

# Theoretical Machine Learning

Lectured by [Zhihua Zhang](#)

L<sup>A</sup>T<sub>E</sub>Xed by [Chengxin Gong](#)

February 26, 2024

## Contents

<a href="#">1</a>	<a href="#">简介</a>	<a href="#">2</a>
<a href="#">2</a>	<a href="#">统计决策理论</a>	<a href="#">2</a>

# 1 简介

- 机器学习的主要任务: 生成、预测、决策. 生成:  $X_1, \dots, X_n \sim F$ , 推断分析  $F$ , 无监督学习, GAN, GPT,  $\dots$ . 预测: 数据对  $(X^{(1)}, Y^{(1)}), \dots, (X^{(n)}, Y^{(n)})$ ,  $X^{(i)} \in \mathbb{R}^d$  输入变量,  $f: \mathcal{X} \rightarrow \mathcal{Y}, x \in \mathcal{X}, y \in \mathcal{Y}$ , 归因, 有监督学习. 决策: 强化学习, Agent  $\leftarrow$  action, state, reward  $\rightarrow$  环境.
- 求解问题的途径: 参数/非参数, 频率 (MLE)/贝叶斯.
- 误差模型: 有监督:  $X = (X_1, \dots, X_d)^T \in \mathbb{R}^d$ , 回归:  $Y \in \mathbb{R}$ ; 分类:  $Y \in \{0, 1\}(\{-1, 1\}, \{1, \dots, M\}, \{0, 1\}^M)$ ;  $X$  随机, Random design(生成模型),  $Y = g(X) + \epsilon \stackrel{\text{or}}{=} g(X, Z), Y^{(i)} = g(X^{(i)}, Z^{(i)})$ ;  $X$  固定  $X = x$ , Fixed design(判别模型),  $Y^{(i)} = g(x^{(i)}, Z^{(i)})$ . 无监督:  $X = g(Z)$ (因子模型:  $X = AZ + \epsilon, Z \in \mathcal{N}(0, 1), \epsilon \sim \mathcal{N}(0, \Sigma)$ ).

# 2 统计决策理论

- Consider a state space  $\Omega$ , data space  $\mathcal{D}$ , model  $\mathcal{P} = \{p(\theta, x)\}$ , action space  $\mathcal{A}$ . Loss function:  $\mathcal{L}: \Omega \times \mathcal{A} \rightarrow [-\infty, +\infty]$ , measurable, nonnegative. A measurable function  $\delta: \mathcal{D} \rightarrow \mathcal{A}$  is called a nonrandomized decision rule. Risk function is defined as  $\mathcal{R}(\theta, \delta) = \int \mathcal{L}(\theta, \delta(x)) dP_\theta(x) = \mathbb{E}_\theta \mathcal{L}(\theta, \delta(X))$ . Randomized decision: for each  $X = x$ ,  $\delta(x)$  is a probability distribution:  $[A|X = x] \sim \delta_x$ . Risk function for  $\delta$ :  $\mathcal{R}(\theta, \delta) = \mathbb{E}_\theta \mathcal{L}(\theta, A) = \mathbb{E}_\theta \mathbb{E}_a \mathcal{L}(\theta, A|X) = \iint \mathcal{L}(\theta, a) d\delta_x(a) dP_\theta(x)$ .
- Example [参数估计]:  $\theta \in \Omega, \mathcal{A} = \Omega, \mathcal{L}(\theta, a) = \|\theta - a\|_2^2 \stackrel{\text{or}}{=} \|\theta - a\|_p^p (p \geq 1) \stackrel{\text{or}}{=} \int \log \frac{P_\theta(x)}{P_a(x)} P_\theta(x) dm(x) (\text{KL})$ .  $\mathcal{R} = \text{Var}(a) + \text{bias}^2(a)$ . Bregmass loss:  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}$  describe any strictly convex differentiable function. Then  $\mathcal{L}_\phi(\theta, a) = \phi(a) - \phi(\theta) - (\phi - a)^T \nabla \phi(a)$ .
- Example [Testing]:  $\mathcal{A} = \{0, 1\}$  with action “0” associated with accepting  $H_0: \theta \in \Omega_0$  and “1”:  $H_1: \theta \in \Omega_1$ .  $\delta_x$  is a Bernolli distribution.  $\mathcal{L}(\theta, a) = I\{a = 1, \theta \in \Omega_0\} + I\{a = 0, \theta \in \Omega_1\}$ . Risk  $\mathcal{R}(\theta, \delta) = \mathbb{P}_\theta(A = 1)1_{\theta \in \Omega_0} + \mathbb{P}_\theta(A = 0)1_{\theta \in \Omega_1}$ .
- A decision rule  $\delta$  is called inadmissible if a competing rule  $\delta^*$  such that  $\mathcal{R}(\theta, \delta^*) \leq \mathcal{R}(\theta, \delta)$  for all  $\theta \in \Omega$  and  $\mathcal{R}(\theta, \delta^*) < \mathcal{R}(\theta, \delta)$  for at least one  $\theta \in \Omega$ . Otherwise,  $\delta$  is admissible.
- The maximum risk  $\bar{\mathcal{R}}(\delta) = \sup_{\theta \in \Omega} \mathcal{R}(\theta, \delta)$  and the Bayes risk  $r(\Lambda, \delta) = \int \mathcal{R}(\theta, \delta) d\Lambda(\theta) = \int \mathcal{L}(\theta, \delta) d\mathbb{P}(x, \theta)$  ( $\Lambda(\theta)$  is a prior). A decision rule that minimizes the Bayes risk is called a Bayes rule, that is,  $\hat{\delta}: r(\Lambda, \hat{\delta}) = \inf_\delta r(\Lambda, \delta)$ . Minimax rule  $\delta^*: \sup_{\theta \in \Omega} \mathcal{R}(\theta, \delta^*) = \inf_\delta \sup_{\theta \in \Omega} \mathcal{R}(\theta, \delta)$ .
- If risk functions for all decision rules are continuous in  $\theta$ , if  $\delta$  is Bayesian for  $\Lambda$  and has finite integrated risk  $r(\Lambda, \delta) < \infty$ , and if the support of  $\Lambda$  is the whole state space  $\Omega$ , then  $\delta$  is admissible.
- $p(\theta|x) = \frac{p_\theta(x)\lambda(\theta)}{\int p_\theta(x)\lambda(\theta)d\theta} := \frac{p_\theta(x)\lambda(\theta)}{m(x)}$ . Define the posterior risk of  $\delta$ :  $r(\delta|X = x) = \int \mathcal{L}(\theta, \delta(x)) d\mathbb{P}(\theta|x)$ . The Bayes risk  $r(\Lambda, \delta)$  satisfies that  $r(\Lambda, \delta) = \int r(\delta|x) dM(x)$ . Let  $\hat{\delta}(x)$  be the value of  $\delta$  that minimizes  $r(\delta|x)$ . Then  $\hat{\delta}$  is the Bayes rule.
- Application to supervised learning. Case 1: Regression.  $(X, Y) \in \mathcal{X} \times \mathcal{Y}, f: \mathcal{X} \rightarrow \mathcal{Y}, \mathcal{A} = \Omega = \mathcal{Y}, \mathcal{D} = \mathcal{X}, \delta = f, \mathcal{L}(Y, f(X)) = \|Y - f(X)\|_p^p, p \geq 1$ , risk  $R_f = \iint \mathcal{L}(y, f(x)) d\mathbb{P}(x, y) = \mathbb{E}[\mathcal{L}(Y, f(X))] = \mathbb{E}[\mathbb{E}\mathcal{L}(Y, f(X))|X]$ . When  $p = 2$ ,  $r(f|X = x) = \int \mathcal{L}(y, f(x)) d\mathbb{P}(y|x) = \int |y - f(x)|^2 d\mathbb{P}(y|x)$ . 回归函数  $g(x) := \int y d\mathbb{P}(y|x) \Rightarrow R_f = \mathbb{E}|Y - f(X)|^2 = \mathbb{E}|Y - g(X) + g(X) - f(X)|^2 = \mathbb{E}|Y - g(X)|^2 + \mathbb{E}|g(X) - f(X)|^2 \geq \mathbb{E}|Y - g(X)|^2$ .
- Case 2: Pattern classification.  $Y \in \{0, 1\}, p_0 = P(Y = 0), p_1 = P(Y = 1) = 1 - p_0, \mathbb{E}[\mathcal{L}(Y, f(X))] = P(Y \neq f(X))$ . The Bayesian rule (predictor) is given by  $f(x) = 1\{\mathbb{P}(Y = 1|X = x) \geq \frac{\mathcal{L}(1,0) - \mathcal{L}(0,0)}{\mathcal{L}(0,1) - \mathcal{L}(1,1)} \mathbb{P}(Y = 0|X = x)\}$ . (Proof:  $\mathbb{E}[\mathcal{L}(Y, f(X))|X = x] = \begin{cases} \mathbb{E}[\mathcal{L}(Y, 0)|X = x] = \mathcal{L}(0,0)\mathbb{P}(Y = 0|X = x) + \mathcal{L}(1,0)\mathbb{P}(Y = 1|X = x) \\ \mathbb{E}[\mathcal{L}(Y, 1)|X = x] = \mathcal{L}(0,1)\mathbb{P}(Y = 0|X = x) + \mathcal{L}(1,1)\mathbb{P}(Y = 1|X = x) \end{cases}$ , 比较大小)
- 联系:  $\mathbb{P}(Y = 1|X = x) = \mathbb{E}(Y|X = x) := g(x)$ (回归),  $f(x) = 1\{g(x) \geq \frac{1}{2}\}$ . Then  $0 \leq \mathbb{P}(\hat{f}(X) \neq Y) - \mathbb{P}(f(X) \neq Y) \leq 2 \int_{\mathcal{X}} |\hat{g}(x) - g(x)| \mu(dx) \leq 2(\int_{\mathcal{X}} |\hat{g}(x) - g(x)|^2 \mu(dx))^{\frac{1}{2}}$ .

- 回到 Case 2.  $f(x) = 1\{\frac{p(x|y=1)}{p(x|y=0)} \geq \frac{p_0(\mathcal{L}(0,1)-\mathcal{L}(0,0))}{p_1(\mathcal{L}(1,0)-\mathcal{L}(1,1))}\}$ , 这与似然比检验 (LRT) 相同: Likelihood  $L(X) := \frac{p(X|Y=1)}{p(X|Y=0)}$ , 形式为  $f(x) = 1\{L(x) \geq \eta\}$ .
- Confusion table:

	$Y = 0$	$Y = 1$
$\hat{Y} = 0$	true negative	false negative
$\hat{Y} = 1$	false positive	true positive

Ture Positive Rate:  $\text{TPR} = \mathbb{P}(\hat{Y} = 1|Y = 1)$ ; False Negative Rate:  $\text{FNR} = 1 - \text{TPR}$ , type II error; False Positive Rate:  $\text{FPR} = \mathbb{P}(\hat{Y} = 1|Y = 0)$ , type I error; True Negative Rate:  $\text{TNR} = 1 - \text{FPR}$ .

- Optimization: maximize TPR subject to  $\text{FPR} \leq \alpha, \alpha \in [0, 1]$ . Randomized rule:  $Q$  return 1 with probability  $Q(x)$  and 0 with probability  $1 - Q(x)$ . Maximize  $\mathbb{E}[Q(x)|Y = 1]$  subject to  $\mathbb{E}[Q(x)|Y = 0] \leq \alpha$ . Suppose the likelihood functions  $p(x|y)$  are continuous. Then the optimal predictor is a deterministic LRT (N-P lemma).