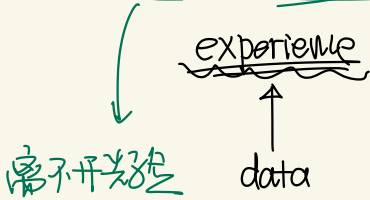


# Summary

## 机器学习理论基础

What is machine learning?



estimation

inference

prediction ...

Machine Learning

limited data

$m$  有限

infinite data

no data

→ probability

+ prior knowledge

performance

$$\text{泛化误差 } R(h) = \mathbb{E}_{x \sim \mathcal{D}} [L(h(x), y)]$$

"概率"

probabilistic / stochastic / agnostic

分类: 0/1

回归:  $(y - y')$

$h \in \mathcal{H} \leftarrow$  假设空间

sampling

PAC (Valiant 1986)

$$\Pr(R(h) \leq \epsilon) \geq 1 - \delta$$

$\mathcal{D}$ : 未知 (Bayesian)

$$S = \{(x_i, y_i)\}_{i=1}^m$$

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m L(h(x_i), y_i)$$

经验风险

Law of Large Number:

$$\hat{R}_S(h) \rightarrow R(h)$$

$m \rightarrow \infty$

有限  $m \rightarrow$  一般

generalization

泛化

$$\mathbb{E}[\hat{R}_S(h)] = R(h)$$

deduction v.s. induction (inductive reasonable)

↑  
exactly

↑  
"probably"

Hume's problem of induction

Assumption

uniformity principle

↳ test, train

i.i.d

minimize  $\hat{R}_S(h) \Rightarrow$  ERM

$$\hat{h}_S^{ERM} = \underset{h \in H}{\operatorname{argmin}} \hat{R}_S(h)$$

H

- ① m 至少为多少时,  $\hat{R}_S(h)$  和  $R(h)$  相差并不会很大 → 样本复杂度
  - ②  $\hat{R}_S(h)$  和  $R(h)$  相差不大而概率至少为多少?
- ⇔ 给定概率,  $\hat{R}_S(h)$  和  $R(h)$  会相差多少?

$m = poly(\frac{1}{\epsilon}, \frac{1}{\delta}, n, size(h))$

$R(h) - \hat{R}_S(h) \leq \text{generalization gap (optimism)}$

Generalization Bound

有限假设空间

$$R(h) \leq \hat{R}_S(h) + \sqrt{\frac{(\log |H| + \log \frac{1}{\delta})}{2m}}$$

w.p.  $\geq 1 - \delta$

$|H|$   
↓

$c(H)$  复杂度

Occam's Razor

简单为王!

- Rademacher Complexity  $R_m(H) = \mathbb{E}_S \left[ \mathbb{E}_\sigma \left[ \sup_h \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \right]$
- Growth Function  $\Pi_H(m) = \max_S \left| \{ (h(x_1), \dots, h(x_m)) : h \in H \} \right|$
- VC-Dimension  $VCdim(H) = \max \{ m : \Pi_H(m) = 2^m \}$

拟合能力

Rademacher Bound:  $\forall g \in G, g(x) \in [0, 1]$

$$\star \mathbb{E}[g(\mathcal{I})] \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + \underbrace{2 R_m(G)} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

$$\Rightarrow R(h) \leq \hat{R}_S(h) + \underbrace{R_m(\mathcal{H})} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

VC dim:  $R(h) \leq \hat{R}_S(h) + \sqrt{\frac{2d \log \frac{1}{\delta}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$   $\hookrightarrow \begin{cases} \text{大: 界松} \\ \text{小: 集合简单} \end{cases}$  trade-off

$d \begin{cases} \rightarrow \text{大: 界松} \\ \rightarrow \text{小: 集合简单} \end{cases}$

Covering Number  $\mathcal{X} \rightarrow \mathcal{Y} \in \mathbb{R} \rightarrow \perp (h(x), y) = (h(x) - y)^2 \leq M^2$

$k = N(\mathcal{H}, \underline{\epsilon})$ : 最少能用多少个半径为  $\underline{\epsilon}$  的球覆盖  $\mathcal{H}$ .  $\mathcal{H}$  bounded

$\exists \{h_1, \dots, h_k\} \subseteq \mathcal{H}$  s.t.  $\forall h \in \mathcal{H}, \exists i \leq k$  s.t.

$$\|h - h_i\|_{\infty} = \max_{x \in \mathcal{X}} |h(x) - h_i(x)| \leq \underline{\epsilon}$$

$\mathcal{H}$ : compact

$\rightarrow N(\mathcal{H}, \underline{\epsilon})$  finite

$$\Pr \left( \sup_{h \in \mathcal{H}} |R(h) - \hat{R}_S(h)| \geq \underline{\epsilon} \right) \leq N\left(\mathcal{H}, \frac{\underline{\epsilon}}{8M}\right) \cdot 2 \exp\left(-\frac{m \underline{\epsilon}^2}{2M^4}\right)$$

Proof: ① Optimism  $L_S = R(h) - \hat{R}_S(h)$

$$\forall h_1, h_2 \in \mathcal{H} \quad \forall S \quad |L_S(h_1) - L_S(h_2)| \leq 4M \cdot \|h_1 - h_2\|_{\infty}$$

$$|L_S(h_1) - L_S(h_2)| \leq |R(h_1) - R(h_2)| + |\hat{R}_S(h_1) - \hat{R}_S(h_2)| \quad |h(x) - y| \leq M$$

$$\begin{aligned} |R(h_1) - R(h_2)| &= \left| \mathbb{E}_{x \sim D} [(h_1(x) - y)^2] - \mathbb{E}_{x \sim D} [(h_2(x) - y)^2] \right| \\ &= \left| \mathbb{E}_{x \sim D} [(h_1(x) - h_2(x)) (h_1(x) + h_2(x) - 2y)] \right| \end{aligned}$$

$$= \sum_{(x,y) \sim (x',y')} \Pr(x,y) \cdot \underline{(h_1(x) - h_2(x))} \cdot (\underline{h_1(x)} + h_2(x) - 2y)$$

$$= \sum_{(x,y) \sim (x',y')} \Pr(x,y) \cdot (\|h_1 - h_2\|_{\infty} \cdot 2M) = 2M \cdot \|h_1 - h_2\|_{\infty}$$

$\{h_1, \dots, h_k\}$

(2)  $\mathcal{H} = \bigcup_{i=1}^k B_i$   $\Pr\left(\sup_{h \in \mathcal{H}} |L_S(h)| \geq \varepsilon\right) \leq \sum_{i=1}^k \Pr\left(\sup_{h \in B_i} |L_S(h)| \geq \varepsilon\right)$

(3)  $\Pr\left(\sup_{h \in B_i} |L_S(h)| \geq \varepsilon\right)$

$$\leq \Pr\left[\bigcap_{i \in [k]} |L_S(h_i)| \geq \frac{\varepsilon}{2}\right] \quad \forall i \in [k]$$

$$\Leftrightarrow \sup_{h \in B_i} |L_S(h)| \geq \varepsilon \Rightarrow |L_S(h_i)| \geq \frac{\varepsilon}{2}$$

$$\Leftrightarrow |L_S(h_i)| < \frac{\varepsilon}{2} \Rightarrow \sup_{h \in B_i} |L_S(h)| < \varepsilon$$

$$|L_S(h)| \leq |L_S(h) - L_S(h_i)| + |L_S(h_i)| \leq 4M \cdot \|h - h_i\|_{\infty} + \frac{\varepsilon}{2}$$

$$\leq 4M \cdot \frac{\varepsilon}{8M} + \frac{\varepsilon}{2} < \varepsilon$$

(4)  $\Pr\left(\sup_{h \in \mathcal{H}} |L_S(h)| \geq \varepsilon\right) \leq \sum_{i=1}^k \Pr\left(|L_S(h_i)| \geq \frac{\varepsilon}{2}\right)$

$$\leq \sum_{i=1}^k \Pr\left(\left|\frac{1}{m} \sum_{j=1}^m (h(x_j) - y_j)^2 - \mathbb{E}\left[\frac{1}{m} \sum_{j=1}^m (h(x_j) - y_j)^2\right]\right| \geq \frac{\varepsilon}{2}\right)$$

$$\leq (k \cdot 2) \exp\left(-\frac{m\varepsilon^2}{2M^4}\right)$$

$$\left[0, \frac{M^2}{m}\right]$$

Covering number:  $N\left(\mathcal{H}, \frac{\varepsilon}{8M}\right)$   $\square$

复杂  $\Rightarrow$  SRM (Structural Risk Minimization)

$$h_S^{\text{SRM}} = \underset{h \in H}{\operatorname{argmin}} \underbrace{\hat{R}_S(h)}_{\text{ERM}} + \underbrace{c(h)}_{\text{正则性}}$$



统计学习

统计学习理论

$$h^* = \underset{h}{\operatorname{argmin}} \underbrace{\hat{R}_S(h)}_{\text{经验风险}} + \underbrace{\lambda c(h)}_{\text{正则项}}$$

trade-off

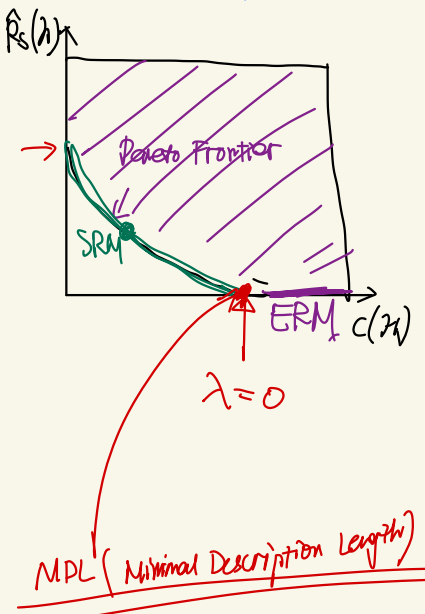
- Soft-margin SVM
- L1-Regularized AdaBoost
- Ridge Regression
- LASSO - -

Agnostic PAC

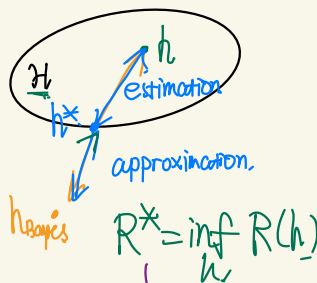
Stochastic learning scenario

bias-variance tradeoff.

$$R = \underbrace{\text{bias}}_{\text{underfitting}} + \underbrace{\text{variance}}_{\text{overfitting}} + \text{noise}$$



SRM:



Probability

Statistics

Inductive

Distribution  $P$  fully specified

What can we say about  $X \sim P$ ?

Data  $D$

Observe data  $X$  from unknown distribution  $P$

What can we conclude about  $P$ ?

i.i.d

Deductive

Estimation

经验误差  $\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{h(x_i) \neq y_i}$

泛化误差  $R(h) = \mathbb{E} [\mathbb{1}_{h(x) \neq y}]$

Generalization bound

$R(h) - \hat{R}_S(h)$  (circled in red)

$\hat{R}_{S, \Phi}(h)$

Optimism

ERM:  $h_S^{\text{ERM}} = \argmin_{h \in H} \hat{R}_S(h)$

$\Downarrow$

SRM:  $h_S^{\text{SRM}} = \argmin_{h \in H} \hat{R}_S(h) + \lambda \cdot C(H)$  ★

统计机器学习

经验误差

0/1-loss

泛化误差

$\|h\|$

复杂度

Rademacher Complexity

Growth Function

VC-dimension

Covering Number

Sauer's Lemma

替代损失

surrogate loss

convex

upper-bound of 0/1-loss

$\Phi$ -loss

核: 支持向量机 (SVM)

min.  $\frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$

s.t.  $y_i (w^T x_i + b) \geq 1 - \xi_i$

$\xi_i \geq 0$

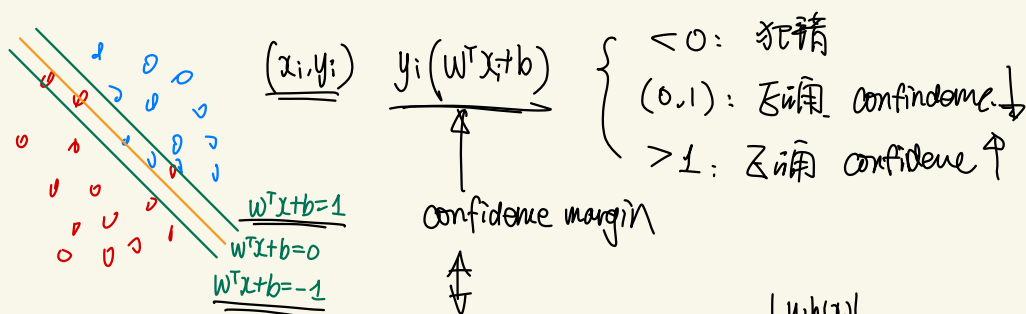
$\Leftrightarrow$

max  $\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j$

s.t.  $\sum_{i=1}^m \alpha_i y_i = 0$

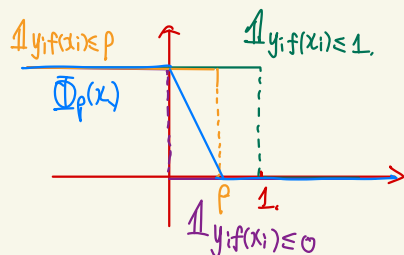
$0 \leq \alpha_i \leq C$

Maximum-Margin



margin

$$\Phi_p(x) = \min\left(1, \max\left(0, 1 - \frac{x}{\rho}\right)\right)$$



①  $\underline{1}_{y_i f(x_i) \leq 0} \leq \Phi_p(x) \leq \underline{1}_{y_i f(x_i) \leq \rho}$

Margin Bound for Binary Classification.

$$R(h) \leq \hat{R}_{S, \rho}(h) + \frac{2}{\rho} \cdot R_m(\mathcal{H}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

w.p.  $\geq 1 - \delta$

② Lipschitz

Hypothesis Space  $\mathcal{H} = \{h: \mathcal{X} \rightarrow \mathcal{Y}\}$

$\downarrow \Phi$  l-lipschitz

Loss Space

$$\mathcal{F} = \{f_h(z) = \Phi(h, z), h \in \mathcal{H}\}$$

$$\mathcal{F} = \Phi \circ \mathcal{H}$$

Talagrand Contraction Lemma

$$R_m(\Phi \circ \mathcal{H}) \leq l \cdot R_m(\mathcal{H})$$

Bounded

$$\|x\|_2 \leq r$$

$$\|w\|_2 \leq 1$$

$$R_m(\mathcal{H}) \leq \sqrt{\frac{r^2 \Lambda^2}{m}}$$

$$\max(1 - y h(x_i), 0) = \xi_i$$

Margin Theory

Kernel Method

$$\text{PDS kernel} \Leftrightarrow K \text{ SPSPD}$$



$$\arg\min_{h \in H} F(h) = \arg\min_{h \in H} \underbrace{G(\|h\|_H)}_{\text{不减} \rightarrow \text{增}} + \underbrace{L(h(x_1), \dots, h(x_m))}_{\text{loss}}$$

$$\Rightarrow \underline{h^* = \sum_{i=1}^m \alpha_i K(x_i, \cdot)}$$

$$\varphi_m(H) = \frac{\Gamma \Lambda}{m} \Rightarrow \underline{R_m(H) \leq \frac{\Delta \sqrt{\Gamma(K)}}{m}}$$

$$\|h\|_H \leq \Gamma \Leftrightarrow K(x_i, x_i) \in \Gamma^2$$

$$R(h) \leq \hat{R}_{S,p}(h) + 2 \cdot \sqrt{\frac{\Gamma^2 \Lambda^2}{m}} + \sqrt{\frac{\log 8}{2m}} \quad \text{w.p.} \geq 1 - \delta$$

## SVM + Kernel Method 核技巧

因地制宜

多分类

$$y = \{1 \dots k\} \quad 1_{h(x) \neq y} \Rightarrow d_H(h(x), y)$$

$$\underline{R(x, y) = h(x, y) - \max_{y' \neq y} h(x, y')}$$

(Multi-class classification)

Multi-Class SVM, Decision Tree  
 OVO, OVR, MVM  $\hookrightarrow$  impurity

Generalization Bound

Rademacher Complexity (Regression)

$$L(y, y') = |y' - y|^2 \quad \hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^n (h(x_i) - y_i)^2$$

VC-dimension

fitting

$\downarrow$   
 pseudo-dimension

fat-dimension (数据相关)

$\epsilon$ -insensitive loss

$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \left( \xi_i + \frac{\xi_i}{2} \right)$$

$$\text{s.t. } (w^T x_i + b) - y_i \leq \epsilon + \xi_i \quad \xi_i \geq 0$$

$$y_i - (w^T x_i + b) \leq \epsilon + \hat{\xi}_i, \quad \hat{\xi}_i \geq 0$$

MSE

Linear Regression

$\rightarrow$  L2: Ridge Regression

$\rightarrow$  L1: LASSO



# Ranking

$$S = \{(\underline{x}_i, \underline{x}'_i, y)\}_{i=1}^m \in \mathcal{X} \times \mathcal{X} \times \{-1, +1\}$$

→ preference  $\{-1, \underline{0}, +1\}$

↓  
scoring

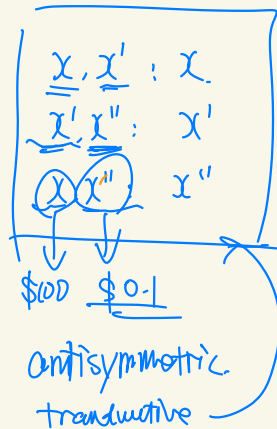
$$h \in H: \begin{aligned} & \underline{h(x_i)} > \underline{h(x'_i)} \\ & \underline{h(x_i)} < \underline{h(x'_i)} \end{aligned}$$

confidence margin:  $\mathbb{E}_{(x, x') \sim D} \left( \underline{y_i (h(x'_i) - h(x_i))} \right) \leq 0$

$$\Rightarrow R(h) = \mathbb{E}_{(x, x') \sim D} \left[ \mathbb{1}_{y_i (h(x'_i) - h(x_i)) \leq 0} \right]$$

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\underline{y_i (h(x'_i) - h(x_i))} \leq 0 \wedge (y_i \neq 0)}$$

$$\hat{R}_{S, \rho}(h) = \frac{1}{m} \sum_{i=1}^m \Phi_{\rho}(\underline{y_i (h(x'_i) - h(x_i))})$$



Ⓜ (Margin Bound for ranking)

$$R(h) \leq \hat{R}_{S, \rho}(h) + \frac{2}{\rho} \left( R_m^{\mathcal{D}_1}(\mathcal{H}) + R_m^{\mathcal{D}_2}(\mathcal{H}) \right) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

fixed  $\rho$

Proof:

$$\mathbb{E} \left[ \underline{\Phi_{\rho}(y_i (h(x'_i) - h(x_i)))} \right] \leq \hat{R}_{S, \rho}(h) + \frac{2}{\rho} R_m(\underline{\Phi_{\rho} \circ F}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

Ⓜ  
R(h)

(x, x', y)  
↓ ↓  
F: loss-space ( $\mathcal{H}$ )  
 $\mathcal{H}$ : hypothesis-space

Talagrand's Lemma:  $R_m(\Phi_{\rho} \circ F) \leq \frac{1}{\rho} R_m(F)$

$$R_m(F) = \frac{1}{m} \mathbb{E}_{S, \sigma} \left[ \left( \sup_{h \in H} \sum_{i=1}^m \sigma_i y_i (h(x'_i) - h(x_i)) \right) \right]$$

Rademacher Variable

$$\leq \frac{1}{m} \mathbb{E}_{S, \sigma} \left[ \sup_{h \in H} \sum_{i=1}^m \sigma_i h(x'_i) + \sup_{h \in H} \sum_{i=1}^m \sigma_i h(x_i) \right]$$

$$= \mathcal{R}_m^{\mathcal{D}_2}(\mathcal{H}) + \mathcal{R}_m^{\mathcal{D}_1}(\mathcal{H})$$

$$\Rightarrow R(h) \leq \widehat{R}_{S,p}(h) + 2 \left( \mathcal{R}_m^{\mathcal{D}_1}(\mathcal{H}) + \mathcal{R}_m^{\mathcal{D}_2}(\mathcal{H}) \right) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad \square$$

$$\hookrightarrow \forall p + \sqrt{(\log \log \frac{2}{\delta})/m}.$$

kernel-based  $\mathcal{R}_m(\mathcal{H}) \leq \sqrt{\frac{r^2 \Lambda^2 / p^2}{m}} \Rightarrow R(h) \leq \widehat{R}_{S,p}(h) + 4 \cdot \sqrt{\frac{r^2 \Lambda^2 / p^2}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$

SVM with Ranking:  $y_i \cdot w(\Phi(x'_i) - \Phi(x_i)) \geq 1 - \xi_i$

$$\Rightarrow \xi_i = \max(1 - y_i(w \cdot (\Phi(x'_i) - \Phi(x_i))), 0)$$

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{s.t. } y_i (w(\Phi(x'_i) - \Phi(x_i))) \geq 1 - \xi_i$$

$$\xi_i \geq 0, \quad \underbrace{\Psi(x'_i, x_i)}_{\forall i \in [m]}$$

$$w \cdot \underbrace{\Psi(x'_i, x_i)}$$

③ Margin Bound

5 N Margin Bound

$\Rightarrow$  Margin Maximization

AdaBoost

$\begin{cases} L_1\text{-margin} & \text{坐标下降} \rightarrow \text{Winnow} \\ L_2\text{-margin} & \text{梯度下降} \rightarrow \text{Perceptron.} \end{cases}$

$$\gamma = \frac{\rho_{\max}}{2}$$

edge

Online Learning

Regret Bound  $\sum_{t=1}^T \ell_t(\hat{y}_t) - \min_{i=1, \dots, N} \sum_{t=1}^T \ell_t(i, y_t)$

$$R_T = \sum_{t=1}^T L(\hat{y}_t, y_t) - \min_{i=1}^N \sum_{t=1}^T L(\hat{y}_{t,i}, y_t)$$

Stability  $|L_{\mathcal{Z}}(h_S) - L_{\mathcal{Z}}(h_{S'})| \leq \beta$

Statistical Learning Theory ★★

统计学习理论

分析问题



损失函数



泛化性



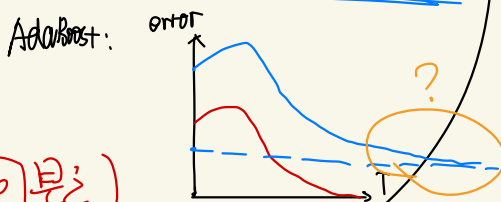
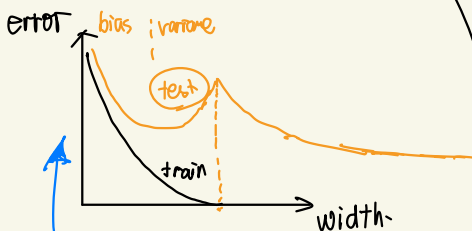
复杂度



算法 (优化问题)

Baum

Double Descent



模型

1. 神经网络

Neural Network Learning:  
Theoretical Foundation  
1998

2. 概率图模型

有向图: Bayesian Network

无向图: Markov Network

检验统计

采样 Sampling

MCMC  
Gibbs (MCM)

Causal Inference