examples
- cat pics
 - text
 - car commands
 - code
 - music
- what you have \rightarrow
- $x_i \in \mathcal{X}$ are the inputs / features / covariates / dependent variables / predictors
- what you predict \rightarrow
- $y_i \in \mathcal{Y}$ are the outputs / labels / responses / independent variables / outcomes

throughout this class we will assume that

$$(x_i, y_i) \stackrel{\text{iid}}{\sim} P(dx, dy) \in \mathcal{M}_+^1(\mathcal{X} \times \mathcal{Y})$$

\rightarrow unknown!

and disintegrate $P(dx, dy) = P_{y|x}(x, dy) \cdot P_x(dx)$

Note: We distinguish between:

- $\mathcal{Y} \subseteq \mathbb{R}^d \rightarrow$ regression
- $\mathcal{Y} \subseteq \mathcal{Z} \rightarrow$ classification

Example (simple linear regression): $\mathcal{X} = \mathbb{R}^d, \mathcal{Y} = \mathbb{R}$

$$P_x \text{ unknown}, P_{y|x} = \mathcal{N}(\beta^* x, 1)$$

in other words $y_i = \beta^* x_i + \varepsilon_i, \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$

Example (Classification): $\mathcal{X} = \mathbb{R}^d, \mathcal{Y} = \{-1, 1\}$

$$P_x \text{ unknown}, P_{y|x}(x, y) = (1 - \pi(x))^{\frac{1-y}{2}} \pi(x)^{\frac{1+y}{2}} \quad \pi: \mathcal{X} \rightarrow [0, 1]$$

$\rightarrow P_{(y=1|x)}$

Decision theory (what the engineer decides).

Find $f: X \rightarrow Y$ such that for a new $(x_{\text{new}}, y_{\text{new}}) \sim P$

$$f(x_{\text{new}}) \approx y_{\text{new}} \implies f(x) \text{ approximates well } y?$$

To perform/evaluate a prediction, we have to fix:

- a measure of performance (loss function).

$$\ell: Y \times Y \rightarrow \mathbb{R}$$

returning the cost $\ell(y, y')$ of predicting $y' \in Y$ when true label is $y \in Y$. This induces

The expected/population Risk (generalization error)
expected loss

$$R(f) = \int_{X \times Y} \ell(y_{\text{new}}, f(x_{\text{new}})) p(dx_{\text{new}}, dy_{\text{new}})$$

a proxy for this quantity is the empirical risk (training error)

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

Example (regression): $\ell(y, y') = \|y' - y\|^2$ (MSE).

$$\Rightarrow R(f) = \int_{x \times y} \|f(x) - y\|^2 P(dx, dy)$$

$$\hat{R}_n(f) = \frac{1}{n} \sum \|f(x_i) - y_i\|^2 \quad (x_i, y_i) \stackrel{\text{iid}}{\sim} P$$

Example (linear regression): for a given $f_\beta: x \mapsto \beta \cdot x$ we

have:

$$R(\beta) = \int_{\mathbb{R}^2} (\beta \cdot x - y)^2 \cdot e^{-\frac{(y - \beta \cdot x)^2}{2}} \cdot \frac{1}{\sqrt{2\pi}} \cdot P_x(dx) dy$$

exercise

$$= \int [(\beta - \beta^*) \cdot x]^2 P_x(dx) + 1$$

Example: (classification, 0-1) $\ell(y, y') = \mathbb{1}(yy' \leq 0)$
 $\mathbb{1}(y \neq y')$

$$R(f) = \int \mathbb{1}(f(x) \neq y) P(dx, dy) = P(f(x) \neq y).$$

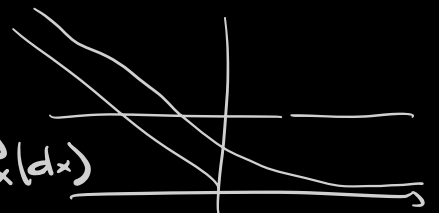
Example (classification, logistic) $\ell(y, y') = \log(1 + \exp(-yy'))$

$$R(f) = \sum_{y \in \{-1, 1\}} \ell_y(1 + \exp(-y f(x))) \underbrace{P_y(x, y)}_{\text{Bernoulli}} P_x(dx)$$

note that sometimes in classification tasks we take $f: X \rightarrow \mathbb{R}$
 and predict with $\text{sign } f(x)$. (the larger $\|f(x)\|$ the more 'sure')

Exercise: setting $\hat{p}(x, y) = \frac{1}{1 + e^{-f(x)}}$ then

$$R(f) = \int \sum_y P_y(x, y) \log(\hat{p}(x, y)) P_x(dx)$$



- a model or hypothesis class

$$\mathcal{F} \subseteq \{f: \mathcal{X} \rightarrow \mathcal{Y}\}$$

the class of functions where we search for the "optimal" approximator.

Usually we work with parametric models, for a param. space Θ :

$$\mathcal{F}_\Theta = \{f_\vartheta: \mathcal{X} \rightarrow \mathcal{Y}, \vartheta \in \Theta\}.$$

→ classes of functions indexed by a parameter ϑ .

Example (linear regression) $\Theta = \mathbb{R}^d$, $\mathcal{F}_\Theta = \{x \mapsto \beta \cdot x, \beta \in \Theta\}$

Example (logistic regression) $\Theta = \mathbb{R}^d$, $\mathcal{F}_\Theta = \{(1 + e^{-\beta \cdot x})^{-1}, \beta \in \Theta\}$

Example (neural networks): (next class).

- an optimization algorithm.

$$A: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{F}$$

$$\mathcal{D}_n(\mathcal{P}) \mapsto A(\mathcal{D}_n)$$

often, such algorithms aim to minimize \hat{R} on \mathcal{F}_Θ

This strategy is called Empirical Risk Minimization (ERM).

Example (linear regression):

Writing $\hat{R}(\beta) = \frac{1}{n} \sum_{j=1}^n (\beta \cdot x_j - y_j)^2 =$

Notation:

$$X = \begin{pmatrix} x_{1,1} & \dots & x_{1,d} \\ \vdots & & \vdots \\ x_{n,1} & \dots & x_{n,d} \end{pmatrix} \in \mathbb{R}^{n \times d} \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

$$= \frac{1}{n} \|X \cdot \vec{\beta} - \vec{y}\|_2^2$$

we find a minimizer by setting

$$D_{\beta} \hat{R}(\beta) = \frac{1}{n} X^T (X \cdot \vec{\beta} - \vec{y}) = 0$$

$$\implies \mathcal{A}(X, \vec{y}) = (X^T X)^{-1} X^T \vec{y}$$

Summary: for a given $P, \mathcal{C}, \mathcal{F}_0, \mathcal{A}$ we study

$$\Delta(\mathcal{D}_n) = \mathcal{R}(\mathcal{A}(\mathcal{D}_n)) - \inf_{f: X \rightarrow Y} \mathcal{R}(f)$$

for instance: $\mathbb{E}_{\mathcal{D}_n}(\Delta(\mathcal{D}_n)) \xrightarrow{n \rightarrow \infty} 0 \quad ?$

$$\Delta(\mathcal{D}_n) \xrightarrow[n \rightarrow \infty]{P} 0 \quad ?$$

(...)

Risk decomposition

Often, studying $\Delta(\mathcal{D}_n)$ is too hard (many moving parts)
so we decompose it into a sum of terms:

$$\underbrace{R(f_{\hat{\theta}}) - \inf_{f: X \rightarrow Y} R(f)}_{\text{population risk}} = \underbrace{R(f_{\hat{\theta}}) - \inf_{\theta \in \Theta} R(f_{\theta})}_{\text{estimation err.}} + \underbrace{\inf_{\theta \in \Theta} R(f_{\theta}) - \inf_{f: X \rightarrow Y} R(f)}_{\text{approximation err.}}$$

ANALYSIS (F)

typically $\xrightarrow[n \rightarrow \infty]{|\Theta| \rightarrow \infty} 0$

typically $\xrightarrow{|\Theta| \rightarrow \infty} 0$

we can further decompose the estimation error as:

$$\underbrace{R(f_{\hat{\theta}}) - R(f_{\theta^*})}_{\text{estimation}} = \underbrace{R(f_{\hat{\theta}}) - \hat{R}(f_{\hat{\theta}})}_{\text{gen / concentr.}} + \underbrace{\hat{R}(f_{\hat{\theta}}) - \hat{R}(f_{\theta^*})}_{\text{empirical optimiz}} + \underbrace{\hat{R}(f_{\theta^*}) - R(f_{\theta^*})}_{\text{gen}}$$

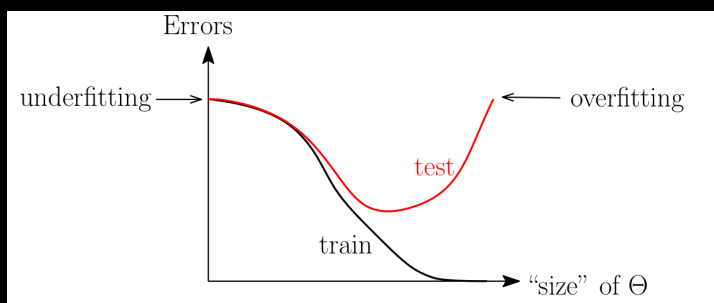
STA (\mathcal{D}_n)

OPT ($\mathcal{F} + \mathcal{F}$)

→ each part of this can "compete" a different element and bound

Time permitting. $y_i = \theta_*^T x_i + \varepsilon_i$ $\mathbb{E}(\varepsilon_i) = 0$ $\text{Var}(\varepsilon_i) = \sigma^2$

$$\Rightarrow R(\theta) = \mathbb{E}(\|X\theta_* + \varepsilon - X\theta\|^2) = \sigma^2 + \frac{1}{n}(\theta - \theta_*)^T X^T X (\theta - \theta_*)$$



$$\begin{aligned} \Rightarrow \mathbb{E}(R(\theta)) - R^* &= \mathbb{E}(\|\hat{\theta} - \mathbb{E}(\hat{\theta}) - \theta_*^*\|_{\Sigma}^2) \\ &= \underbrace{\mathbb{E}(\|\hat{\theta} - \mathbb{E}(\hat{\theta})\|_{\Sigma}^2)}_{\text{Variance}} + 0 + \underbrace{\|\mathbb{E}(\hat{\theta}) - \theta_*^*\|_{\Sigma}^2}_{\text{Bias}} \end{aligned}$$

Program: - General, finite width neural networks

- unsatisfactory
- Definition
 - Expressivity
 - Training
 - Generalization bounds when $n \rightarrow \infty$

- ∞ -width NN 2 regimes: NTK, MF

- more satisf. but still not done!
- NTK:
 - linear behavior of training
 - generalization analysis
 - MF:
 - nonlinear training behavior
 - optimality
 - Differences between two regimes.
 - Extensions to RL, RNNs, generative

Logistics: - website contains up-to-date info on class

- notes published on website
- exam: presentation for those who come to class, the oral exam.
- THU class: attending (not this week).
- discuss open problems.