

Covariate Shift and Causal Inference.

Recap last time.

Generalization in RKHS

Bias + variance

$$\mathbb{E}[\text{Error}] \leq \frac{\lambda p^2}{n} \sum_{j=1}^n \frac{\lambda \mu_j}{(\lambda + \mu_j)^2} + \frac{\sigma^2}{n} \sum_{j=1}^n \frac{\mu_j^2}{(\lambda + \mu_j)^2}$$

noise tolerance $\leq \left(\frac{\lambda p^2 + \sigma^2}{n} \right) \sum_{j=1}^n \min\left(1, \frac{\mu_j}{\lambda}\right)$

Problem Set 4 due today
 Pset 5: should be posted
 Thursday

← Could bound
 e.g. for power law
 decay μ_j

Generalization from \hat{p}_n to p_i

$$\hat{S}_n + \lambda_n I \geq S^*$$

← will revisit later

Covariate shift: useful condition for domain adaptation + causal inference

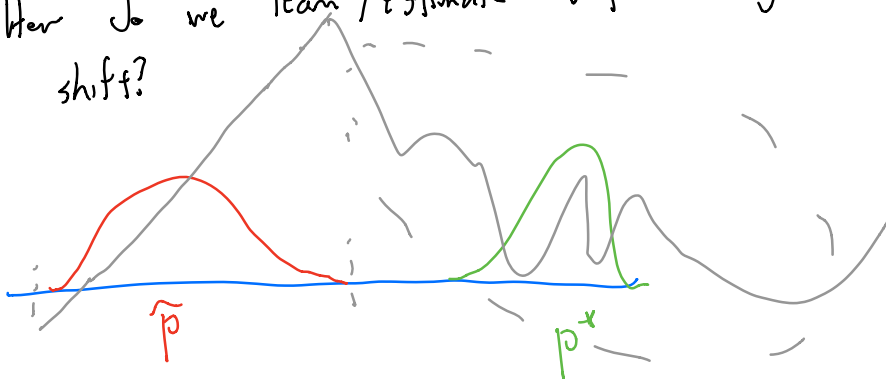
Test distⁿ: $p^*(x, y)$ p_t (target)

Train distⁿ: $\tilde{p}(x, y)$ p_s (source)

Assumption (Covariate Shift)

We assume that $\tilde{p}(y|x) = p^*(y|x) \forall x$.

Question: How do we learn/estimate things taking advantage of
 covariate shift?



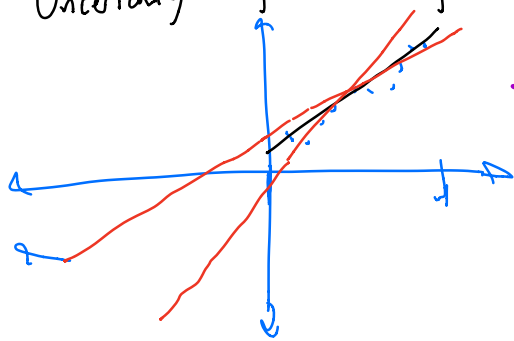
Additional assumptions,

(1) well-specified / realizable assumption

$\theta \in \Theta$

$\hat{p}(y|x) = p^*(y|x) = p_{\theta^*}(y|x)$ for some $\theta^* \in \Theta$

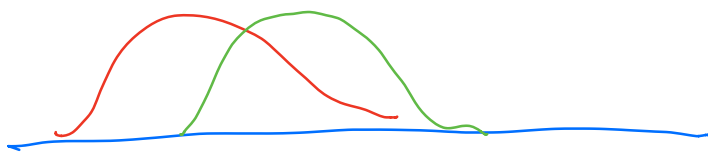
→ Easy to extrapolate from \bar{p} to p^*
BUT. Uncertainty might be higher for p^* than \bar{p}



See this for
linear regression with
 \bar{y}, \bar{y}^*

(2) assume overlap b/t \bar{p} and p^*
 - Use \bar{p} to "simulate" p^*

[This lecture]



Importance Weighting,

Assume $p^*(x)$ and $\bar{p}(x)$ are known.

Traditionally, $(x_1, y_1), \dots, (x_n, y_n) \sim \bar{p}$

$L_n(\theta) = \frac{1}{n} \sum_{i=1}^n l(\theta; x_i, y_i)$

minimize $L_n(\theta)$

traditional
MLE/ERM

$$L_n(\theta) = \tilde{L}(\theta) = \mathbb{E}_{x,y \sim \tilde{p}} [l(\theta; x, y)]$$

$p^* \neq \tilde{p}$: Core about $L^*(\theta) = \mathbb{E}_{x,y \sim p^*} [l(\theta; x, y)]$

Observation: $\mathbb{E}_{p^*} [l(\theta; x, y)] = \mathbb{E}_{\tilde{p}} \left[\frac{p^*(x)}{\tilde{p}(x)} l(\theta; x, y) \right]$ |

assuming covariate shift holds.

Instead of minimizing $L_n(\theta)$, minimize

$$L_n^*(\theta) = \frac{1}{n} \sum_{i=1}^n \frac{p^*(x_i)}{\tilde{p}(x_i)} l(\theta; x_i, y_i)$$

propensity-weighted loss

could be very large!

Pf. $\mathbb{E}_{x,y \sim p^*} [l(\theta; x, y)]$

$$= \int p^*(x, y) l(\theta; x, y) dx dy$$

$$= \int \tilde{p}(x, y) \frac{p^*(x, y)}{\tilde{p}(x, y)} l(\theta; x, y) dx dy$$

$$= \mathbb{E}_{x,y \sim \tilde{p}} \left[\frac{p^*(x, y)}{\tilde{p}(x, y)} l(\theta; x, y) \right]$$

$$= \mathbb{E}_{x,y \sim \tilde{p}} \left[\frac{p^*(x)}{\tilde{p}(x)} \frac{p^*(y|x)}{\tilde{p}(y|x)} l(\theta; x, y) \right]$$

$$= \mathbb{E}_{\tilde{p}} \left[\frac{p^*(x)}{\tilde{p}(x)} l(\theta; x, y) \right] \quad \square$$

Variance of estimator,

$$\mathbb{E}[L_n^*(\theta)] = L^*(\theta) \text{ for any fixed } \theta$$

over draw of training data

$$L_n^*(\theta) = \frac{1}{n} \sum_{i=1}^n \frac{p^*(x_i)}{\tilde{p}(x_i)} \ell(\theta; x_i, y_i)$$

each term is independent and has expectation

$$\mathbb{E}_{\tilde{p}} \left[\frac{p^*(x)}{\tilde{p}(x)} \ell(\theta; x, y) \right] = L^*(\theta)$$

Variance: $\frac{1}{n} \cdot (\text{Variance of individual term})$

$$\text{Var}_{\tilde{p}} \left[\frac{p^*(x)}{\tilde{p}(x)} \ell(\theta; x, y) \right]$$

$$\leq \mathbb{E}_{\tilde{p}} \left[\frac{p^*(x)^2}{\tilde{p}(x)^2} \ell(\theta; x, y)^2 \right]$$

$$\leq B^2 \mathbb{E}_{\tilde{p}} \left[\frac{p^*(x)^2}{\tilde{p}(x)^2} \right]$$

Assume $|\ell| \leq B$
almost surely

$$= B^2 \cdot (D_{\chi^2}(\tilde{p} \| p^*) + 1)$$

$$D_{\chi^2}(p \| q) = \int \frac{(q(x) - p(x))^2}{p(x)} dx$$

$$= \int \frac{q(x)^2}{p(x)} dx - 1$$

Upper Variance upper-bounded in terms of D_{χ^2} .

Understanding χ^2 -divergence,

Claim: $KL \leq \chi^2$

$$\int p(x) \log \frac{p(x)}{q(x)} dx$$

$$= \int p(x) (\log p(x) - \log q(x)) dx$$

$$\begin{aligned}
 & \int p(x) \frac{p(x) - q(x)}{q(x)} dx \\
 &= \int \frac{p(x)^2}{q(x)} dx - \int p(x) dx \\
 &= \int \frac{p(x)^2}{q(x)} dx - 1 = D_{\chi^2}(q \| p). \quad \text{⑩}
 \end{aligned}$$

For any concave $h(x)$, $h(y) - h(x) \leq (y-x)h'(x)$.

χ^2 for Gaussians

$$N(\mu, I), N(\mu', I)$$

$$D_{\chi^2}(N(\mu, I), N(\mu', I)) = \exp(\|\mu - \mu'\|_2^2) - 1$$

→ χ^2 distance is exponential in distance b/w means!!

Connection Causal inference

Estimating treatment effects

Patient: covariates/features: X
treatment : $T \in \{0, 1\}$
outcome : Y

Ex. Vaccines
 T : 1 if got vaccine
0 else

Y : 1 if got COVID
0 else

Goal Estimate treatment effect \rightarrow WRONG

$$\cancel{E[Y | T=1] - E[Y | T=0]}$$

$T=1$: 65+
 $T=0$: 18-64

Solution Potential outcomes framework

$X, T, Y(0), Y(1)$

what "would have"

What "would have" happened if $T=1$
 happened if $T=0$

Only observe $(X, T, Y(T))$

(Average) Treatment Effect:

$$\mathbb{E}_{\bar{p}}[Y(1)] - \mathbb{E}_{\bar{p}}[Y(0)]$$

we'll estimate via reduction to
 covariate shift problem.

Definition (Unconfoundedness)

$$Y(0), Y(1) \perp\!\!\!\perp T \mid X$$

"pre-treatment"
 information

$$\begin{aligned} \mathbb{E}_{\bar{p}}[Y(1)] &= \mathbb{E}_{X \sim \bar{p}} \left[\mathbb{E}_{\bar{p}}[Y(1) \mid X] \right] \\ &= \mathbb{E}_{X \sim \bar{p}} \left[\mathbb{E}_{\bar{p}}[Y(1) \mid X, T=1] \right] \\ &= \mathbb{E}_{X, Y \sim p_1^*} [Y(T)] \end{aligned}$$

where $p_1^*(x, t, y) = \bar{p}(x) \mathbb{I}[t=1] \bar{p}(y \mid x, t=1)$

↳ Key point: \bar{p}, p_1^* satisfy covariate shift

since $\bar{p}(y|x, e) = p^*(y|x, e)$

$$\Rightarrow \mathbb{E}_{p_1^*}[\gamma(T)] = \mathbb{E}_{\bar{p}} \left[\frac{p_1^*(x, e)}{\bar{p}(x, e)} \gamma(T) \right]$$

$$= \mathbb{E}_{\bar{p}} \left[\frac{p_1^*(x) p_1^*(t|x)}{\bar{p}(x) \bar{p}(t|x)} \gamma(T) \right]$$

$$= \mathbb{E}_{\bar{p}} \left[\frac{\mathbb{I}[t=1]}{\bar{p}(t=1|x)} \gamma(T) \right]$$

only involves $\gamma(T)$, hence empirically observable!!

Can do same thing for $p_0^* = \bar{p}(x) \mathbb{I}[t=0] \tilde{p}(y|x, t=0)$
($\gamma_0^*(x, y)$)

$$ATE = \mathbb{E}_{\bar{p}} [\gamma(1) - \gamma(0)]$$

$$= \mathbb{E}_{\bar{p}} \left[\left(\frac{\mathbb{I}[T=1]}{\bar{p}(T=1|x)} - \frac{\mathbb{I}[T=0]}{\bar{p}(T=0|x)} \right) \gamma(T) \right]$$

(X_i, T_i, Y_i)

$$\checkmark \pi(x) = p(t=1|x)$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n \left(\frac{t_i}{\pi(x_i)} - \frac{(1-t_i)}{1-\pi(x_i)} \right) y_i$$

"propensity to treatment"
 inverse propensity weighted estimator

Recipe. Given observational data,

① Fit model to estimate $\pi(x)$

② Estimate ATE via inverse propensity weighting

Next time Fencer recipe called "Jacoby robust" estimation that improves on IPW.