Pattern Information Processing²²¹ Covariate Shift Adaptation

> Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/

Common Assumption in Supervised Learning

Goal of supervised learning: From training samples $\{(x_i, y_i)\}_{i=1}^n$, predict outputs of unseen test samples



222

We always assume

Training and test samples are drawn from the same distribution

$$P_{train}(\boldsymbol{x}, y) = P_{test}(\boldsymbol{x}, y)$$

Is this assumption really true?

Not Always True!

223

- Less women in face dataset than reality.
 More criticisms in survey sampling than reality.
- Sample generation mechanism varies over time.



The Yale Face Database B

Covariate Shift

224

However, no chance for generalization if training and test samples have nothing in common.

$$P_{train}(\boldsymbol{x}, y) \neq P_{test}(\boldsymbol{x}, y)$$

Covariate shift:

- Input distribution changes $P_{train}(\boldsymbol{x}) \neq P_{test}(\boldsymbol{x})$
- Functional relation remains unchanged

$$P_{train}(y|\boldsymbol{x}) = P_{test}(y|\boldsymbol{x})$$

Examples of Covariate Shift ²²⁵

(Weak) extrapolation: Predict output values outside training region



Organization

- 1. Linear regression under covariate shift
- 2. Parameter learning
- 3. Importance estimation
- 4. Model selection



Covariate Shift

To illustrate the effect of covariate shift, let's focus on linear extrapolation



Generalization Error = Bias + Variance

$$\mathbb{E}_{\epsilon} \int \left(\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})\right)^2 p_{test}(\boldsymbol{x}) d\boldsymbol{x}$$
$$= \int \left(\mathbb{E}_{\epsilon} \hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})\right)^2 p_{test}(\boldsymbol{x}) d\boldsymbol{x} \qquad \text{Bias}$$
$$+ \mathbb{E}_{\epsilon} \int \left(\mathbb{E}_{\epsilon} \hat{f}(\boldsymbol{x}) - \hat{f}(\boldsymbol{x})\right)^2 p_{test}(\boldsymbol{x}) d\boldsymbol{x} \qquad \text{Variance}$$



 \mathbb{E}_{ϵ} : expectation over noise

Model Specification

229

Model is said to be correctly specified if

$$\exists \boldsymbol{\alpha}^*, \ \hat{f}(\boldsymbol{x}; \boldsymbol{\alpha}^*) = f(\boldsymbol{x})$$

In practice, our model may not be correct.
 Therefore, we need a theory for misspecified models!

Ordinary Least-Squares

$$\min_{\boldsymbol{\alpha}} \left[\sum_{i=1}^{n} \left(\hat{f}(\boldsymbol{x}_{i}) - y_{i} \right)^{2} \right]$$



Law of Large Numbers

Sample average converges to the population mean:

$$\frac{1}{n} \sum_{i=1}^{n} A(\boldsymbol{x}_{i}) \longrightarrow \int A(\boldsymbol{x}) p_{train}(\boldsymbol{x}) d\boldsymbol{x}$$
$$\boldsymbol{x}_{i} \stackrel{i.i.d.}{\sim} p_{train}(\boldsymbol{x})$$

We want to estimate the expectation over test input points only using training input points $\{x_i\}_{i=1}^n$.

$$\int A(t) p_{test}(t) dt \qquad t \sim$$

 $p_{test}(\boldsymbol{x})$



Importance-weighted average:

$$\begin{array}{l} \frac{1}{n} \sum_{i=1}^{n} \frac{p_{test}(\boldsymbol{x}_i)}{p_{train}(\boldsymbol{x}_i)} A(\boldsymbol{x}_i) \longrightarrow \int \frac{p_{test}(\boldsymbol{x})}{p_{train}(\boldsymbol{x})} A(\boldsymbol{x}) p_{train}(\boldsymbol{x}) d\boldsymbol{x} \\ \\ \boldsymbol{x}_i \stackrel{i.i.d.}{\sim} p_{train}(\boldsymbol{x}) \qquad = \int A(\boldsymbol{x}) p_{test}(\boldsymbol{x}) d\boldsymbol{x} \\ \\ \boldsymbol{t} \sim p_{test}(\boldsymbol{x}) \qquad \text{(cf. importance sampling)} \end{array}$$

Importance-Weighted LS ²³³

$$\min_{\boldsymbol{\alpha}} \left[\sum_{i=1}^{n} \frac{p_{test}(\boldsymbol{x}_i)}{p_{train}(\boldsymbol{x}_i)} \left(\hat{f}(\boldsymbol{x}_i) - y_i \right)^2 \right]$$

 $p_{train}(oldsymbol{x}), p_{test}(oldsymbol{x})$:Assumed strictly positive

- Even for misspedified models, IWLS minimizes bias asymptotically.
- We need to estimate importance in practice.



Organization

- 1. Linear regression under covariate shift
- 2. Parameter learning
- 3. Importance estimation
- 4. Model selection



Importance Estimation

$$w(\boldsymbol{x}_i) = \frac{p_{test}(\boldsymbol{x}_i)}{p_{train}(\boldsymbol{x}_i)}$$

- Assumption: We have training inputs {x_i^{train}}_{i=1}^{n_{train}} and test inputs {x_i^{test}}_{i=1}^{n_{test}}.
 Naïve approach: Estimate p_{train}(x) and p_{test}(x) separately, and take the ratio of the density estimates
- This does not work well since density estimation is hard in high dimensions.

Modeling Importance Function²³⁶

$$w(\boldsymbol{x}) = rac{p_{test}(\boldsymbol{x})}{p_{train}(\boldsymbol{x})}$$

We use a linear model:

$$\begin{split} \widehat{w}(\boldsymbol{x}) &= \sum_{i=1}^{t} \theta_{i} \phi_{i}(\boldsymbol{x}) \quad \theta_{i}, \phi_{i}(\boldsymbol{x}) \geq 0 \\ \text{Test density is approximated by} \\ \widehat{p}_{test}(\boldsymbol{x}) &= \widehat{w}(\boldsymbol{x}) p_{train}(\boldsymbol{x}) \\ \text{Idea: Learn } \{\theta_{i}\}_{i=1}^{t} \text{ so that } \widehat{p}_{test}(\boldsymbol{x}) \text{ well} \\ \text{approximates } p_{test}(\boldsymbol{x}). \end{split}$$

237 **Kullback-Leibler Divergence** $\min_{\{\theta_i\}_{i=1}^t} KL[p_{test}(\boldsymbol{x})||\widehat{p}_{test}(\boldsymbol{x})]$ $\widehat{p}_{test}(\boldsymbol{x}) = \widehat{w}(\boldsymbol{x})p_{train}(\boldsymbol{x})$ $KL[p_{test}(\boldsymbol{x})||\widehat{w}(\boldsymbol{x})p_{train}(\boldsymbol{x})|$ $= \int p_{test}(\boldsymbol{x}) \log \frac{p_{test}(\boldsymbol{x})}{\widehat{w}(\boldsymbol{x}) p_{train}(\boldsymbol{x})} d\boldsymbol{x}$ $= \int p_{test}(\boldsymbol{x}) \log \frac{p_{test}(\boldsymbol{x})}{p_{train}(\boldsymbol{x})} d\boldsymbol{x} \quad \text{(constant)}$ $-\int p_{test}(\boldsymbol{x})\log\widehat{w}(\boldsymbol{x})d\boldsymbol{x}$ (relevant)

Learning Importance Function²³⁸

Thus

$$\min_{\theta_i\}_{i=1}^t} KL[\widehat{w}(\boldsymbol{x})p_{train}(\boldsymbol{x})||p_{test}(\boldsymbol{x})]$$

$$\bigoplus_{\{\theta_i\}_{i=1}^t} \int p_{test}(\boldsymbol{x}) \log \widehat{w}(\boldsymbol{x}) d\boldsymbol{x}$$
(objective function)

Since $\widehat{p}_{test}(\boldsymbol{x}) = \widehat{w}(\boldsymbol{x})p_{train}(\boldsymbol{x})$ is density,

$$\int \widehat{w}(\boldsymbol{x}) p_{train}(\boldsymbol{x}) d\boldsymbol{x} = 1$$
(constraint)

KLIEP (Kullback-Leibler ²³⁹ Importance Estimation Procedure)



Convexity: unique global solution is available
 Sparse solution: prediction is fast!



Model Selection of KLIEP 241

How to choose tuning parameters (such as Gaussian width)?

Likelihood cross-validation:

- Divide test samples $\{ m{x}_i^{test} \}_{i=1}^{n_{test}}$ into \mathcal{X} and \mathcal{X}' .
- Learn importance from \mathcal{X} .
- Estimate the likelihood using \mathcal{X}' .

$$\frac{1}{|\mathcal{X}'|} \sum_{\boldsymbol{x} \in \mathcal{X}'} \log \widehat{w}_{\mathcal{X}}(\boldsymbol{x})$$

This gives an unbiased estimate of KL (up to an irrelevant constant).



Organization

- 1. Linear regression under covariate shift
- 2. Parameter learning
- 3. Importance estimation
- 4. Model selection



Model Selection

Choice of models is crucial:



We want to determine the model so that generalization error is minimized:

$$G = \int \left(\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x}) \right)^2 p_{test}(\boldsymbol{x}) d\boldsymbol{x}$$

Generalization Error Estimation²⁴⁴ $G = \int \left(\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x})\right)^2 p_{test}(\boldsymbol{x}) d\boldsymbol{x}$

Generalization error is not accessible since the target function f(x) is unknown.

Instead, we use a generalization error estimate.



Cross-Validation

245

- Divide training samples into k groups.
- **Train a learning machine with** k = 1 groups.
- Validate the trained machine using the rest.
- Repeat this for all combinations and output the mean validation error.



CV is almost unbiased without covariate shift.
 But, CV is heavily biased under covariate shift!

Importance-Weighted CV (IWC⅔)⁶

When testing the classifier in CV process, we also importance-weight the test error.



IWCV gives almost unbiased estimates of generalization error even under covariate shift

Example of IWCV



- IWCV gives better estimates of generalization error than CV.
- Model selection by IWCV outperforms CV!

Summary

- Covariate shift: input distribution varies but functional relation remains unchanged
- Importance weighting for adaptation.
 - IW least-squares: consistent
 - KLIEP: direct importance estimation
 - IW cross-validation: unbiased
- Further reading:

Sugiyama & Kawanabe Machine Learning in Non-Stationary Environments, MIT Press, 2012

