

SEMISUPERVISED
LEARNING

33

CHAPTER CONTENTS

Manifold Regularization	375
Manifold Structure Brought by Input Samples	375
Computing the Solution	377
Covariate Shift Adaptation	378
Importance Weighted Learning	378
Relative Importance Weighted Learning	382
Importance Weighted Cross Validation	382
Importance Estimation	383
Class-balance Change Adaptation	385
Class-balance Weighted Learning	385
Class-balance Estimation	386

Supervised learning is performed based on input-output paired training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$. However, gathering many input-output paired samples is often expensive in practical applications. On the other hand, *input-only* samples $\{\mathbf{x}_i\}_{i=n+1}^{n+n'}$ can be easily collected abundantly. For example, in web page classification, class label y_i (such as “sport”, “computer”, and “politics”) should be *manually* given after carefully investigating web page \mathbf{x}_i , which requires a huge amount of human labor. On the other hand, input-only web pages $\{\mathbf{x}_i\}_{i=n+1}^{n+n'}$ can be automatically collected by crawlers.

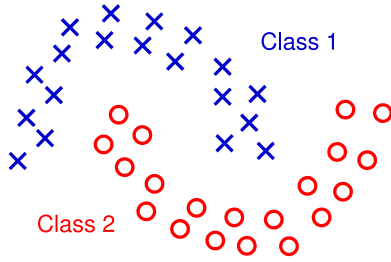
In this chapter, methods of *semisupervised learning* [28, 101] are introduced, which utilize input-only samples $\{\mathbf{x}_i\}_{i=n+1}^{n+n'}$ in addition to input-output paired samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$.

33.1 MANIFOLD REGULARIZATION

In this section, a method of semisupervised learning based on *manifold regularization* is introduced.

33.1.1 MANIFOLD STRUCTURE BROUGHT BY INPUT
SAMPLES

Supervised learning from input-output paired samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ can be regarded as estimating the conditional density $p(y|\mathbf{x})$. On the other hand, unsupervised

**FIGURE 33.1**

Semisupervised classification. Samples in the same cluster are assumed to belong to the same class.

learning from input-only samples $\{\mathbf{x}_i\}_{i=n+1}^{n+n'}$ can be regarded as estimating the marginal density $p(\mathbf{x})$. Thus, without any assumption, input-only samples $\{\mathbf{x}_i\}_{i=n+1}^{n+n'}$ do not help improve the estimation of the conditional density $p(y|\mathbf{x})$. For this reason, semisupervised learning imposes a certain assumption between $p(\mathbf{x})$ and $p(y|\mathbf{x})$ and utilizes the estimation of $p(\mathbf{x})$ to improve the accuracy of estimating $p(y|\mathbf{x})$.

Below, a semisupervised learning method based on a *manifold* assumption is introduced. Mathematically, a manifold is a topological space that can be locally approximated by Euclidean space. On the other hand, a manifold is just regarded as a local region in the context of semisupervised learning. More specifically, the manifold assumption means that input samples appear only on manifolds and output values change *smoothly* on the manifolds. In the case of classification, this means that samples in the same cluster belong to the same class (Fig. 33.1).

The Gaussian kernel model introduced in Section 21.2 actually utilizes this manifold assumption (see Fig. 21.5):

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^n \theta_j K(\mathbf{x}, \mathbf{x}_j), \quad K(\mathbf{x}, \mathbf{c}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}\|^2}{2h^2}\right).$$

That is, by locating smooth Gaussian functions on input samples $\{\mathbf{x}_i\}_{i=1}^n$, a smooth function over the input manifold can be learned. In semisupervised learning, the above model may be augmented to locate Gaussian kernels also on input-only samples $\{\mathbf{x}_i\}_{i=n+1}^{n+n'}$:

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^{n+n'} \theta_j K(\mathbf{x}, \mathbf{x}_j). \quad (33.1)$$

The parameters in model (33.1) are learned so that output at input samples, $\{f_{\theta}(\mathbf{x}_i)\}_{i=1}^{n+n'}$, is similar to each other. For example, in the case of ℓ_2 -regularized

LS, the optimization problem is given as

$$\min_{\boldsymbol{\theta}} \left[\frac{1}{2} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 + \frac{\nu}{4} \sum_{i,i'=1}^{n+n'} W_{i,i'} (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - f_{\boldsymbol{\theta}}(\mathbf{x}_{i'}))^2 \right], \quad (33.2)$$

where the first and second terms correspond to ℓ_2 -regularized LS and the third term is called the *Laplacian regularizer*, following the terminology in *spectral graph theory* [29]. $\nu \geq 0$ is the regularization parameter for semisupervised learning that controls the smoothness on the manifolds. $W_{i,i'} \geq 0$ denotes the similarity between \mathbf{x}_i and $\mathbf{x}_{i'}$, which takes a large value if \mathbf{x}_i and $\mathbf{x}_{i'}$ are similar, and a small value if \mathbf{x}_i and $\mathbf{x}_{i'}$ are dissimilar. Popular choices of the similarity measure are described in Fig. 35.8. Below the similarity matrix \mathbf{W} is assumed to be symmetric (i.e. $W_{i,i'} = W_{i',i}$).

33.1.2 COMPUTING THE SOLUTION

Here, how to compute the solution of Laplacian-regularized LS is explained. Let \mathbf{D} be the diagonal matrix whose diagonal elements are given by the row-sums of matrix \mathbf{W} :

$$\mathbf{D} = \text{diag} \left(\sum_{i=1}^{n+n'} W_{1,i}, \dots, \sum_{i=1}^{n+n'} W_{n+n',i} \right),$$

and let $\mathbf{L} = \mathbf{D} - \mathbf{W}$. Then the third term in Eq. (33.2) can be rewritten as

$$\begin{aligned} & \sum_{i,i'=1}^{n+n'} W_{i,i'} (f_{\boldsymbol{\theta}}(\mathbf{x}_i) - f_{\boldsymbol{\theta}}(\mathbf{x}_{i'}))^2 \\ &= \sum_{i=1}^{n+n'} D_{i,i} f_{\boldsymbol{\theta}}(\mathbf{x}_i)^2 - 2 \sum_{i,i'=1}^{n+n'} W_{i,i'} f_{\boldsymbol{\theta}}(\mathbf{x}_i) f_{\boldsymbol{\theta}}(\mathbf{x}_{i'}) + \sum_{i'=1}^{n+n'} D_{i',i'} f_{\boldsymbol{\theta}}(\mathbf{x}_{i'})^2 \\ &= 2 \sum_{i,i'=1}^{n+n'} L_{i,i'} f_{\boldsymbol{\theta}}(\mathbf{x}_i) f_{\boldsymbol{\theta}}(\mathbf{x}_{i'}). \end{aligned}$$

Thus, Eq. (33.2) for kernel model (33.1) can be reduced to the generalized ℓ_2 -regularized LS as follows:

$$\begin{aligned} \min_{\boldsymbol{\theta}} \left[\frac{1}{2} \sum_{i=1}^n \left(\sum_{j=1}^{n+n'} \theta_j K(\mathbf{x}_i, \mathbf{x}_j) - y_i \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n+n'} \theta_j^2 \right. \\ \left. + \frac{\nu}{2} \sum_{j,j'=1}^{n+n'} \theta_j \theta_{j'} \sum_{i,i'=1}^{n+n'} L_{i,i'} K(\mathbf{x}_i, \mathbf{x}_j) K(\mathbf{x}_{i'}, \mathbf{x}_{j'}) \right]. \end{aligned}$$

This can be compactly rewritten as

$$\min_{\theta} \left[\frac{1}{2} \|\mathbf{K}\theta - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\theta\|^2 + \frac{\nu}{2} \theta^\top \mathbf{K} \mathbf{L} \mathbf{K} \theta \right],$$

where

$$\mathbf{K} = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_{n+n'}) \\ \vdots & \ddots & \vdots \\ K(\mathbf{x}_{n+n'}, \mathbf{x}_1) & \cdots & K(\mathbf{x}_{n+n'}, \mathbf{x}_{n+n'}) \end{pmatrix},$$

$$\theta = (\theta_1, \dots, \theta_n, \theta_{n+1}, \dots, \theta_{n+n'})^\top,$$

$$\mathbf{y} = (y_1, \dots, y_n, \underbrace{0, \dots, 0}_{n'})^\top.$$

Note that the vector \mathbf{y} defined above contains n' zeros since $\{y_i\}_{i=n+1}^{n+n'}$ are not available. The solution $\hat{\theta}$ can be obtained analytically as

$$\hat{\theta} = (\mathbf{K}^2 + \lambda \mathbf{I} + \nu \mathbf{K} \mathbf{L} \mathbf{K})^{-1} \mathbf{K} \mathbf{y}.$$

A MATLAB code for Laplacian-regularized LS is provided in [Fig. 33.2](#), and its behavior is illustrated in [Fig. 33.3](#). This is an extremely difficult classification problem where only two labeled training samples are available. Laplacian-regularized LS gives a decision boundary that separates two point clouds, which is appropriate when the cluster assumption holds. On the other hand, ordinary LS gives a decision boundary that goes through the middle of the two labeled training samples, which is not reasonable under the cluster assumption.

In the above explanation, Laplacian regularization was applied to LS. However, the idea of Laplacian regularization is generic and it can be combined with various regression and classification techniques introduced in [Part 4](#).

33.2 COVARIATE SHIFT ADAPTATION

The Laplace regularization method introduced above utilized the manifold structure of input samples. In this section, a semisupervised learning method called *covariate shift* adaptation is introduced [[101](#)], which explicitly takes into account the probability distributions of $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$. A covariate is another name for an input variable, and covariate shift refers to the situation where $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ follow different probability distributions, but the conditional density $p(y|\mathbf{x})$ is unchanged.

33.2.1 IMPORTANCE WEIGHTED LEARNING

An example of covariate shift in regression is illustrated in [Fig. 33.4](#), where the learning target function $f(x)$ is unchanged. However, $\{x_i\}_{i=1}^n$ are distributed

```

n=200; a=linspace(0,pi,n/2);
u=-10*[cos(a)+0.5 cos(a)-0.5]'+randn(n,1);
v=10*[sin(a) -sin(a)]'+randn(n,1);
x=[u v]; y=zeros(n,1); y(1)=1; y(n)=-1;
x2=sum(x.^2,2); hh=2*1^2;
k=exp(-( repmat(x2,1,n)+repmat(x2',n,1)-2*x*x')/hh); w=k;
t=(k^2+1*eye(n)+10*k*(diag(sum(w))-w)*k)\(k*y);

m=100; X=linspace(-20,20,m)'; X2=X.^2;
U=exp(-( repmat(u.^2,1,m)+repmat(X2',n,1)-2*u*X')/hh);
V=exp(-( repmat(v.^2,1,m)+repmat(X2',n,1)-2*v*X')/hh);
figure(1); clf; hold on; axis([-20 20 -20 20]);
colormap([1 0.7 1; 0.7 1 1]);
contourf(X,X,sign(V'*(U.*repmat(t,1,m)))));
plot(x(y==1,1),x(y==1,2),'bo');
plot(x(y==-1,1),x(y==-1,2),'rx');
plot(x(y==0,1),x(y==0,2),'k. ');

```

FIGURE 33.2

MATLAB code for Laplacian-regularized LS.

around $x = 1$, while $\{x'_{i'}\}_{i'=1}^{n'}$ are distributed around $x = 2$. Thus, this is a (weak) *extrapolation* problem.

Fig. 33.6(a) illustrates the function obtained by fitting the straight line model,

$$f_{\theta}(x) = \theta_1 + \theta_2 x, \quad (33.3)$$

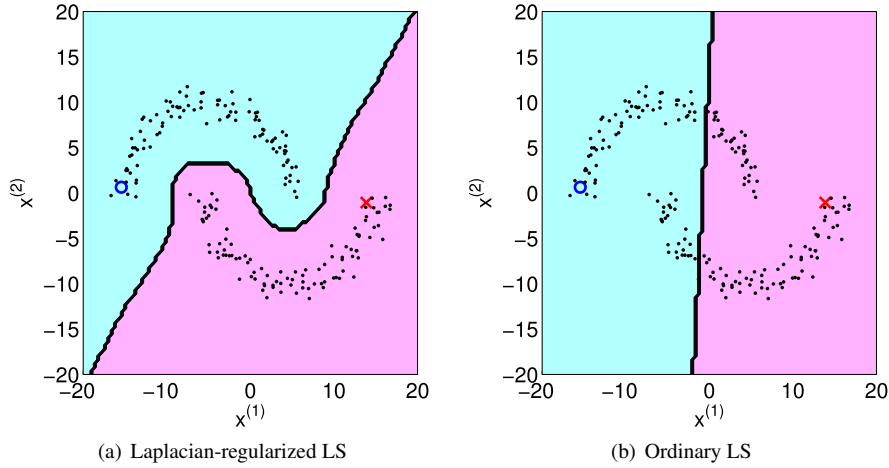
to samples $\{(x_i, y_i)\}_{i=1}^n$ by ordinary LS:

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n (f_{\theta}(x_i) - y_i)^2.$$

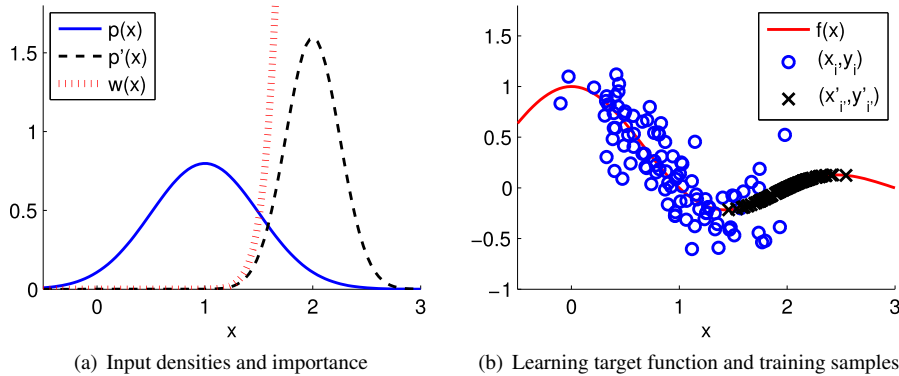
This shows that, although the samples $\{(x_i, y_i)\}_{i=1}^n$ are appropriately fitted by LS, prediction of output values at $\{x'_{i'}\}_{i'=1}^{n'}$ is poor.

In such a covariate shift situation, it intuitively seems that only using samples $\{(x_i, y_i)\}_{i=1}^n$ near $\{x'_{i'}\}_{i'=1}^{n'}$ can provide good prediction of output values at $\{x'_{i'}\}_{i'=1}^{n'}$. This intuitive idea can be more formally realized by *importance weighting*. More specifically, *importance weighted LS* is given by

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n w(x_i) (f_{\theta}(x_i) - y_i)^2,$$

**FIGURE 33.3**

Examples of Laplacian-regularized LS compared with ordinary LS. Dots denote unlabeled training samples.

**FIGURE 33.4**

Covariate shift in regression. Input distributions change, but the input-output relation is unchanged.

where $w(x)$ is called the importance function which is defined as the ratio of probability density $p'(x)$ for $\{x'_i\}_{i=1}^{n'}$ and probability density $p(x)$ for $\{x_i\}_{i=1}^n$:

$$w(x) = \frac{p'(x)}{p(x)}.$$

```

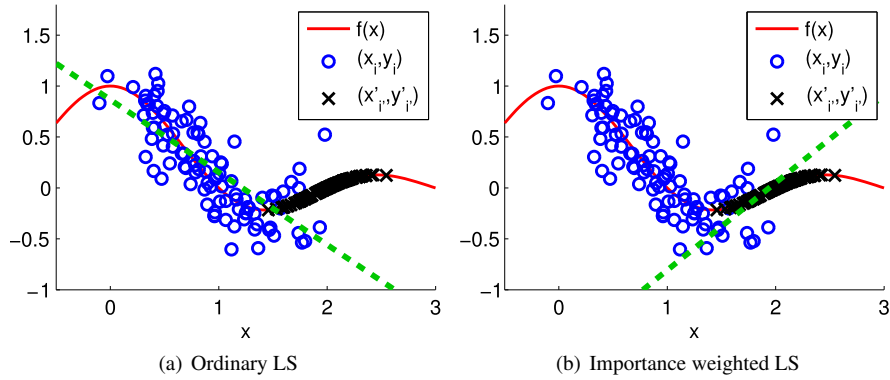
n=100; u=randn(n,1)/4+2; x=randn(n,1)/2+1;
w=2*exp(-8*(x-2).^2+2*(x-1).^2); %w=ones(n,1);

y=sin(pi*x)/(pi*x)+0.1*randn(n,1);
x(:,2)=1; t=(x*(repmat(w,1,2).*x))\'(x*(w.*y));
X=linspace(-1,3,100); Y=sin(pi*X)/(pi*X);
u(:,2)=1; v=u*t;
figure(1); clf; hold on;
plot(x(:,1),y,'bo'); plot(X,Y,'r-'); plot(u(:,1),v,'kx');

```

FIGURE 33.5

MATLAB code for importance weighted LS.

**FIGURE 33.6**

Example of LS learning under covariate shift. The dashed lines denote learned functions.

This importance weighted learning is based on the idea of *importance sampling* (see Section 19.2), which approximates the expectation with respect to $p'(x)$ by importance weighted average with respect to $p(x)$:

$$\int g(x)p'(x)dx = \int g(x)\frac{p'(x)}{p(x)}p(x)dx \approx \frac{1}{n} \sum_{i=1}^n g(x_i)w(x_i).$$

A MATLAB code of importance weighted LS for straight line model (33.3) is provided in Fig. 33.5, and its behavior is illustrated in Fig. 33.6(b). This shows that importance weighting contributes to improving the accuracy of predicting output values at $\{x'_i\}_{i=1}^n$.

In the above example, importance weighting was applied to LS. However, the idea of importance weighting is generic and it can be combined with various regression and classification techniques introduced in [Part 4](#).

33.2.2 RELATIVE IMPORTANCE WEIGHTED LEARNING

As illustrated in [Fig. 33.6\(b\)](#), prediction performance under covariate shift can be improved by importance weighted learning.

[Fig. 33.4\(a\)](#) illustrates the importance function $w(x)$, which is monotone increasing as x grows. This means that, among all training samples $\{(x_i, y_i)\}_{i=1}^n$, only a few of them with large x have large importance values and the importance values of others are negligibly small. Thus, importance weighted learning is actually performed only from a few training samples, which can be unstable in practice.

This instability is caused by the fact that the importance function $w(\mathbf{x})$ takes large values for some \mathbf{x} . Thus, the instability problem is expected to be mitigated if a smooth variant of the importance function, such as the *relative importance* [121] defined as

$$w_\beta(\mathbf{x}) = \frac{p'(\mathbf{x})}{\beta p'(\mathbf{x}) + (1 - \beta)p(\mathbf{x})},$$

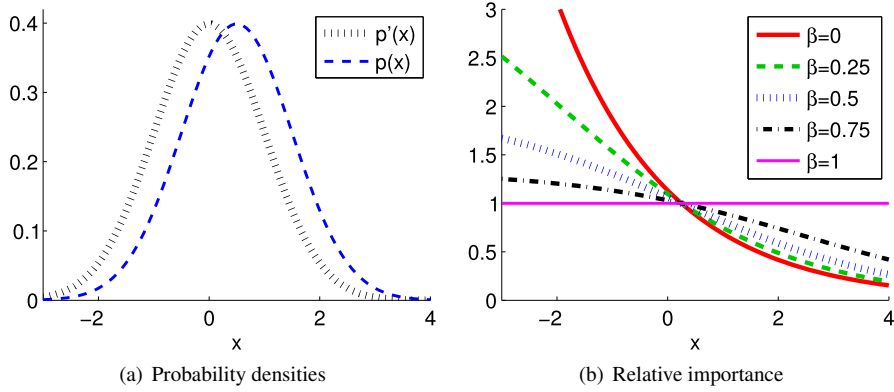
is employed, where $\beta \in [0, 1]$ controls the smoothness. The relative importance function $w_\beta(\mathbf{x})$ is reduced to the ordinary importance function $p'(\mathbf{x})/p(\mathbf{x})$ if $\beta = 0$. It gets smoother if β is increased, and it yields the uniform weight $w_\beta(\mathbf{x}) = 1$ if $\beta = 1$ ([Fig. 33.7](#)). Since the importance is always non-negative, the relative importance is no larger than $1/\beta$:

$$w_\beta(\mathbf{x}) = \frac{1}{\beta + (1 - \beta) \frac{p(\mathbf{x})}{p'(\mathbf{x})}} \leq \frac{1}{\beta}.$$

33.2.3 IMPORTANCE WEIGHTED CROSS VALIDATION

The performance of relative importance weighted learning depends on the choice of smoothness β . Also, choice of model $f_\theta(\mathbf{x})$ and regularization parameters significantly affects the final performance.

In [Section 23.3](#), a model selection method based on *cross validation* was introduced, which estimates the prediction error for test samples. Actually, the validity of cross validation is ensured only when $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ follow the same probability distributions. Thus, performing cross validation under covariate shift can produce undesired solutions. Under covariate shift, a variant called *importance weighted cross validation* [102] is useful. The algorithm of importance weighted cross validation is summarized in [Fig. 33.8](#), which is essentially the same as ordinary cross validation, but the validation error is computed with importance weights.

**FIGURE 33.7**

Relative importance when $p'(x)$ is the Gaussian density with expectation 0 and variance 1 and $p(x)$ is the Gaussian density with expectation 0.5 and variance 1.

33.2.4 IMPORTANCE ESTIMATION

In (relative) importance weighted learning and cross validation, importance weights are needed. However, the importance function is usually unknown and only samples $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ are available. A naive way to estimate the importance function $w(\mathbf{x}) = p'(\mathbf{x})/p(\mathbf{x})$ is to separately estimate density $p'(\mathbf{x})$ from $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ and $p(\mathbf{x})$ from $\{\mathbf{x}_i\}_{i=1}^n$ and then compute the ratio of estimated densities. However, such a two-step approach is unreliable because division by an estimated density can magnify the estimation error significantly. Here, a direct importance estimator without going through density estimation is introduced.

Let us model the relative importance function $w_\beta(\mathbf{x})$ by a linear-in-parameter model:

$$w_\alpha(\mathbf{x}) = \sum_{j=1}^b \alpha_j \psi_j(\mathbf{x}) = \alpha^\top \boldsymbol{\psi}(\mathbf{x}),$$

where $\alpha = (\alpha_1, \dots, \alpha_b)^\top$ is the parameter vector and $\boldsymbol{\psi}(\mathbf{x}) = (\psi_1(\mathbf{x}), \dots, \psi_b(\mathbf{x}))^\top$ is the vector of basis functions. The parameter α is learned so that the following $J(\alpha)$ is minimized:

$$\begin{aligned} J(\alpha) &= \frac{1}{2} \int (w_\alpha(\mathbf{x}) - w_\beta(\mathbf{x}))^2 (\beta p'(\mathbf{x}) + (1 - \beta)p(\mathbf{x})) d\mathbf{x} \\ &= \frac{1}{2} \int \alpha^\top \boldsymbol{\psi}(\mathbf{x}) \boldsymbol{\psi}(\mathbf{x})^\top \alpha (\beta p'(\mathbf{x}) + (1 - \beta)p(\mathbf{x})) d\mathbf{x} \\ &\quad - \int \alpha^\top \boldsymbol{\psi}(\mathbf{x}) p'(\mathbf{x}) d\mathbf{x} + C, \end{aligned}$$

1. Prepare candidates of models: $\{\mathcal{M}_j\}_j$.
2. Split training samples $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ into t disjoint subsets of (approximately) the same size: $\{\mathcal{D}_\ell\}_{\ell=1}^t$.
3. For each model candidate \mathcal{M}_j
 - (a) For each split $\ell = 1, \dots, t$

- i. Obtain learned function $\hat{f}_j^{(\ell)}(\mathbf{x})$ using model \mathcal{M}_j from all training samples without \mathcal{D}_ℓ .
- ii. Compute the average prediction error $\hat{G}_j^{(\ell)}$ of $\hat{f}_j^{(\ell)}(\mathbf{x})$ for holdout samples \mathcal{D}_ℓ .
 - Regression (squared loss):

$$\hat{G}_j^{(\ell)} = \frac{1}{|\mathcal{D}_\ell|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_\ell} w(\mathbf{x}) (y - \hat{f}_j^{(\ell)}(\mathbf{x}))^2,$$

where $|\mathcal{D}_\ell|$ denotes the number of elements in the set \mathcal{D}_ℓ .

- Classification (0/1-loss):

$$\hat{G}_j^{(\ell)} = \frac{1}{|\mathcal{D}_\ell|} \sum_{(\mathbf{x}, y) \in \mathcal{D}_\ell} \frac{w(\mathbf{x})}{2} \left(1 - \text{sign}(\hat{f}_j^{(\ell)}(\mathbf{x})y) \right).$$

- (b) Compute the average prediction error \hat{G}_j over all t splits:

$$\hat{G}_j = \frac{1}{t} \sum_{\ell=1}^t \hat{G}_j^{(\ell)}.$$

4. Choose the model $\mathcal{M}_{\hat{j}}$ that minimizes the average prediction error:

$$\hat{j} = \underset{j}{\text{argmin}} \hat{G}_j.$$

5. Obtain the final function approximator using chosen model $\mathcal{M}_{\hat{j}}$ from all training samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$.

FIGURE 33.8

Algorithm of importance weighted cross validation.

where $C = \frac{1}{2} \int w_\beta(\mathbf{x}) p'(\mathbf{x}) d\mathbf{x}$ is independent of α and thus can be ignored. Approximating the expectations by sample averages and adding the ℓ_2 -regularizer yield the following optimization problem:

$$\min_{\alpha} \left[\frac{1}{2} \alpha^{\top} \widehat{G}_{\beta} \alpha - \alpha^{\top} \widehat{h} + \frac{\lambda}{2} \|\alpha\|^2 \right],$$

where

$$\begin{aligned} \widehat{G}_{\beta} &= \frac{\beta}{n'} \sum_{i'=1}^{n'} \psi(\mathbf{x}'_{i'}) \psi(\mathbf{x}'_{i'})^{\top} + \frac{1-\beta}{n} \sum_{i=1}^n \psi(\mathbf{x}_i) \psi(\mathbf{x}_i)^{\top}, \\ \widehat{h} &= \frac{1}{n'} \sum_{i'=1}^{n'} \psi(\mathbf{x}'_{i'}). \end{aligned}$$

The minimizer $\widehat{\alpha}$ can be obtained analytically as

$$\widehat{\alpha} = (\widehat{G} + \lambda \mathbf{I})^{-1} \widehat{h}.$$

This method is called *LS relative density ratio estimation* [121]. The regularization parameter λ and parameters included in basis functions ψ can be optimized by cross validation with respect to the squared error J .

A MATLAB code of LS relative density ratio estimation for the Gaussian kernel model,

$$w_{\alpha}(\mathbf{x}) = \sum_{j=1}^n \alpha_j \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_j\|^2}{2h^2}\right),$$

is provided in Fig. 33.9, and its behavior is illustrated in Fig. 33.10. This shows that the true relative importance function is nicely estimated.

33.3 CLASS-BALANCE CHANGE ADAPTATION

In the previous section, semisupervised learning methods for covariate shift adaptation were introduced, which can be naturally applied in regression. On the other hand, in classification, *class-balance change* is a natural situation, where the class-prior probabilities differ in $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$ (see Fig. 33.11), but the class-conditional probability $p(\mathbf{x}|y)$ remains unchanged. In this section, adaptation methods for class-balance change are introduced.

33.3.1 CLASS-BALANCE WEIGHTED LEARNING

The bias caused by class-balance change can be canceled by *class-balance weighted learning*.

More specifically, the class-prior ratio $p'(y)/p(y)$ is used as a weighting factor, where $p(y)$ and $p'(y)$ are the class-prior probabilities for $\{\mathbf{x}_i\}_{i=1}^n$ and $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$, respectively. For example, in the case of LS, the learning criterion is given by

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \frac{p'(y_i)}{p(y_i)} (f_{\theta}(\mathbf{x}_i) - y_i)^2,$$

```

n=300; x=randn(n,1); y=randn(n,1)+0.5;
hhs=2*[1 5 10].^2; ls=10.^[-3 -2 -1]; m=5; b=0.5;
x2=x.^2; xx= repmat(x2,1,n)+ repmat(x2',n,1)-2*x*x';
y2=y.^2; yx= repmat(y2,1,n)+ repmat(x2',n,1)-2*y*x';
u=mod(randperm(n),m)+1; v=mod(randperm(n),m)+1;

for hk=1:length(hhs)
    hh=hhs(hk); k=exp(-xx/hh); r=exp(-yx/hh);
    for i=1:m
        ki=k(u~=i,:); ri=r(v~=i,:); h=mean(ki)';
        kc=k(u==i,:); rj=r(v==i,:);
        G=b*ki'*ki/sum(u~=i)+(1-b)*ri'*ri/sum(v~=i);
        for lk=1:length(ls)
            l=ls(lk); a=(G+l*eye(n))\h; kca=kc*a;
            g(hk,lk,i)=b*mean(kca.^2)+(1-b)*mean((rj*a).^2);
            g(hk,lk,i)=g(hk,lk,i)/2-mean(kca);
        end, end, end
    [gl,ggl]=min(mean(g,3),[],2); [ghl,gghl]=min(gl);
    L=ls(ggl(gghl)); HH=hhs(gghl);
    k=exp(-xx/HH); r=exp(-yx/HH);
    s=r*((b*k'*k/n+(1-b)*r'*r/n+L*eye(n))\ (mean(k)'));
    figure(1); clf; hold on; plot(y,s,'rx');

```

FIGURE 33.9

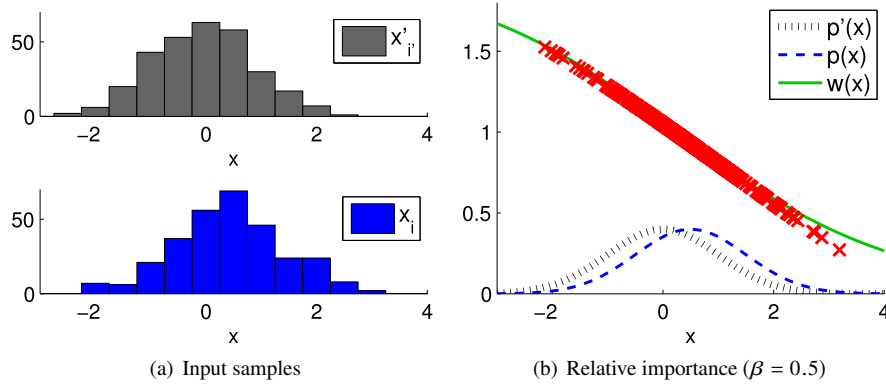
MATLAB code for LS relative density ratio estimation for Gaussian kernel model.

which is called *class-balance weighted LS*. Beyond LS, this class-balance change adaptation technique can be applied to various classification methods introduced in [Part 4](#).

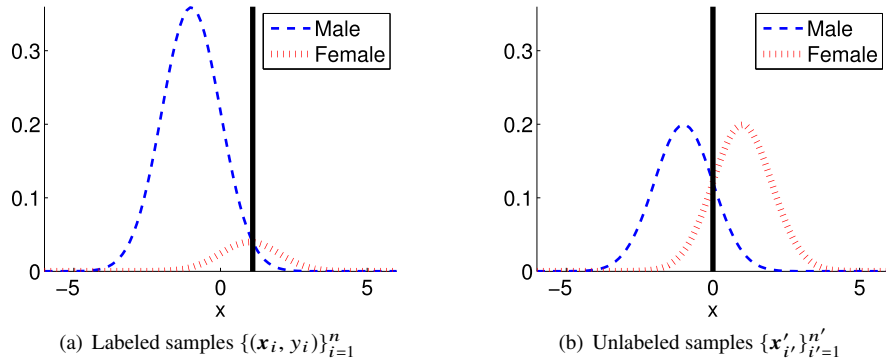
Model selection of class-balance weighted learning may be carried out by *class-balance weighted cross validation*, which is essentially the same as *importance weighted cross validation* introduced in [Section 33.2.3](#), but $p'(y_i)/p(y_i)$ is used as a weight.

33.3.2 CLASS-BALANCE ESTIMATION

To use class-balance weighted learning, the class-prior probabilities $p(y)$ and $p'(y)$ are needed, which are often unknown in practice. For labeled samples $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, estimating the class prior $p(y)$ is straightforward by n_y/n , where n_y denotes the number of samples in class y . However, $p'(y)$ cannot be estimated naively because of

**FIGURE 33.10**

Example of LS relative density ratio estimation. \times 's in the right plot show estimated relative importance values at $\{\mathbf{x}_i\}_{i=1}^n$.

**FIGURE 33.11**

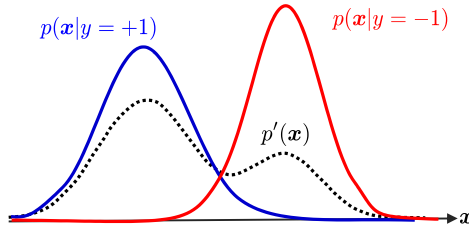
Class-balance change, which affects the decision boundary.

lack of output values $\{y'_{i'}\}_{i'=1}^{n'}$ for $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$. Below, a practical estimator of $p'(y)$ is introduced [38].

The basic idea of estimating $p'(y)$ is to fit a mixture $q_\pi(\mathbf{x})$ of classwise input densities $p(\mathbf{x}|y)$ to $p'(\mathbf{x})$ (Fig. 33.12):

$$q_\pi(\mathbf{x}) = \sum_{y=1}^c \pi_y p(\mathbf{x}|y),$$

where c denotes the number of classes and the coefficient π_y corresponds to $p'(y)$.

**FIGURE 33.12**

Class-prior estimation by distribution matching.

Matching of q_π to p' can be performed, for example, under the *KL divergence* (see Section 14.2),

$$\text{KL}(p' \| q_\pi) = \int p'(x) \log \frac{p'(x)}{q_\pi(x)} dx,$$

or the *L_2 -distance*:

$$L_2(p', q_\pi) = \int (p'(x) - q_\pi(x))^2 dx.$$

The KL divergence can be accurately estimated by *KL density ratio estimation* explained in Section 38.3, while the L_2 -distance can be accurately estimated by *LS density difference estimation* introduced in Section 39.1.3.

Another useful distance measure is the *energy distance* [106], which is the weighted squared distance between characteristic functions:

$$D_E(p', q_\pi) = \int_{\mathbb{R}^d} \|\varphi_{p'}(t) - \varphi_{q_\pi}(t)\|^2 \left(\frac{\pi^{\frac{d+1}{2}}}{\Gamma(\frac{d+1}{2})} \|t\|^{d+1} \right)^{-1} dt,$$

where $\|\cdot\|$ denotes the Euclidean norm, φ_p denotes the *characteristic function* (see Section 2.4.3) of p , $\Gamma(\cdot)$ is the *gamma function* (see Section 4.3), and d denotes the dimensionality of \mathbf{x} . Thanks to the careful design of the weight function, the energy distance can be equivalently expressed as

$$\begin{aligned} D_E(p', q_\pi) &= 2\mathbb{E}_{\mathbf{x}' \sim p', \mathbf{x} \sim q_\pi} \|\mathbf{x}' - \mathbf{x}\| - \mathbb{E}_{\mathbf{x}', \tilde{\mathbf{x}}' \sim p'} \|\mathbf{x}' - \tilde{\mathbf{x}}'\| - \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim q_\pi} \|\mathbf{x} - \tilde{\mathbf{x}}\| \\ &= 2\boldsymbol{\pi}^\top \mathbf{b} - \boldsymbol{\pi}^\top \mathbf{A} \boldsymbol{\pi} + C, \end{aligned}$$

where $\mathbb{E}_{\mathbf{x}' \sim p'}$ denotes the expectation with respect to \mathbf{x}' following density p' and C is a constant independent of $\boldsymbol{\pi}$. \mathbf{A} is the $c \times c$ symmetric matrix and \mathbf{b} is the c -dimensional vector defined as

$$A_{y, \tilde{y}} = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}|y), \tilde{\mathbf{x}} \sim p(\mathbf{x}|\tilde{y})} \|\mathbf{x} - \tilde{\mathbf{x}}\|,$$

```

x=[randn(90,1)-2; randn(10,1)+2] 2*randn(100,1)];
x(:,3)=1; y=[ones(90,1); 2*ones(10,1)]; n=length(y);
X=[randn(10,1)-2; randn(90,1)+2] 2*randn(100,1)];
X(:,3)=1; Y=[ones(10,1); 2*ones(90,1)]; N=length(Y);
x2=sum(x.^2,2); X2=sum(X.^2,2);
xx=sqrt(repmat(x2,1,n)+repmat(x2',n,1)-2*x*x');
xX=sqrt(repmat(x2,1,N)+repmat(X2',n,1)-2*x*X');
for i=1:2
    s(i)=sum(y==i)/n; b(i)=mean(mean(xX(y==i,:)));
    for j=1:2
        A(i,j)=mean(mean(xx(y==i,y==j)));
    end, end
v=(A(1,2)-A(2,2)-b(1)+b(2))/(2*A(1,2)-A(1,1)-A(2,2));
v=min(1,max(0,v)); v(2)=1-v; w=v(y)./s(y); z=2*y-3;
u=x\z; t=(x*(repmat(w',1,size(x,2)).*x))\'(x*(w'.*z));
figure(1); clf; hold on
plot([-5 5],-(u(3)+[-5 5]*u(1))/u(2),'k--');
plot([-5 5],-(t(3)+[-5 5]*t(1))/t(2),'g-');
plot(X(Y==1,1),X(Y==1,2),'bo');
plot(X(Y==2,1),X(Y==2,2),'rx');
legend('Unweighted','Weighted'); axis([-5 5 -10 10])

```

FIGURE 33.13

MATLAB code for class-balance weighted LS.

$$b_y = \mathbb{E}_{\mathbf{x}' \sim p', \mathbf{x} \sim p(\mathbf{x}|y)} \|\mathbf{x}' - \mathbf{x}\|.$$

Although $D_E(p', q_\pi)$ is a concave function with respect to $\pi = (\pi_1, \dots, \pi_c)^\top$, it is a convex function with respect to π_1, \dots, π_{c-1} for $\pi_c = 1 - \sum_{y=1}^{c-1} \pi_y$ and thus its minimizer can be easily obtained [61]. For example, when $c = 2$ with $\pi_1 = \pi$ and $\pi_2 = 1 - \pi$, $D_E(p', q_\pi)$ can be expressed as a function of π up to a constant as

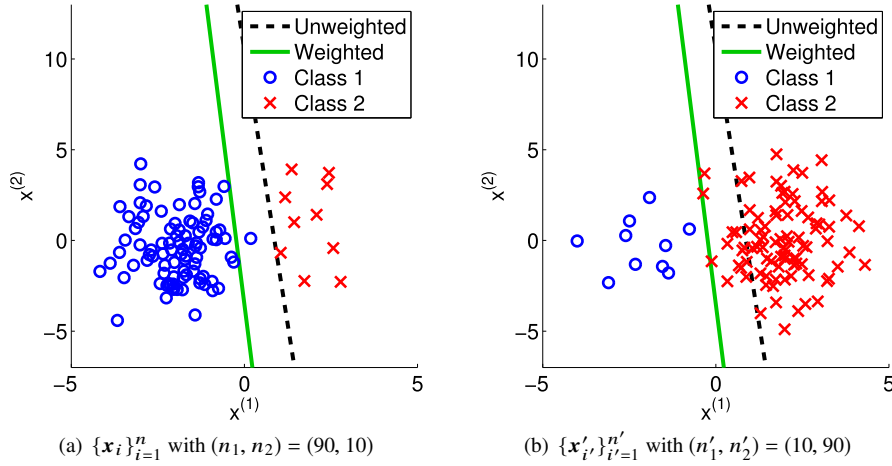
$$J(\pi) = a\pi^2 - 2b\pi,$$

where

$$\begin{aligned} a &= 2A_{1,2} - A_{1,1} - A_{2,2} = D_E(p(\mathbf{x}|y=1), p(\mathbf{x}|y=2)) \geq 0, \\ b &= A_{1,2} - A_{2,2} - b_1 + b_2. \end{aligned}$$

Since $J(\pi)$ is convex with respect to π , its minimizer is given analytically as $\min(1, \max(0, b/a))$. Note that $A_{y,\bar{y}}$ and b_y can be empirically approximated as

$$\hat{A}_{y,\bar{y}} = \frac{1}{n_y n_{\bar{y}}} \sum_{i:y_i=y} \sum_{\bar{i}:y_{\bar{i}}=\bar{y}} \|\mathbf{x}_i - \mathbf{x}_{\bar{i}}\|,$$

**FIGURE 33.14**

Example of class-balance weighted LS. The test class priors are estimated as $\widehat{p}'(y = 1) = 0.18$ and $\widehat{p}'(y = 2) = 0.82$, which are used as weights in class-balance weighted LS.

$$\widehat{b}_y = \frac{1}{n'n_y} \sum_{i'=1}^{n'} \sum_{i:y_i=y} \|\mathbf{x}'_{i'} - \mathbf{x}_i\|.$$

A MATLAB code for class-balance weighted LS is provided in Fig. 33.13, and its behavior is illustrated in Fig. 33.14. This shows that the class prior can be estimated reasonably well and class-balance weighted learning contributes to improving the classification accuracy for test input points $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'}$.