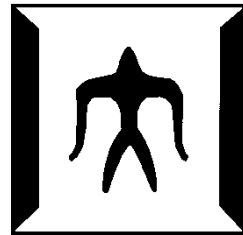


Learning under Class-Prior Change



Christo du Plessis

Department of Computer Science
Tokyo Institute of Technology

Outline

- 1. Motivating Example**
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. Class-prior Change Correction
6. HomeWork

Motivating Example (1)

- A certain medical test for a rare disease has a high accuracy:
 - If the disease is present, the test gives a positive result 90% of the time
 - If the disease is not present, the test gives a negative result 90% of the time
- The disease is quite rare and *only 5% of the population has the disease*
- How likely is the disease if the test result is positive?
 - Around **0.9**, **0.5**, or **0.3**?



Motivating Example (2)

- Frequency of the disease in the population

$$P(A) = 5\% \quad P(\bar{A}) = 1 - P(A) = 95\%$$

A Disease occurs
 \bar{A} Disease does not occur

- Frequency of a positive result when the disease is present:

$$P(B|A) = 90\%$$

B Test is positive
 \bar{B} Test is negative

- Frequency of a negative result when disease is not present

$$P(\bar{B}|\bar{A}) = 90\% \quad P(B|\bar{A}) = 10\%$$

- Frequency that the disease occurs when the test gives a positive answer: $P(A|B)$ **0.9, 0.5, or 0.3?**

Motivating Example (3)

- This can be computed with Bayes' rule

$$\begin{array}{c} \text{Posterior} \\ \downarrow \\ \boxed{P(A|B)} = \frac{P(B|A) \boxed{P(A)}}{P(B|A)P(A) + P(B|\bar{A})P(\bar{A})} \end{array}$$

Prior
↓



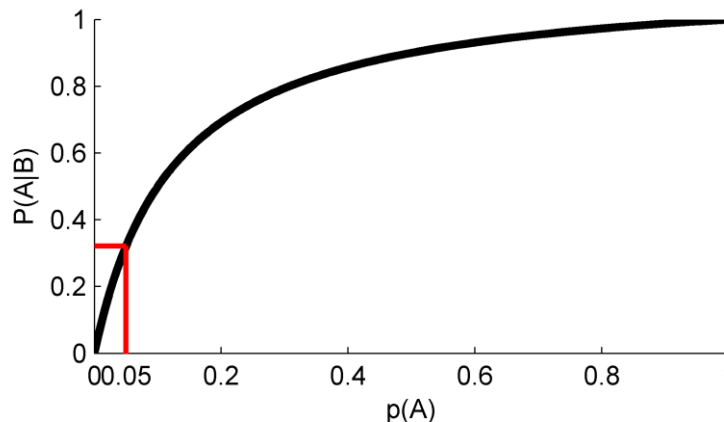
Rev. Thomas Bayes, English statistician and minister

- Substituting the values in the previous slide gives

$$\begin{aligned} P(A|B) &= \frac{0.9 \times 0.05}{0.9 \times 0.05 + 0.1 \times 0.95} \\ &= 32.13\% \end{aligned}$$

Motivating Example (4)

- The result is counterintuitive: much lower than commonly expected
- This is due to the low class prior



$$P(A) = 5\%$$

- **Conclusion:** When doing inference, it is important to take into account the effect of the **class prior**!
- In this lecture, we will discuss the effect of the class prior on classification

Outline

1. Motivating Example
- 2. Classification and Risk**
 - **Risk minimization for Classification**
 - Risk and Class-prior
3. Class-prior Change
4. Class-prior Change Mitigation
5. Class-prior Change Correction
6. Homework

Classification

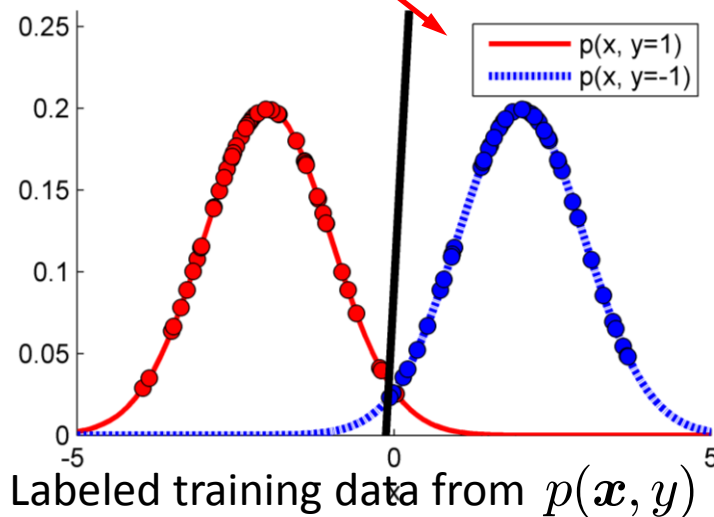
■ Training data:

$$\mathcal{X}_{\text{tr}} := \{\mathbf{x}, y\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y)$$

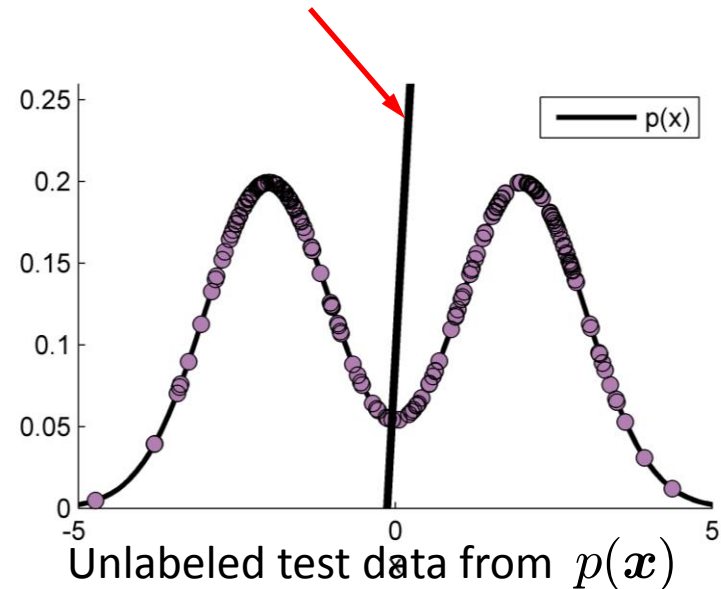
\mathbf{x} Feature
 $y (\in \{-1, 1\})$ Class label
i.i.d: Independently and identically distributed

■ **Goal:** Learn a rule to classify the labeled and unlabeled samples

Decision boundary learned from labeled samples



Learned decision boundary applied on unlabeled samples



■ According to what criterion should the decision boundary be selected?

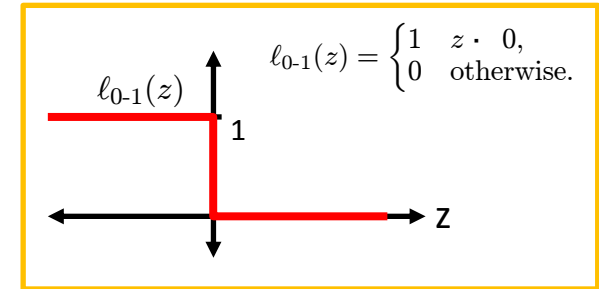
Risk and Classification (1)

- $f(\mathbf{x})$ is a decision boundary: $f(\mathbf{x}) \geq 0$: Class 1
 $f(\mathbf{x}) < 0$: Class -1

- Risk defined as

Misclassification cost **Class prior**

$$R(f) = c_+ p(y = 1) R_1(f) + c_- [1 - p(y = 1)] R_{-1}(f) \quad f: \text{Discriminant}$$



$$R_1(f) = \int \ell_{0-1}(f(\mathbf{x})) p(\mathbf{x}|y = 1) d\mathbf{x}$$

False Negative Rate

$$R_{-1}(f) = \int \ell_{0-1}(-f(\mathbf{x})) p(\mathbf{x}|y = -1) d\mathbf{x}$$

False Positive Rate

$$\ell_{0-1}(f(\mathbf{x})) = \begin{cases} 1 & f(\mathbf{x}) < 0, \\ 0 & f(\mathbf{x}) \geq 0. \end{cases}$$

$$\ell_{0-1}(-f(\mathbf{x})) = \begin{cases} 0 & f(\mathbf{x}) < 0, \\ 1 & f(\mathbf{x}) \geq 0. \end{cases}$$

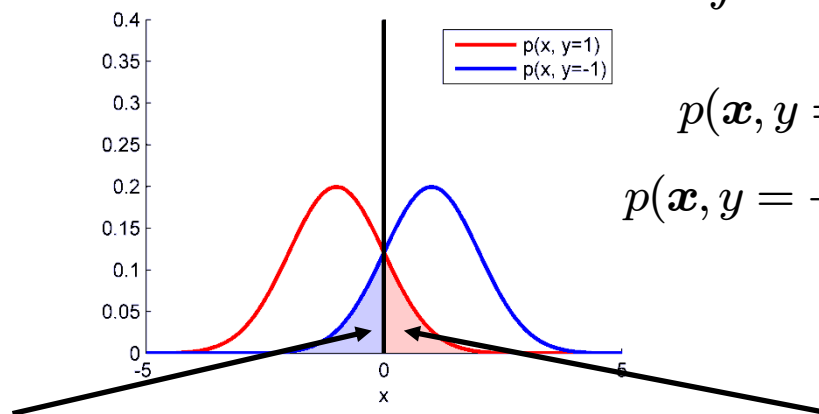
Risk and Classification (2)

10

$$R(f) = p(y = 1)R_1(f) + [1 - p(y = 1)] R_{-1}(f)$$

$$R_1(f) = \int \ell_{0-1}(f(\mathbf{x}))p(\mathbf{x}|y = 1)d\mathbf{x}$$

$$R_{-1}(f) = \int \ell_{0-1}(-f(\mathbf{x}))p(\mathbf{x}|y = -1)d\mathbf{x}$$



$$p(\mathbf{x}, y = 1) = p(\mathbf{x}|y = 1)p(y = 1)$$

$$p(\mathbf{x}, y = -1) = p(\mathbf{x}|y = -1)p(y = -1)$$

$$[1 - p(y = 1)] \int_{f(\mathbf{x}) > 0} p(\mathbf{x}|y = -1)d\mathbf{x} \quad p(y = 1) \int_{f(\mathbf{x}) < 0} p(\mathbf{x}|y = 1)d\mathbf{x}$$

- Decision boundary should minimize the risk

$$f^* = \arg \min_f R(f) \quad R^* = R(f^*)$$

- When $c_+ = c_- = 1$, risk is the misclassification rate
- What is the optimal f^* that minimizes the risk?

Optimal classifier

- When $c_+ = c_- = 1$, the optimal discriminant is

$$f(\mathbf{x}) = \text{sign} [p(y = 1|\mathbf{x}) - p(y = -1|\mathbf{x})]$$

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}|y)p(y)}{\sum_{y'} p(\mathbf{x}|y')p(y')}$$

$p(y)$

Class prior

$p(\mathbf{x}|y)$

Class-conditional density

$p(y|\mathbf{x})$

Posterior

- $f(x)$ is the Bayes-optimal classifier and $R^* = R(f^*)$ is the Bayes-optimal Risk

Outline

1. Motivating Example
- 2. Classification and Risk**
 - Risk minimization for Classification
 - **Risk and Class-prior**
3. Class-prior Change
4. Class-prior Change Mitigation
5. Class-prior Change Correction
6. HomeWork

Bayes Risk vs. Class prior

■ Recall, Bayes Risk is

We will use the symbol π to denote a class prior $p(y = 1)$ from now on

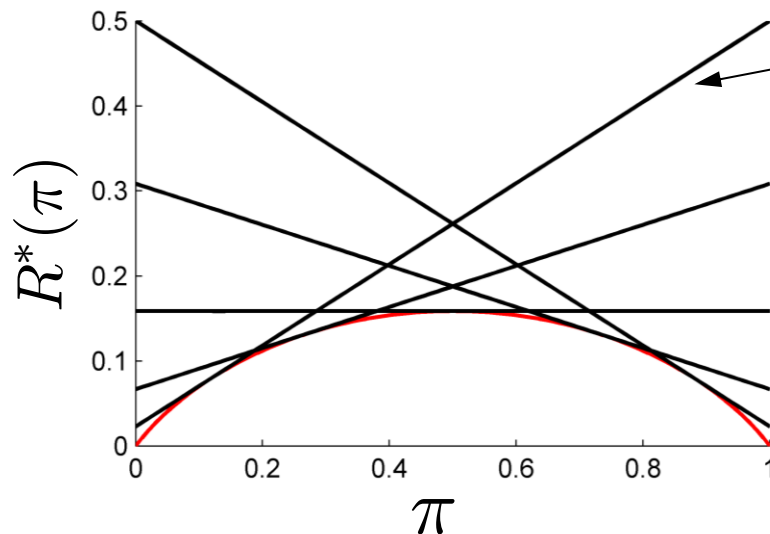
$$R^*(\pi) = \min_f \pi R_1(f) + (1 - \pi) R_{-1}(f)$$

$$R_1(f) = \int \ell_{0-1}(f(\mathbf{x})p(\mathbf{x}|y = 1)d\mathbf{x}$$

False Negative Rate

$$R_{-1}(f) = \int \ell_{0-1}(-f(\mathbf{x})p(\mathbf{x}|y = -1)d\mathbf{x}$$

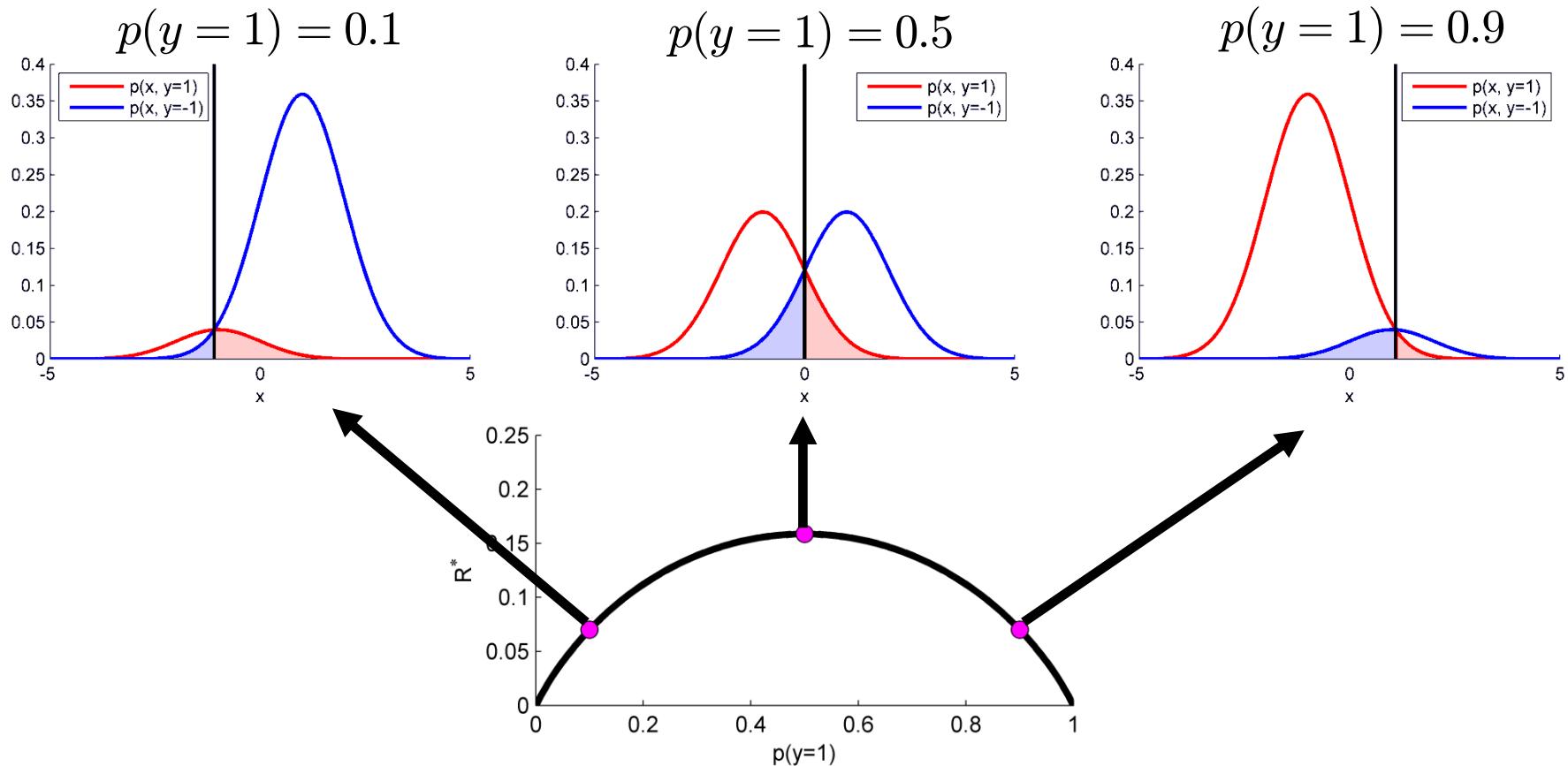
False Positive Rate



Function is *concave* w.r.t. π :

- Minimum of linear functions

Example



- The decision boundary changes when the class prior changes
- Bayes risk is dependent on the class prior

Conclusion: Section 2

- The misclassification rate is a weighted combination of the false negative and false positive rate
 - Weighted by the class priors

- The classifier that minimizes the risk is

$$f(\mathbf{x}) = \text{sign} [p(y = 1|\mathbf{x}) - p(y = -1|\mathbf{x})]$$

- The optimal risk is a *concave function* of the class prior

Outline

1. Motivating Example
2. Classification and Risk
3. **Class-prior Change**
 - **Class-prior Change**
 - Causes of prior Change
 - Dataset shift
 - Selection Bias
 - Class-prior change and Risk
4. Class-prior Change Mitigation
5. Class-prior Change Correction
6. Homework

Class-prior Change

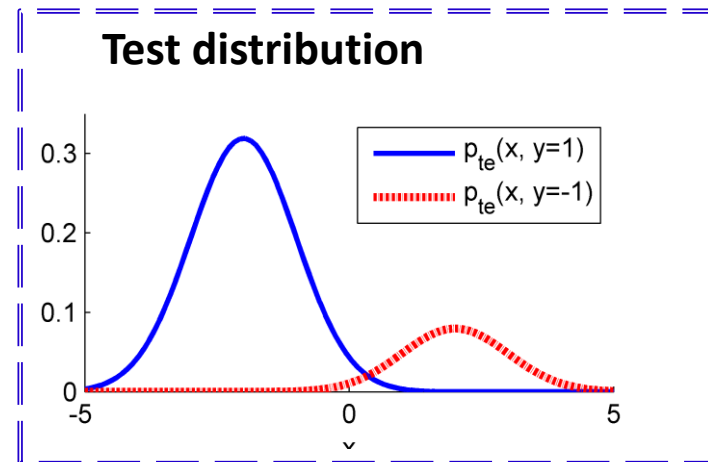
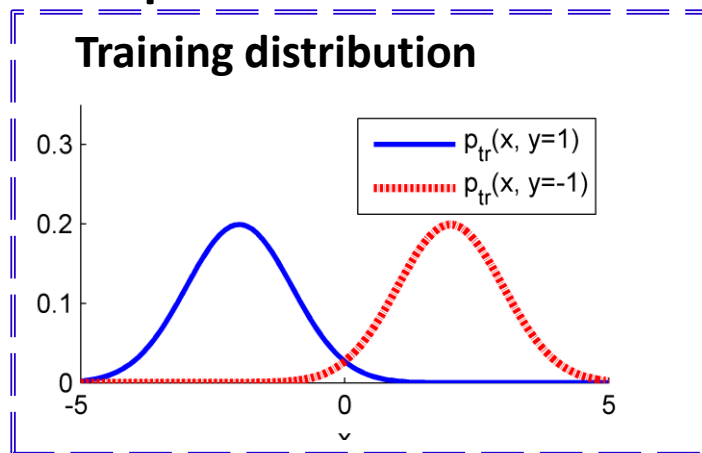
- Class prior between the training and test data differ:

$$p_{\text{te}}(\mathbf{x}, y) = p(\mathbf{x}|y)p_{\text{te}}(y) \quad p_{\text{tr}}(\mathbf{x}, y) = p(\mathbf{x}|y)p_{\text{tr}}(y)$$

$$p_{\text{tr}}(y) \neq p_{\text{te}}(y) \quad \text{Class priors differ}$$

$$p(\mathbf{x}|y) \quad \text{Same class-conditional density}$$

- Example:



Outline

1. Motivating Example
2. Classification and Risk
3. **Class-prior Change**
 - Class-prior Change
 - **Causes of Class-prior Change**
 - **Dataset shift**
 - **Selection Bias**
 - Class-prior Change and Risk
4. Class-prior Change Mitigation
5. Class-prior Change Correction
6. Homework

Why may the dataset change?

■ Dataset shift

- Natural change in the dataset between training and test
- Example: Face images selected in a laboratory compared to the real world

Training dataset:



(Olivetti dataset)

Class balance:

Male: 18/20

Female: 2/20

■ Selection bias (next slide)

Selection bias

- Samples drawn for the test dataset may be drawn in a biased way

- Selection bias model

$$(\mathbf{x}, y, s) \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y, s)$$

\mathbf{x} Feature
 $y(\in \{-1, 1\})$ Class label
 $s(\in \{0, 1\})$ Selection of samples

- When $s = 1$, the sample is in the test set, when $s = 0$, the sample is not in the test set
- Training distribution: $p_{\text{tr}}(\mathbf{x}, y) = p(\mathbf{x}, y)$
- Test distribution: $p_{\text{te}}(\mathbf{x}, y) = p(\mathbf{x}, y | s = 1)$

Selection bias (2)

■ Three possibilities:

1. No selection bias

s is independent of x and y

$$p(s = 1|x, y) = p(s = 1)$$

No selection bias

2. Covariate shift

s is independent of y given x

$$p(s = 1|x, y) = p(s = 1|x)$$

3. Class-prior change

s is independent of x given y

$$p(s = 1|x, y) = p(s = 1|y)$$

Covariate shift

- $p(s = 1|\mathbf{x}, y) = p(s = 1|\mathbf{x})$ implies that

$$p_{\text{tr}}(y|\mathbf{x}) = p_{\text{te}}(y|\mathbf{x})$$

■ Proof:

$$\begin{aligned}
 p_{\text{te}}(y|\mathbf{x}) &= p(y|\mathbf{x}, s = 1) = \frac{p(\mathbf{x}, y, s = 1)}{p(\mathbf{x}, s = 1)} \\
 &= \frac{p(s = 1|\mathbf{x}, y)p(\mathbf{x}, y)}{p(\mathbf{x}, s = 1)} \\
 &= \frac{p(s = 1|\mathbf{x})p(\mathbf{x}, y)}{p(\mathbf{x}, s = 1)} \\
 &= \frac{p(s = 1|\mathbf{x})}{p(s = 1|\mathbf{x})} p(y|\mathbf{x}) = p_{\text{tr}}(y|\mathbf{x})
 \end{aligned}$$

$$\begin{aligned}
 p_{\text{te}}(\mathbf{x}, y) &= p(\mathbf{x}, y|s = 1) \\
 p_{\text{tr}}(\mathbf{x}, y) &= p(\mathbf{x}, y)
 \end{aligned}$$

- Covariate shift occurs in practice!
- Can be corrected for in the semi-supervised setup
- See **lecture 13**

Class-prior Change

- $p(s = 1|\mathbf{x}, y) = p(s = 1|y)$ implies that

$$p_{\text{tr}}(\mathbf{x}|y) = p_{\text{te}}(\mathbf{x}|y)$$

■ Proof:

$$\begin{aligned}
 p_{\text{te}}(\mathbf{x}|y) &= p(\mathbf{x}|y, s = 1) = \frac{p(\mathbf{x}, y, s = 1)}{p(y, s = 1)} \\
 &= \frac{p(s = 1|\mathbf{x}, y)p(\mathbf{x}, y)}{p(s = 1|y)p(y)} \\
 &= \frac{p(s = 1|y)}{p(s = 1|y)} \frac{p(\mathbf{x}, y)}{p(y)} = p(\mathbf{x}|y)
 \end{aligned}$$

$$\begin{aligned}
 p_{\text{te}}(\mathbf{x}, y) &= p(\mathbf{x}, y|s = 1) \\
 p_{\text{tr}}(\mathbf{x}, y) &= p(\mathbf{x}, y)
 \end{aligned}$$

- Class-prior change may be due to selection bias
- We discuss methods to mitigate the effect of class-prior change in the supervised setup
- We discuss a correction for class-prior change in the semi-supervised setup

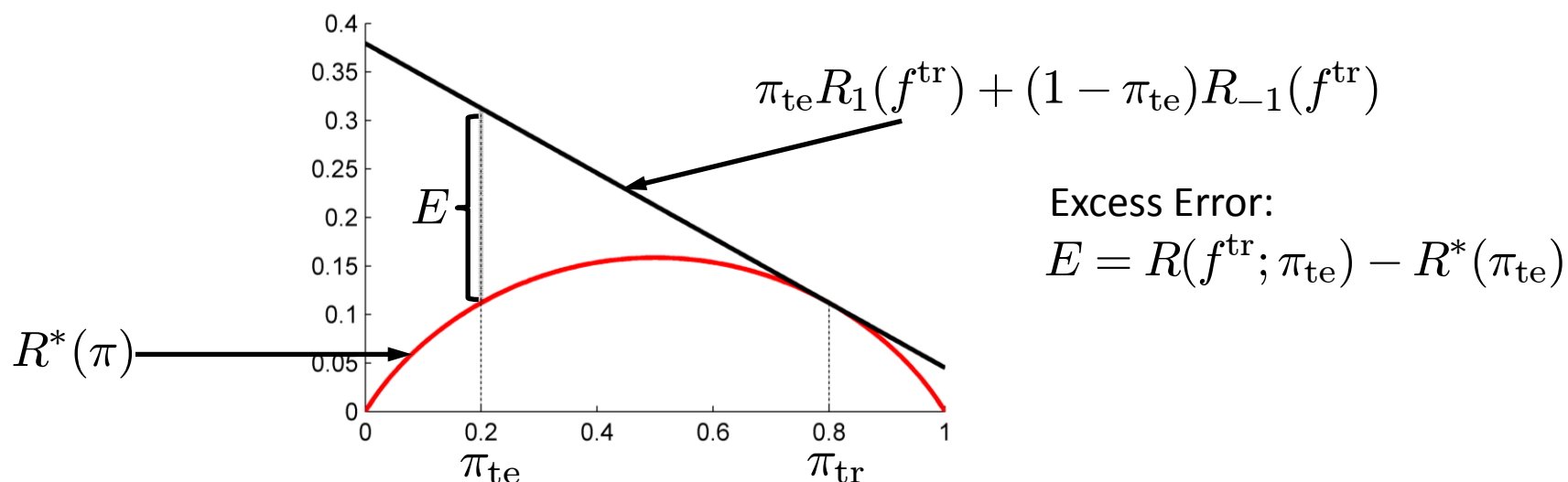
Outline

1. Motivating Example
2. Classification and Risk
3. **Class-prior Change**
 - Class-prior Change
 - Causes of Class-prior Change
 - Dataset shift
 - Selection Bias
 - **Class-prior Change and Risk**
4. Class-prior Change Mitigation
5. Class-prior Change Correction
6. Homework

Effect of Class-prior Change on Risk

- Training: $f_{\text{tr}} = \arg \min_f R(f; \pi_{\text{tr}})$ $R_1(f_{\text{tr}})$ $R_{-1}(f_{\text{tr}})$
- At point π_{tr} this is the same as $R^*(\pi_{\text{tr}})$
- When f_{tr} is applied to a dataset with class prior π_{te} the error is

$$\pi_{\text{te}} R_1(f^{\text{tr}}) + (1 - \pi_{\text{te}}) R_{-1}(f^{\text{tr}}) \geq \min_f R(\pi_{\text{te}})$$



Conclusion: Section 3

- *Class-prior change* may occur between the training and test data:

$$p_{\text{te}}(\mathbf{x}, y) = p(\mathbf{x}|y)p_{\text{te}}(y) \quad p_{\text{tr}}(\mathbf{x}, y) = p(\mathbf{x}|y)p_{\text{tr}}(y)$$
$$p_{\text{tr}}(y) \neq p_{\text{te}}(y)$$

- This may be due to *dataset shift* or *sample selection bias*
- When a classifier is selected according to π_{tr} and applied on a dataset with π_{te} , the Risk is linear and tangent to the optimal risk curve at π_{tr}

Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. **Class-prior Change Mitigation**
 - **The minimax criterion**
5. Class-prior Change Correction
6. Homework

Class-prior change mitigation

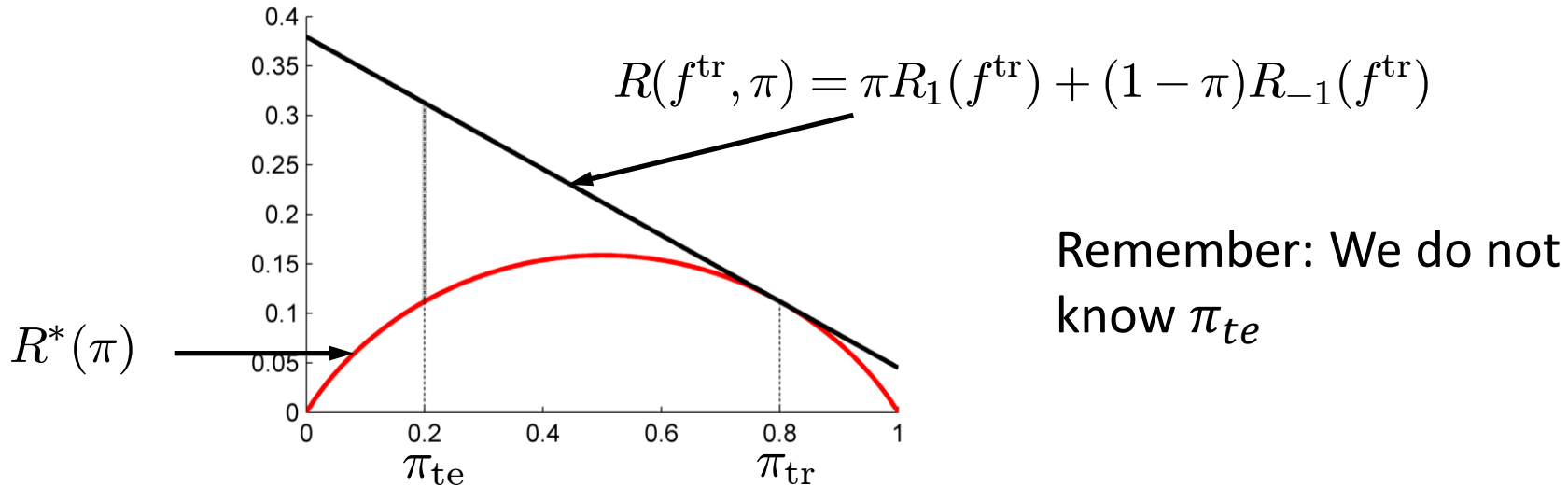
- Class-prior change can have an adverse effect on the classification accuracy
- In practice, the test class prior π_{te} is unknown
 - We can therefore not correct for the effect of class-prior change
- Can we mitigate the effect of class prior change?



Mitigate (definition): to make less severe

Minimax Criterion (1)

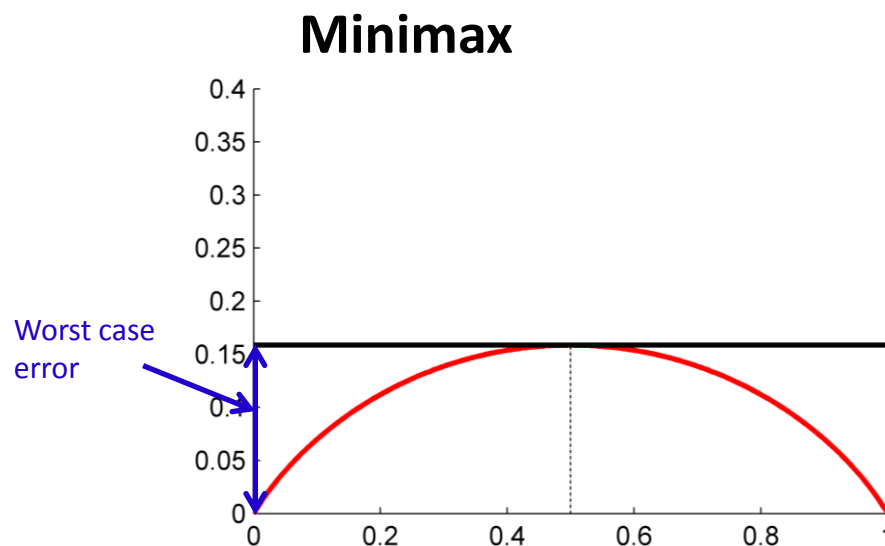
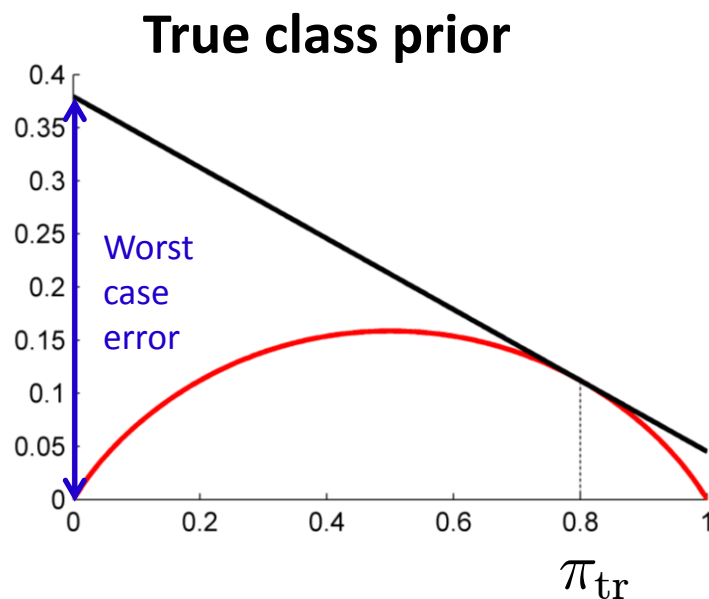
- Recall this figure:



- The black line is the misclassification rate according to the new class prior when trained with the old class prior
- This line is always tangent to the optimal risk $R^*(\pi)$

Minimax Criterion (2)

- Why not select this line so that it does not change w.r.t. the new class prior?
- In other words, the tangent should be 0
- Since $R^*(\pi)$ is concave, this would occur at the maximum



Further Reading

- The minimax criterion is described in
 - “*Pattern Classification*”, 2nd Edition (Richard O. Duda, Peter E. Hart, David G. Stork), p.g. 26.
 - (Ookayama Main Lib. B1F - Books 548.13/D)
 - “*Detection, estimation, and linear modulation theory*” (Van Trees, Harry L.) 1968
 - Ookayama Main Lib. B1F ; Compact Shelving - Y000998
 - The minimax criterion is discussed, and an extension introduced in “*Minimax Regret Classifier for Imprecise Class Distributions*” (Alaiz-Rodríguez, Rocío, Alicia Guerrero-Curieses, and Jesús Cid-Sueiro)
 - <http://jmlr.org/papers/volume8/alaiz-rodriguez07a/alaiz-rodriguez07a.pdf>

Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. **Class-prior Change Correction**
 - **Classifier reweighting**
 - Class-prior Estimation in the Semi-Supervised Setup
 - Class-prior Estimation via Distribution Matching
 - Distribution matching via L_2 -distance minimization
 - Direct L_2 -distance Estimation
 - Example
6. Homework

Correction for Class-prior change

- Recall that a cost-sensitive classifier minimizes

$$R(f) = c_+ \pi_{\text{tr}} R_1(f) + c_- [1 - \pi_{\text{tr}}] R_{-1}(f)$$

False negative rate

False positive rate

- Misclassification rate according to π_{te} can be obtained by weights:

$$c_+ = \frac{\pi_{te}}{\pi_{tr}} \quad c_- = \frac{1 - \pi_{te}}{1 - \pi_{tr}}$$

- Libraries such as libSVM allows specification of cost
- **Problem:** Test class prior π_{te} is often unknown

Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. **Class-prior Change Correction**
 - Classifier reweighting
 - **Class-prior Estimation in the Semi-Supervised Setup**
 - Class-prior Estimation via Distribution Matching
 - Distribution matching via L_2 -distance minimization
 - Direct L_2 -distance Estimation
 - Example
6. Homework

Semi-supervised setup

- In many situations, unlabeled data in addition to labeled data is available

$$\mathcal{X}_{\text{tr}} := \{\mathbf{x}, y\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p_{\text{tr}}(\mathbf{x}, y)$$

$$\mathcal{X}_{\text{te}} := \{\mathbf{x}_i\}_{i=1}^{n'} \stackrel{\text{i.i.d.}}{\sim} p_{\text{te}}(\mathbf{x})$$

- In the class-prior change assumption, the two distributions shares a class-conditional density:

$$p_{\text{tr}}(\mathbf{x}, y) = \underbrace{p(\mathbf{x}|y)}_{\text{Shared}} p_{\text{tr}}(y) \quad p_{\text{te}}(\mathbf{x}, y) = \underbrace{p(\mathbf{x}|y)}_{\text{Shared}} p_{\text{te}}(y)$$

- We wish to estimate the class prior of the unlabeled dataset $p_{\text{te}}(y)$
- This is difficult, because no labeled samples are available

Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. **Class-prior Change Correction**
 - Classifier reweighting
 - Class-prior Estimation in the Semi-Supervised Setup
 - **Class-prior Estimation via Distribution Matching**
 - Distribution matching via L_2 -distance minimization
 - Direct L_2 -distance Estimation
 - Example
6. Homework

Distribution matching framework

- Lets model the test input distribution $p_{\text{te}}(\mathbf{x})$ in terms of:

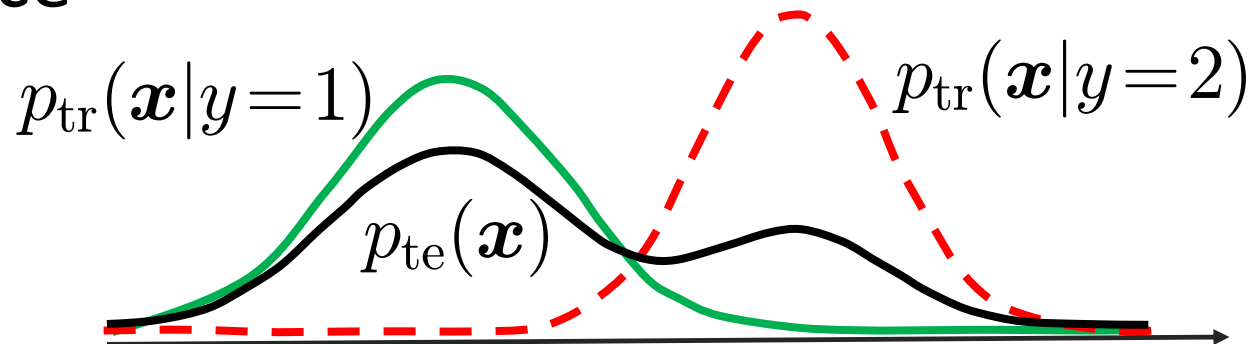
- The training class-conditional distribution $p_{\text{tr}}(\mathbf{x}|y)$
- The test class priors $\pi_y = p_{\text{te}}(y)$

$$q_{\text{te}}(\mathbf{x}) = \sum_{y=1}^c \pi_y p_{\text{tr}}(\mathbf{x}|y)$$

Since $p_{\text{tr}}(\mathbf{x}|y) = p_{\text{te}}(\mathbf{x}|y)$:

$$p_{\text{te}}(\mathbf{x}) = \sum_{y=1}^c p_{\text{tr}}(\mathbf{x}|y) p_{\text{te}}(y)$$

- **Problem:** Match $q_{\text{te}}(\mathbf{x})$ to $p_{\text{te}}(\mathbf{x})$ under some divergence



Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. **Class-prior Change Correction**
 - Classifier reweighting
 - Class-prior Estimation in the Semi-Supervised Setup
 - Class-prior Estimation via Distribution Matching
 - **Distribution matching via L_2 -distance minimization**
 - Direct L_2 -distance Estimation
 - Example
6. Homework

Distribution Matching via L_2 -distance Minimization

- The similarity between two distributions can be measured by the L_2 -distance:

$$L_2(p_{te}, q_{te}) = \frac{1}{2} \int [p_{te}(\mathbf{x}) - q_{te}(\mathbf{x})]^2 d\mathbf{x}$$

- The class prior can therefore be selected as
$$(\pi_1, \dots, \pi_c) = \arg \min_{\pi} L_2(p_{te}(\mathbf{x}), q_{te}(\mathbf{x}))$$
- The L_2 distance can be estimated by first estimating the densities $p_{te}(x)$ and $q_{te}(x)$
 - Not good since density estimation is a difficult problem

Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. **Class-prior Change Correction**
 - Classifier reweighting
 - Class-prior Estimation in the Semi-Supervised Setup
 - Class-prior Estimation via Distribution Matching
 - Distribution matching via L_2 -distance minimization
 - **Direct L_2 -distance Estimation**
 - Example
6. Homework

Direct L_2 estimation

- Expectations can be estimated via sample averages:

$$\int f(\mathbf{x})p(\mathbf{x})d\mathbf{x} = \mathbb{E}_p[f(\mathbf{x})] \approx \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i) \quad \{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

- Can we estimate the L_2 -distance in terms of sample averages?
- It is possible by obtaining a lower-bound that is linear in the densities

Direct L_2 estimation

- We use the following inequality:

$$\frac{1}{2}(t - y)^2 \geq 0, \quad \longrightarrow \quad \frac{1}{2}t^2 \geq ty - \frac{1}{2}y^2$$

- Applying this pointwise gives $w(\mathbf{x})$ role of y

$$\frac{1}{2} [p_{\text{te}}(\mathbf{x}) - q_{\text{te}}(\mathbf{x})]^2 \geq w(\mathbf{x}) [p_{\text{te}}(\mathbf{x}) - q_{\text{te}}(\mathbf{x})] - \frac{1}{2}w(\mathbf{x})^2$$

- Integrating and selecting the tightest lower-bound gives

$$\begin{aligned} \frac{1}{2} \int [p_{\text{te}}(\mathbf{x}) - q_{\text{te}}(\mathbf{x})]^2 d\mathbf{x} \\ \geq \sup_w \int w(\mathbf{x}) [p_{\text{te}}(\mathbf{x}) - q_{\text{te}}(\mathbf{x})] d\mathbf{x} - \frac{1}{2} \int w(\mathbf{x})^2 d\mathbf{x} \end{aligned}$$

Direct L_2 estimation

- This lower-bound can then be estimated via sample averages
- Lets model $w(\mathbf{x})$ with a linear model

$$w(\mathbf{x}) = \sum_{\ell=1}^b \alpha_{\ell} \varphi_{\ell}(\mathbf{x}) \quad \varphi_{\ell}(\mathbf{x}) = \exp \left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{c}_{\ell}\|^2 \right)$$

- The L_2 -distance lower bound can be written in terms of expectations

$$L_2(p_{\text{te}}, q_{\text{te}}) \geq \sup_w \mathbb{E}_{p_{\text{te}}} [w(\mathbf{x})] - \sum_{c=1}^{n_c} \mathbb{E}_{\pi_y p_{\text{te}}} [w(\mathbf{x})] - \frac{1}{2} \int w(\mathbf{x})^2 d\mathbf{x}$$

Direct L_2 estimation

$$L_2(p_{te}, q_{te}) \geq \sup_w \mathbb{E}_{p_{te}} [w(\mathbf{x})] - \sum_{c=1}^{n_c} \mathbb{E}_{\pi_y p_{te}} [w(\mathbf{x})] - \frac{1}{2} \int w(\mathbf{x})^2 d\mathbf{x}$$

- The expectations can be estimated via sample averages:

$$\hat{\mathbf{h}}_{te} = \frac{1}{n} \sum_{i=1}^n \varphi(\mathbf{x}_i)$$

$$\hat{\mathbf{h}}_c = \frac{1}{n_c} \sum_{i=1, y_i=c}^n \varphi(\mathbf{x}_i)$$

$$\mathbf{H} = \int \varphi(\mathbf{x}) \varphi(\mathbf{x})^\top d\mathbf{x}$$

$$\varphi(\mathbf{x}) = [\varphi_1(\mathbf{x}) \quad \varphi_2(\mathbf{x}) \quad \dots \quad \varphi_\ell(\mathbf{x})]$$

- Which gives an objective function of:

$$\hat{L}_2(\{\pi_y\}_{y=1}^c) \approx \max_{\alpha} \alpha^\top \hat{\mathbf{h}}_{te} - \sum_{y=1}^c \theta_y \alpha^\top \hat{\mathbf{h}}_y - \frac{1}{2} \alpha^\top \mathbf{H} \alpha$$

Direct L_2 estimation

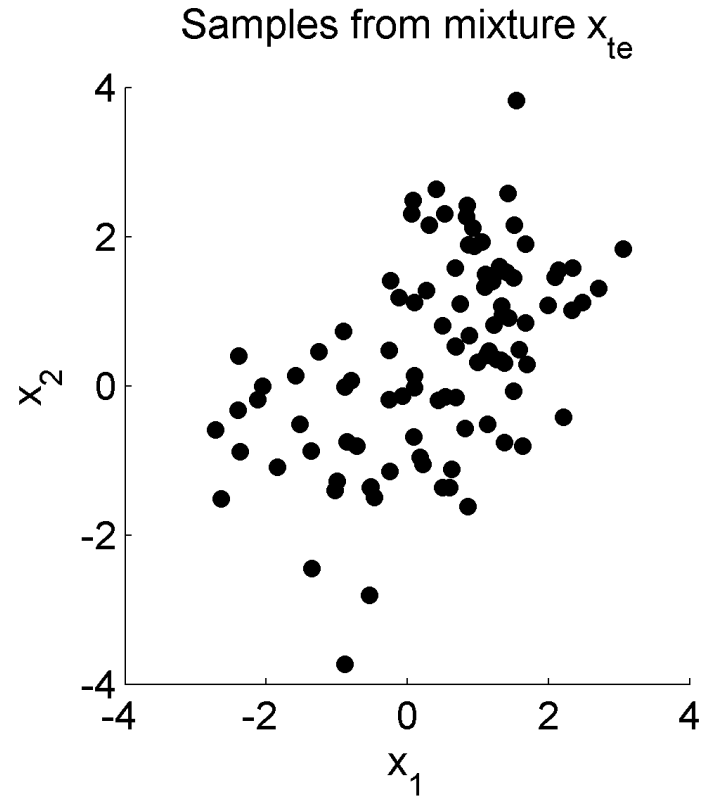
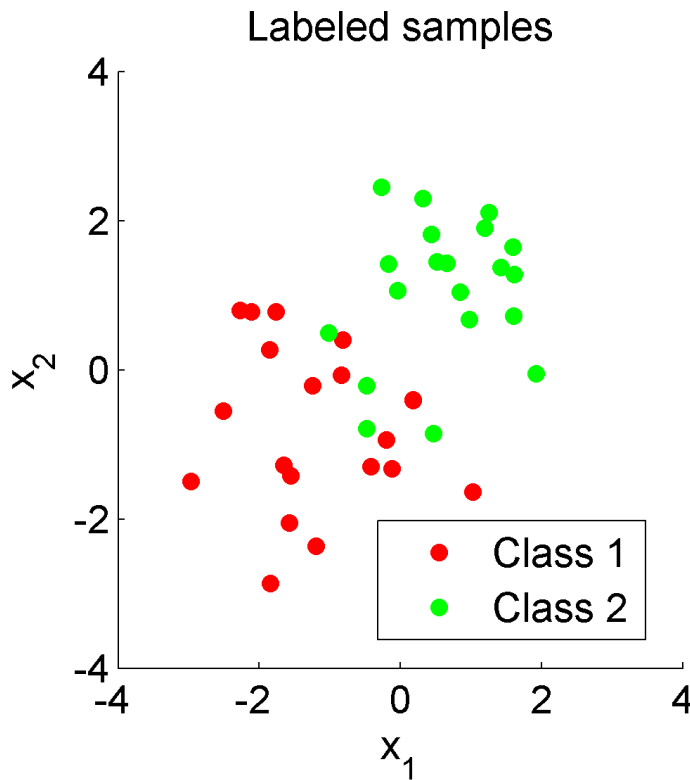
- Minimizing the L_2 distance estimate w.r.t. $\{\pi_y\}_{y=1}^c$ gives an estimate of the class prior
- This can then be used to reweight a classifier

Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. **Class-prior Change Correction**
 - Classifier reweighting
 - Class-prior Estimation in the Semi-Supervised Setup
 - Class-prior Estimation via Distribution Matching
 - Distribution matching via L_2 -distance minimization
 - Direct L_2 -distance Estimation
 - **Example**
6. Homework

Example

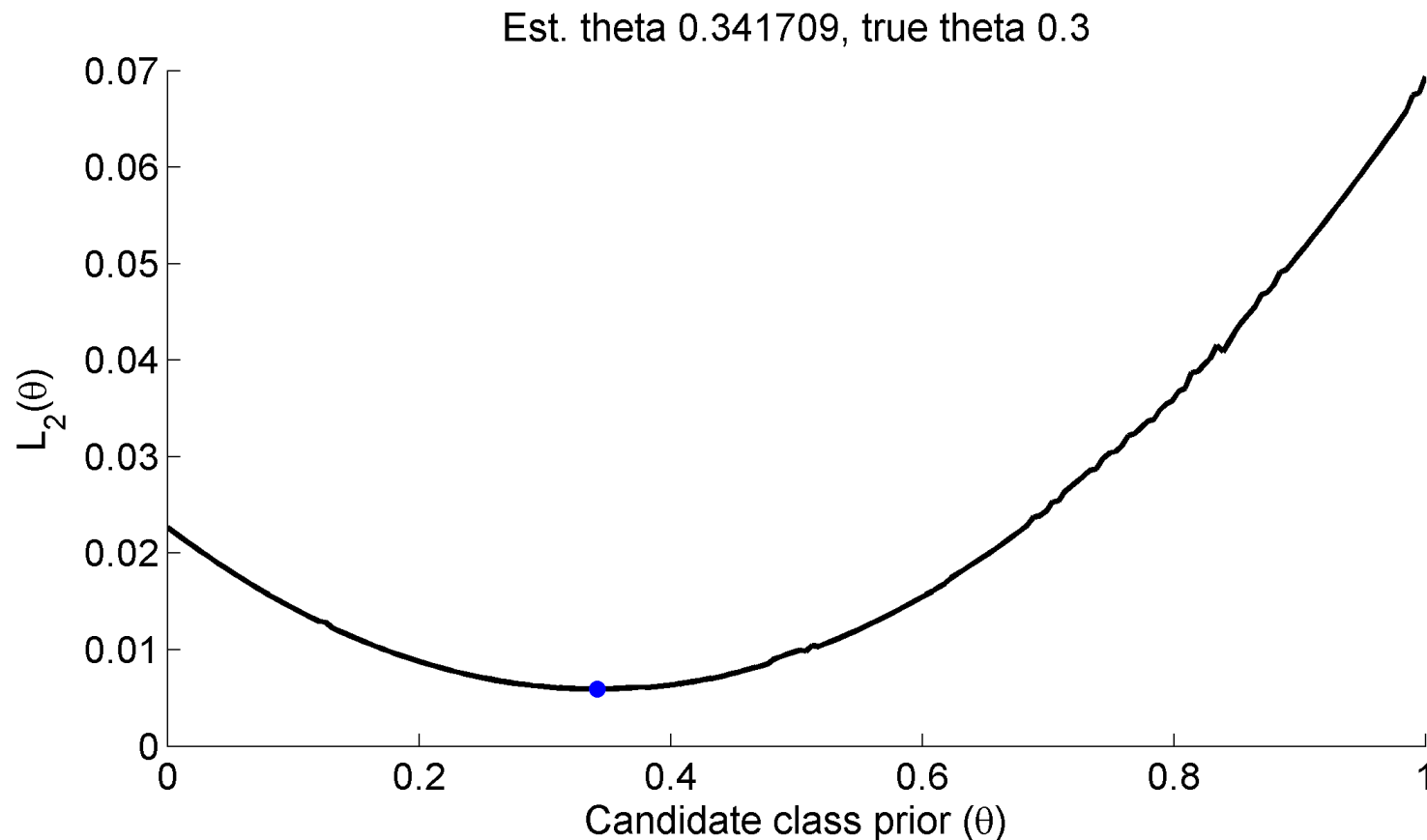
- Samples from two Gaussians with different means:



- The true class prior is $p_{te}(y = 1) = 0.3$

Example (2)

- The L_2 -distance estimated from samples is given below:



- Minimum of L_2 distance is near the true class prior
- Difference is due to estimation from a small set I

Conclusion: Section 5

- By *reweighting* the risk, a classifier can be trained
- The *reweighting factor* depends on the *unknown class prior*
- In a *semi-supervised setup*, the *unknown class prior can be estimated*
- Estimation is possible by matching a model of the test input density to the true test input density

Further Reading

- “Adjusting the outputs of a classifier to new a priori probabilities: a simple procedure” (Saerens, M., Latinne, P., and Decaestecker, C.)
 - Neurocomputation 14 (2002)
 - Introduced estimation of the class prior for re-adjustment of the classifier
- “Semi-supervised learning of class balance under class-prior change by distribution matching.” (du Plessis, M. C. & Sugiyama, M.)
 - Estimation of the class prior via Pearson divergence matching
- “Density-difference estimation” (Sugiyama, M., Suzuki, T., Kanamori, T., du Plessis, M. C., Liu, S., & Takeuchi, I.)
 - Estimation of the class prior via L_2 -distance estimation (discussed here)

Outline

1. Motivating Example
2. Classification and Risk
3. Class-prior Change
4. Class-prior Change Mitigation
5. Class-prior Change Correction
6. **Take-home message and homework**

Take-home Message

- The class-prior may change between the training and test phase
- In *supervised learning*, the minimax approach can be used:
 - Minimizes the worst case result
- If the test class-prior is known, the classifier can be selected by reweighting
- In *semi-supervised learning*, the test class prior can be estimated

Homework

- Please hand in your reports now!

- Homework:

Write you opinion about the special lecture today.

Directly submit the printed report to the lecturer next week.