

半教師付き学習(16章) と転移学習(18章)

杉山将・本多淳也

sugi@k.u-tokyo.ac.jp, jhonda@k.u-tokyo.ac.jp

<http://www.ms.k.u-tokyo.ac.jp>

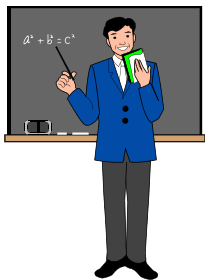
教師付き分類

- **教師付き分類**: クラスラベル付きの訓練データから学習

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

- 精度の良い学習結果を得るためには、多くの訓練データが必要
- しかし、ラベル付きデータの収集にコスト（人手など）がかかる場合、多数のラベル付きデータを集められない

講義の流れ



1. 半教師付き学習(16章)
2. 転移学習(18章)

- ラベルなしのデータは簡単に手に入ることがある

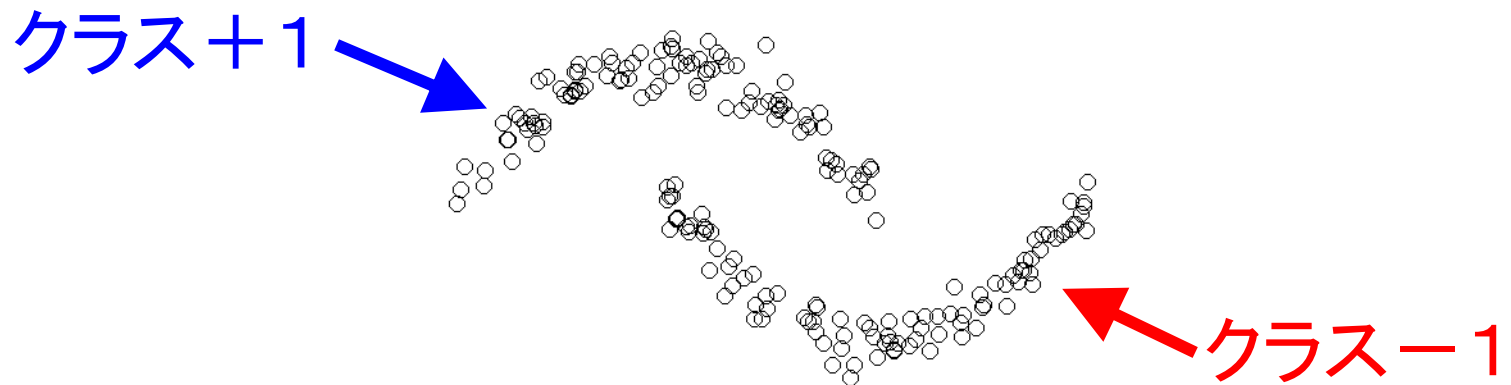
$$\{x_i\}_{i=n+1}^{n+n'}$$

- 顔認識: 顔画像は容易に大量入手できるが、ラベル(性別, 年齢など)付けは人手が必要
- ウェブページ分類: ウェブページは自動的に大量入手できるが、ラベル(政治, 芸能など)付けは人手が必要
- 半教師付き学習: ラベルなしデータも活用し、学習の精度を向上させる

半教師付き学習の大前提

5

- 何も仮定をおかなければ, ラベルなしデータは一般に役に立たない
- 典型的な仮定:
 - 同じ“かたまり”に属するデータは, 同じラベルを持つ



教師付き分類の場合

クラス +1



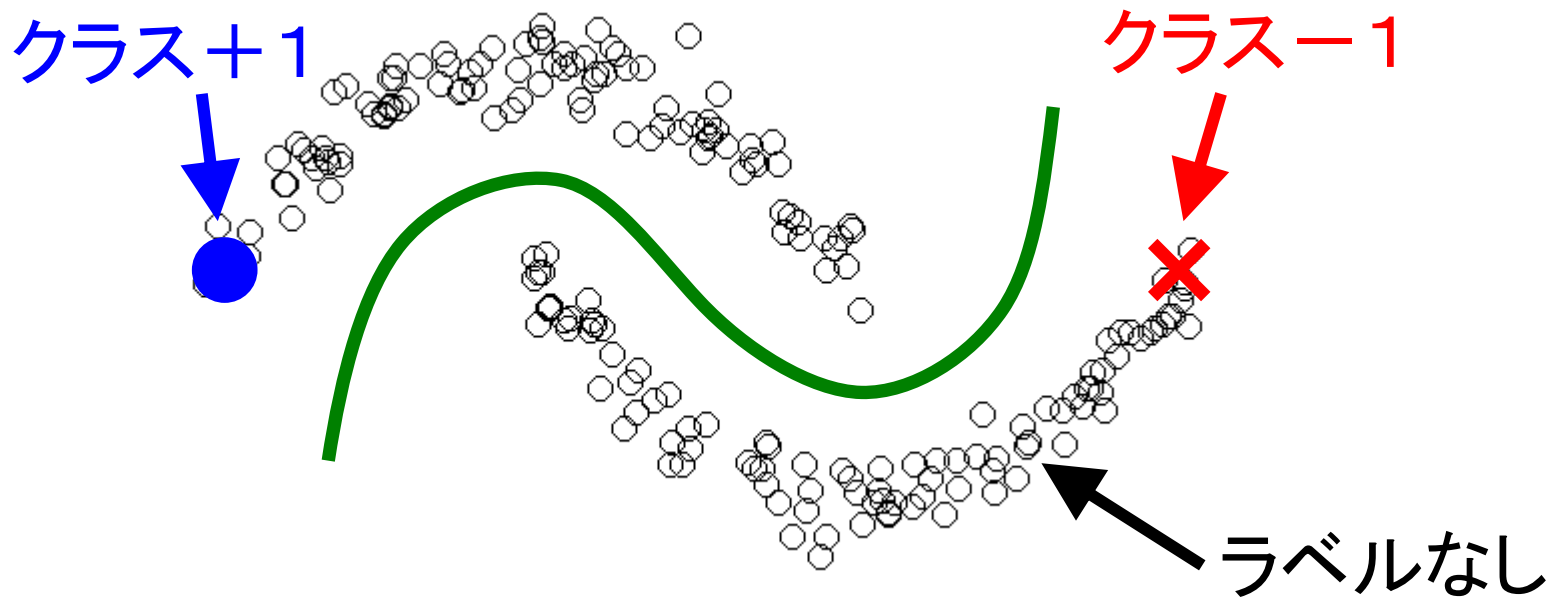
クラス -1



真ん中で分けるのが自然

半教師付き学習の場合

7



ラベルなしデータがなす領域に
沿って分けるのが自然(多様体仮定)

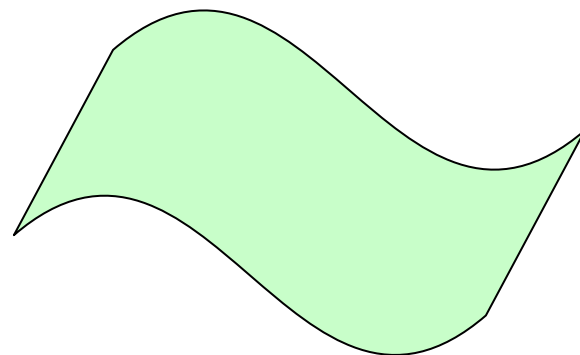
多様体とは

■ 数学的な定義:

- 局所的にユークリッド空間とみなせる空間

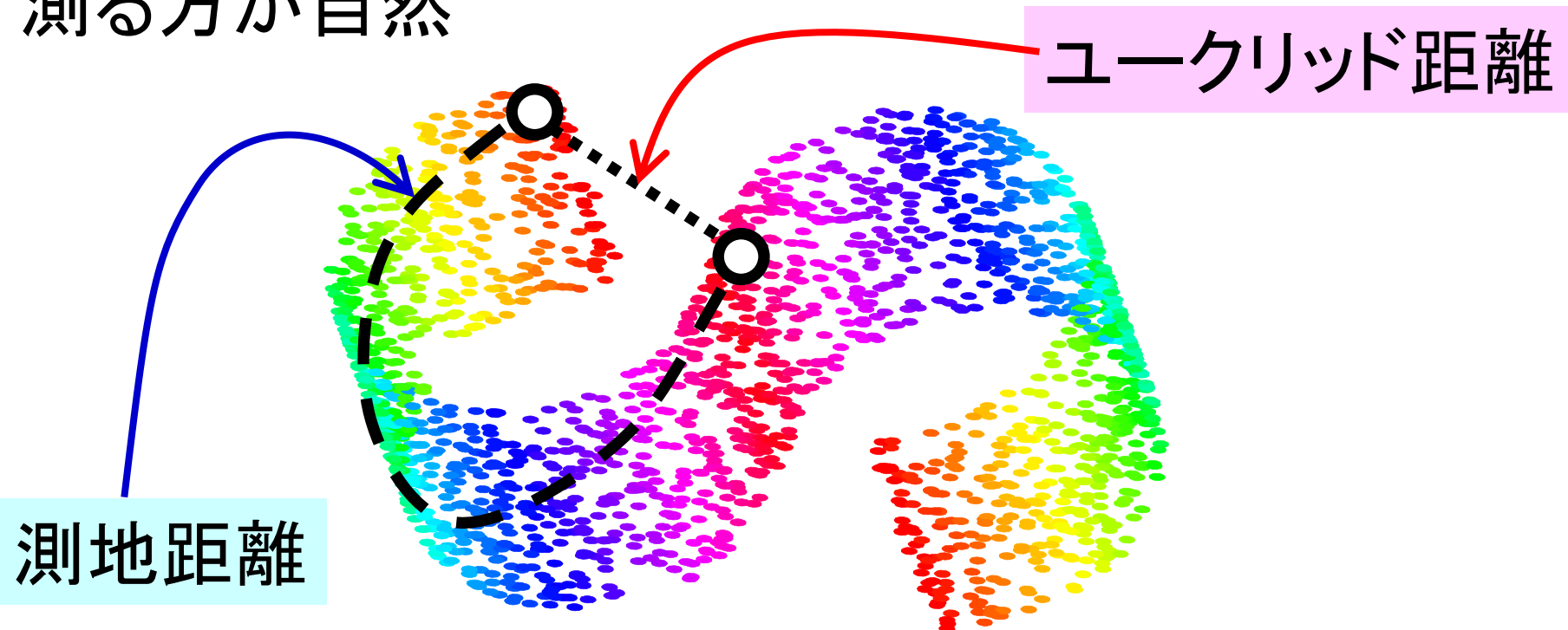
■ 機械学習的な解釈:

- 線形空間の非線形拡張



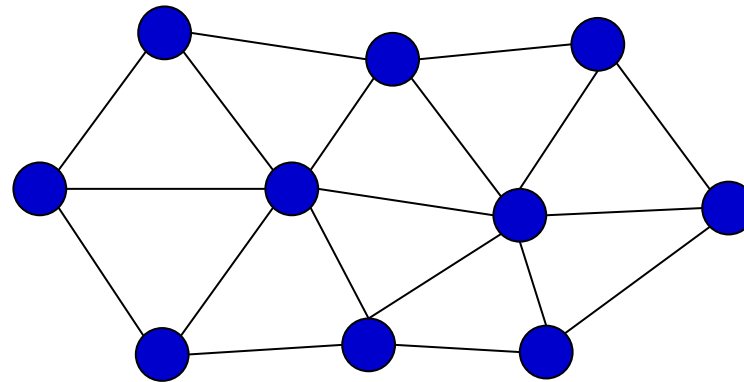
多様体を考える利点

- 高次元データは、空間全体に分布するのではなく、**低次元の多様体上**に分布することが多い。
- ユークリッド距離よりも、多様体上の**測地線**に沿った距離（多様体上での最短距離）で近さを測る方が自然



近傍グラフ

- 多様体は連続的な集合
- 与えられるデータは離散的な集合
- 多様体を, データから作られる
近傍グラフで近似



$$W_{i,i'} = \begin{cases} 1 : x_i \text{ と } x_{i'} \text{ が } k \text{ 近傍} \\ 0 : \text{それ以外} \end{cases}$$

ガウスクアーネルを使うこともある

正則化に基づく半教師付き学習 11

- ラベル付き訓練データ: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$
- ラベルなし訓練データ: $\{\mathbf{x}_i\}_{i=n+1}^{n+n'}$

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

訓練出力に
対する適合の良さ

正則化

ラプラス正則化(後述): 近傍の
入力点間の出力の滑らかさ

$$W_{i,i'} = \begin{cases} 1 & : \mathbf{x}_i \text{ と } \mathbf{x}_{i'} \text{ が } k \text{ 近傍} \\ 0 & : \text{それ以外} \end{cases}$$

■ ラプラス行列 $L = D - W$ および

$a_1, a_2, \dots, a_m \in \mathbb{R}$ について次式を証明せよ

(W は対称行列)

$$\sum_{i,i'=1}^m W_{i,i'} (a_i - a_{i'})^2 = 2 \sum_{i,i'=1}^m L_{i,i'} a_i a_{i'}$$

$$D = \text{diag} \left(\sum_{i=1}^m W_{1,i}, \dots, \sum_{i=1}^m W_{m,i} \right)$$

解答例

13

$$\sum_{i,i'=1}^m W_{i,i'} (a_i - a_{i'})^2$$

$$= 2 \sum_{i,i'=1}^m W_{i,i'} a_i^2 - 2 \sum_{i,i'=1}^m W_{i,i'} a_i a_{i'}$$

$$= 2 \sum_{i=1}^m D_{i,i} a_i^2 - 2 \sum_{i,i'=1}^m W_{i,i'} a_i a_{i'}$$

$$= 2 \sum_{i,i'=1}^m D_{i,i'} a_i a_{i'} - 2 \sum_{i,i'=1}^m W_{i,i'} a_i a_{i'}$$

$$= 2 \sum_{i,i'=1}^m L_{i,i'} a_i a_{i'}$$

ラプラス正則化の変形

14

$$\sum_{i,i'=1}^m W_{i,i'} (a_i - a_{i'})^2 = 2 \sum_{i,i'=1}^m L_{i,i'} a_i a_{i'}$$

■ 線形モデル $f_{\theta}(x) = \theta^{\top} \phi(x) = \phi(x)^{\top} \theta$ を用いると

$$\sum_{i,i'=1}^{n+n'} W_{i,i'} \left(f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_{i'}) \right)^2$$

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \cdots & \phi_b(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_{n+n'}) & \cdots & \phi_b(\mathbf{x}_{n+n'}) \end{pmatrix}$$

$$= \sum_{i,i'=1}^{n+n'} \theta^{\top} \phi(\mathbf{x}_i) L_{i,i'} \phi(\mathbf{x}_{i'})^{\top} \theta$$

$$= \theta^{\top} \left(\sum_{i,i'=1}^{n+n'} \phi(\mathbf{x}_i) L_{i,i'} \phi(\mathbf{x}_{i'})^{\top} \right) \theta = \theta^{\top} \Phi^{\top} L \Phi \theta$$

ラプラス正則化最小二乗分類の 15 解の求め方

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left[\sum_{i=1}^n \left(f_{\theta}(\mathbf{x}_i) - y_i \right)^2 + \lambda \|\theta\|^2 + \nu \sum_{i,i'=1}^{n+n'} W_{i,i'} \left(f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_{i'}) \right)^2 \right]$$

■ 線形モデル $f_{\theta}(\mathbf{x}) = \theta^{\top} \phi(\mathbf{x})$ に対して

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left[\|\tilde{\Phi}\theta - \mathbf{y}\|^2 + \lambda \|\theta\|^2 + 2\nu \theta^{\top} \Phi^{\top} \mathbf{L} \Phi \theta \right]$$

$$= (\tilde{\Phi}^{\top} \tilde{\Phi} + \lambda \mathbf{I} + 2\nu \Phi^{\top} \mathbf{L} \Phi)^{-1} \tilde{\Phi}^{\top} \mathbf{y}$$

$$\Phi = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \cdots & \phi_b(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_{n+n'}) & \cdots & \phi_b(\mathbf{x}_{n+n'}) \end{pmatrix}$$

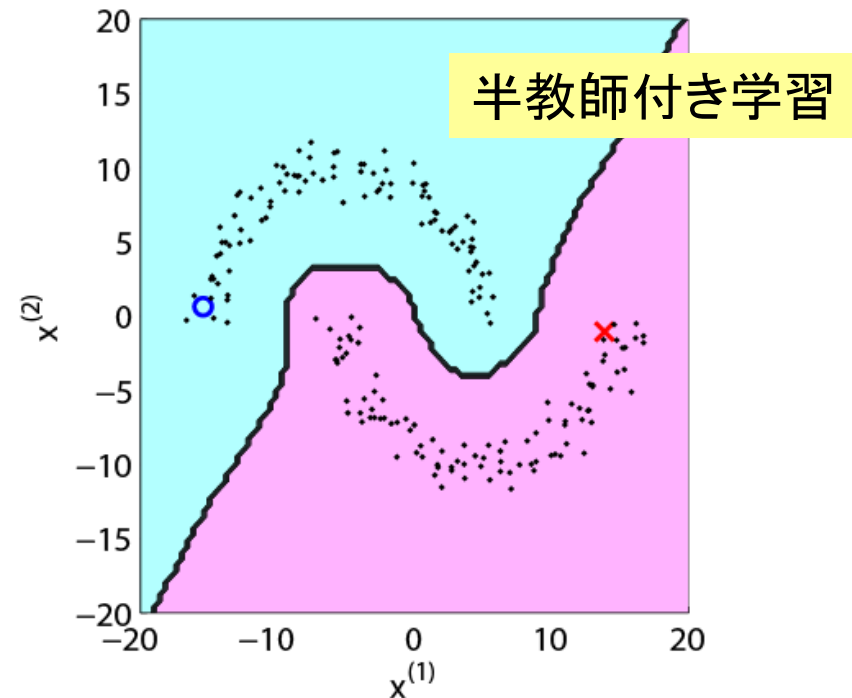
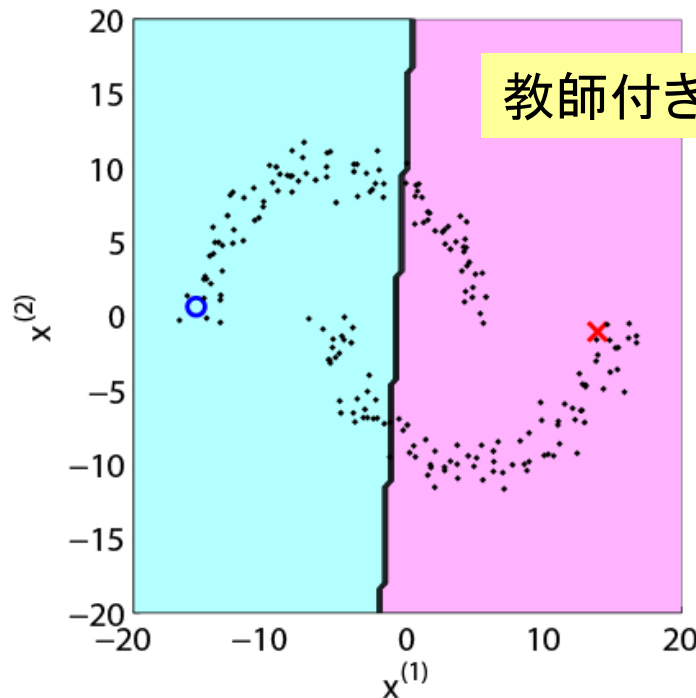
$$\tilde{\Phi} = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \cdots & \phi_b(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\mathbf{x}_n) & \cdots & \phi_b(\mathbf{x}_n) \end{pmatrix}$$

$$\mathbf{y} = (y_1, \dots, y_n)^{\top}$$

実行例

■ ガウスカーネルモデル $f_{\theta}(x) = \sum_{j=1}^{n+n'} \theta_j \exp \left(-\frac{\|x - x_j\|^2}{2h^2} \right)$

に対するラプラス正則化最小二乗分類



■ ラプラス正則化により, ラベルなしデータの多様体に沿った分離境界が得られている

- ラベル付きデータは大量に集めることが困難
- ラベルなしデータは容易に得られる
- 同じ「かたまり」に属するデータは同じラベルを持つという仮定のもと, ラベルなしデータを活用
 - グラフ・ラプラス行列を用いて, ラベルなしデータがなす多様体にそってラベルを伝播
- ラプラス正則化は, 最小二乗分類以外の分類法や, 回帰にも適用可能

講義の流れ



1. 半教師付き学習(16章)

2. 転移学習(18章)

A) 共変量シフト

B) クラスバランス変化

転移学習 (Transfer Learning) 19

- **教師付き学習の大前提**: 訓練標本とテスト標本が同じ確率分布に従う
 - この仮定のもと, 訓練標本から関数を学習し, テスト入力に対する出力を予測
- 応用問題によっては, 訓練標本とテスト標本の確率分布が異なる事がある
- そのような状況でも, 訓練標本から得られる知見を転移させて, テスト入力に対する出力を予測したい

講義の流れ



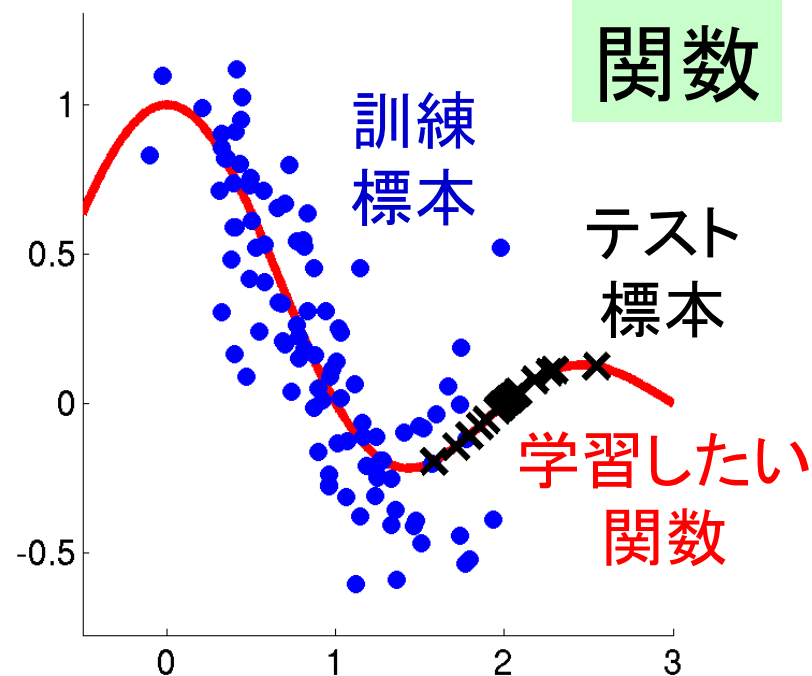
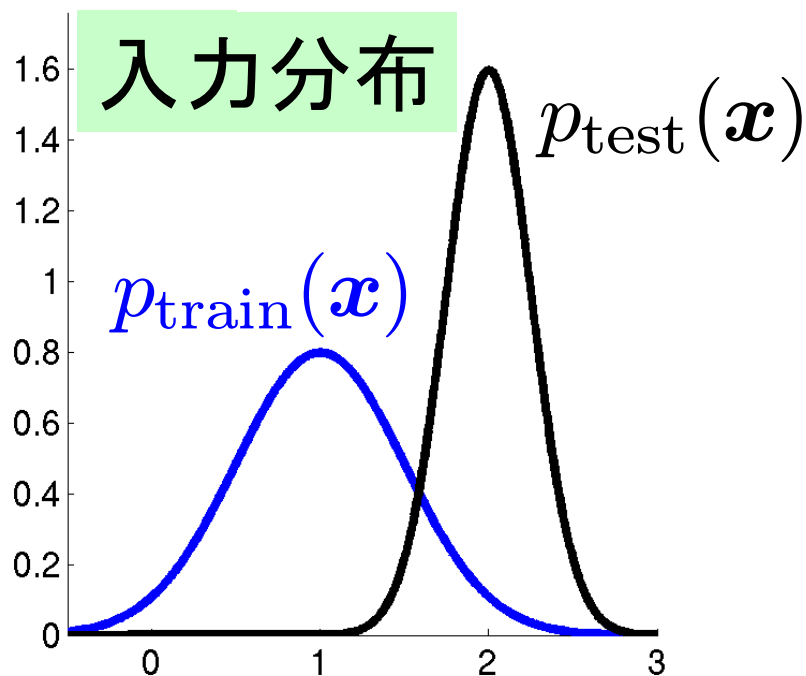
1. 半教師付き学習(16章)

2. 転移学習(18章)

A) 共変量シフト

B) クラスバランス変化

- **共変量**: 入力 of 別名
- **共変量シフト**: 訓練時とテスト時で入力分布が変化するが, 入出力関数は変わらない
- **外挿問題**が典型的な例



■ 医療データ:

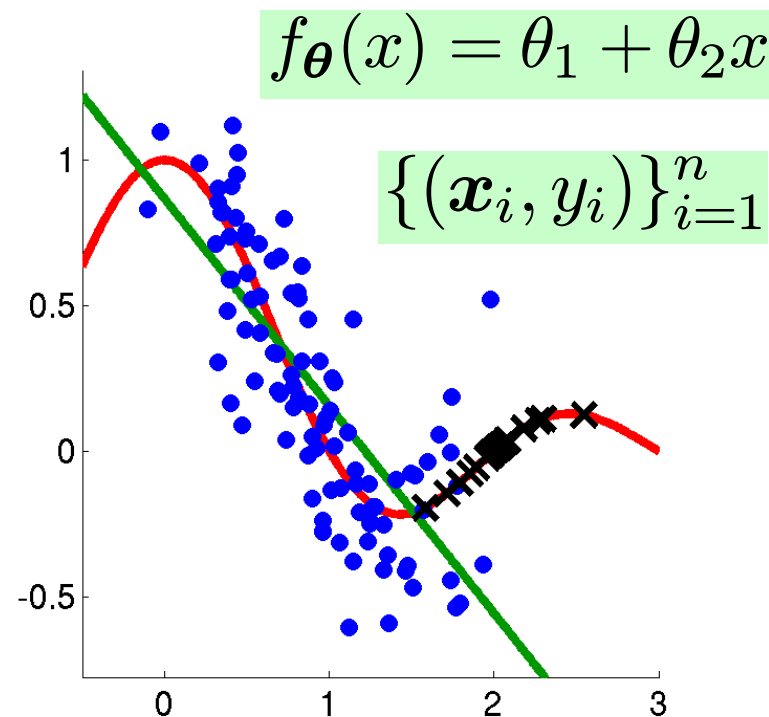
- 訓練データ: 何らかの理由で検査を受けることになった患者の診断結果
- テストデータ: 人間ドック等を受けた一般人の検査データ

■ スпамフィルタ:

- 訓練データ: データベースに蓄積されたスパム/正常メール
- テストデータ: 個々のユーザーに応じたスパム/正常メール

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

- 通常の設定では、最小二乗法は**一**致性を持つ。
つまり、 $n \rightarrow \infty$ で(モデル内で)最適な解が求まる
- しかし共変量シフト下では、モデルが真の関数を含まない限り最小二乗法は**一**致性をもたない



最小二乗法の一貫性：大数の法則⁴

- 標本平均は真の期待値に収束：

$$\frac{1}{n} \sum_{i=1}^n \text{loss}_{\theta}(\mathbf{x}_i) \longrightarrow \int \text{loss}_{\theta}(\mathbf{x}) p_{\text{train}}(\mathbf{x}) d\mathbf{x}$$

$$\mathbf{x}_i \stackrel{\text{i.i.d.}}{\sim} p_{\text{train}}(\mathbf{x})$$

- 共変量シフトがなければ，このことから一貫性が保証される
- 訓練データ $\{\mathbf{x}_i\}_{i=1}^n$ を用いて，テストデータに対する期待値を最小にしたい！

$$\int \text{loss}_{\theta}(\mathbf{x}) p_{\text{test}}(\mathbf{x}) d\mathbf{x}$$

重点(importance)サンプリング 25

- **重要度**: テスト入力と訓練入力の密度の比

$$\frac{p_{\text{test}}(\boldsymbol{x})}{p_{\text{train}}(\boldsymbol{x})}$$

- 通常のアVERAGE損失のかわりに**重要度重み付き平均**を考えると, 大数の法則により

$$\frac{1}{n} \sum_{i=1}^n \frac{p_{\text{test}}(\boldsymbol{x}_i)}{p_{\text{train}}(\boldsymbol{x}_i)} \text{loss}_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \quad \boldsymbol{x}_i \stackrel{\text{i.i.d.}}{\sim} p_{\text{train}}(\boldsymbol{x})$$

$$\longrightarrow \int \frac{p_{\text{test}}(\boldsymbol{x})}{p_{\text{train}}(\boldsymbol{x})} \text{loss}_{\boldsymbol{\theta}}(\boldsymbol{x}) p_{\text{train}}(\boldsymbol{x}) d\boldsymbol{x}$$

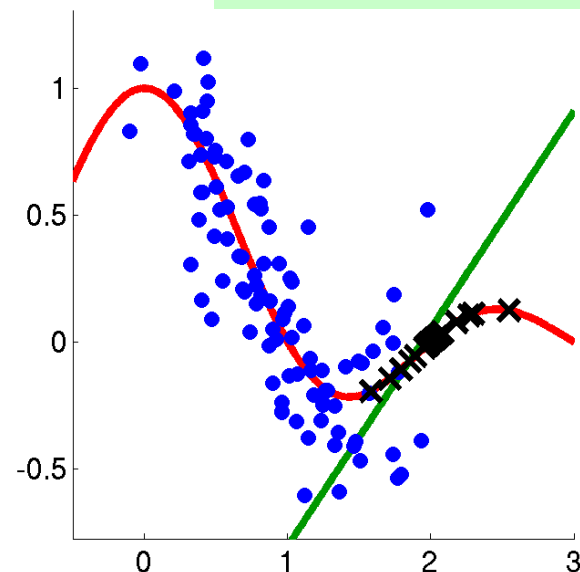
$$= \int \text{loss}_{\boldsymbol{\theta}}(\boldsymbol{x}) p_{\text{test}}(\boldsymbol{x}) d\boldsymbol{x}$$

$$\min_{\theta} \sum_{i=1}^n \frac{p_{\text{test}}(\mathbf{x}_i)}{p_{\text{train}}(\mathbf{x}_i)} \left(f_{\theta}(\mathbf{x}_i) - y_i \right)^2$$

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

$$\{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

- 重要度重み付き最小二乗法は
共変量シフト下でも一貫性を持つ
- 重要度重み付けは多くの学習法
に適用できる:
 - サポートベクトルマシン
 - ロジスティック回帰
 - 条件付き確率場など



- 重要度重み付き最小二乗法では, **重要度**

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

が必要

- **単純な方法**: 確率密度 $p_{\text{train}}(\boldsymbol{x})$ と $p_{\text{test}}(\boldsymbol{x})$ を 訓練データとテストデータからそれぞれ推定
- しかし, 推定した確率密度での割り算は推定誤差を増幅させてしまう
- 重要度を直接推定したい

■ データ:

- 訓練入力: $\{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} p_{\text{train}}(\mathbf{x})$

- テスト入力: $\{\mathbf{x}'_{i'}\}_{i'=1}^{n'} \stackrel{i.i.d.}{\sim} p_{\text{test}}(\mathbf{x})$

■ 真の重要度 $w(\mathbf{x})$ との二乗誤差を最小にする
ように重要度モデル $w_{\alpha}(\mathbf{x})$ を学習: $\min_{\alpha} [J(\alpha)]$

$$J(\alpha) = \frac{1}{2} \int \left(w_{\alpha}(\mathbf{x}) - w(\mathbf{x}) \right)^2 p_{\text{train}}(\mathbf{x}) d\mathbf{x}$$

$$w(\mathbf{x}) = \frac{p_{\text{test}}(\mathbf{x})}{p_{\text{train}}(\mathbf{x})}$$

$$= \frac{1}{2} \int w_{\alpha}(\mathbf{x})^2 p_{\text{train}}(\mathbf{x}) d\mathbf{x} - \int w_{\alpha}(\mathbf{x}) p_{\text{test}}(\mathbf{x}) d\mathbf{x} + C$$

$$\approx \frac{1}{2n} \sum_{i=1}^n w_{\alpha}(\mathbf{x}_i)^2 - \frac{1}{n'} \sum_{i'=1}^{n'} w_{\alpha}(\mathbf{x}'_{i'}) + C$$

■ 重要度モデル: $w_{\alpha}(\mathbf{x}) = \sum_{j=1}^{n'} \alpha_j \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'_j\|^2}{2\sigma^2} \right)$

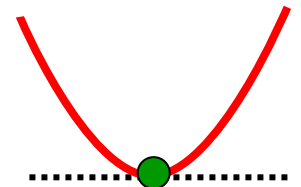
■ 最適化規準: $\min_{\alpha} \left[\frac{1}{2} \alpha^{\top} \hat{G} \alpha - \hat{h}^{\top} \alpha + \underbrace{\frac{\lambda}{2} \alpha^{\top} \alpha}_{\text{正則化項}} \right]$

$$\hat{G}_{j,j'} = \frac{1}{n} \sum_{i=1}^n \exp \left(-\frac{\|\mathbf{x}_i - \mathbf{x}'_j\|^2}{2\sigma^2} \right) \exp \left(-\frac{\|\mathbf{x}_i - \mathbf{x}'_{j'}\|^2}{2\sigma^2} \right)$$

$$\hat{h}_j = \frac{1}{n'} \sum_{i'=1}^{n'} \exp \left(-\frac{\|\mathbf{x}'_{i'} - \mathbf{x}'_j\|^2}{2\sigma^2} \right)$$

■ 大域的最適解が解析的に計算可能:

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



■ σ, λ は二乗誤差 J に関する交差確認により決定

重要度の推定例

■ ガウスカーネル重要度モデル

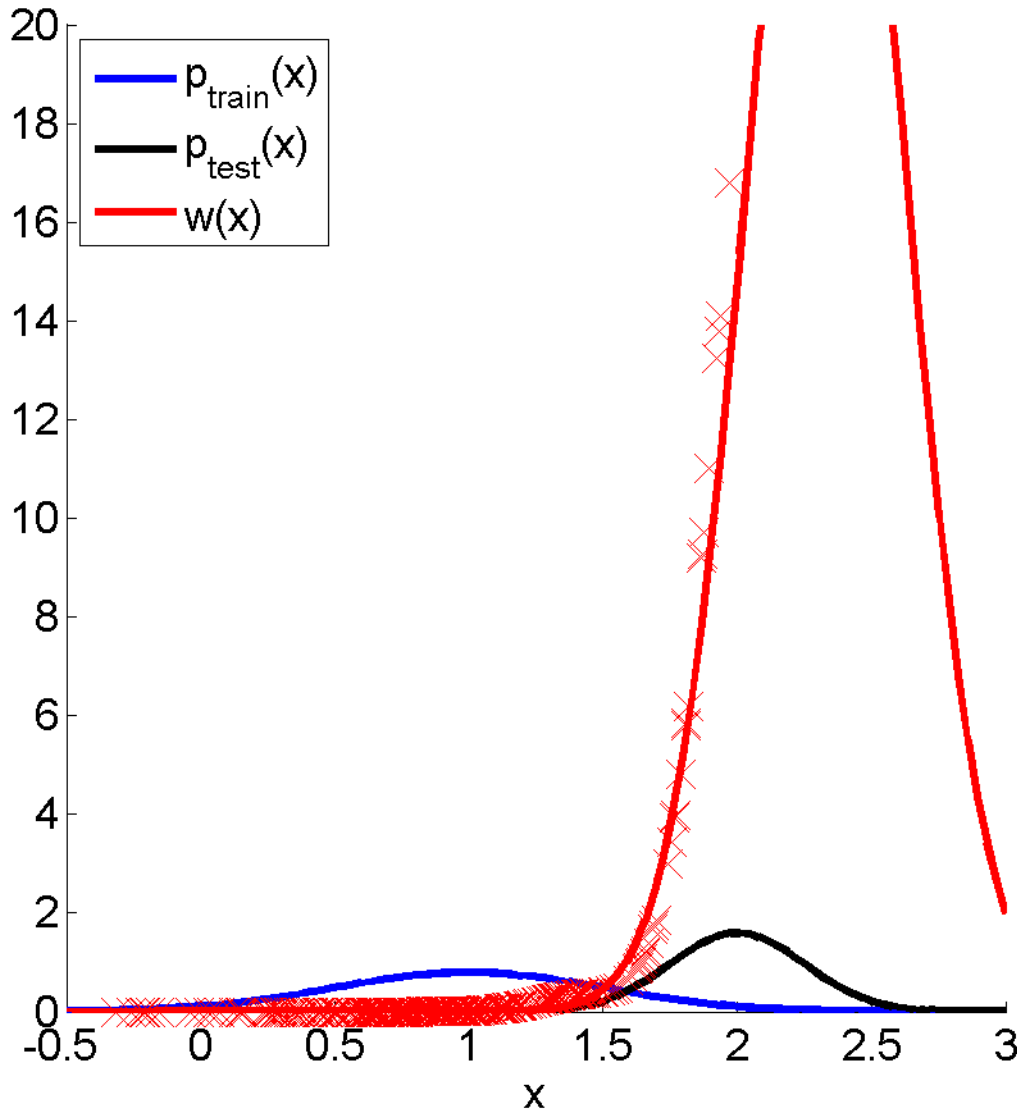
$$w_{\alpha}(\boldsymbol{x}) = \sum_{j=1}^{n'} \alpha_j \exp \left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'_j\|^2}{2\sigma^2} \right)$$

に対する最小二乗重要度推定法

```
clear all; rand('state',0); randn('state',0);  
n=100; x=randn(n,1)/4+1; u=randn(n,1)/2;  
x2=x.^2; xx=repmat(x2,1,n)+repmat(x2',n,1)-2*x*x';  
u2=u.^2; ux=repmat(u2,1,n)+repmat(x2',n,1)-2*u*x';  
  
k=exp(-xx/0.1); r=exp(-ux/0.1);  
w=r*((r'*r/n+0.1*eye(n))\ (mean(k)'));  
figure(1); clf; hold on; plot(u,w,'rx');
```

重要度の推定例(続き)

31



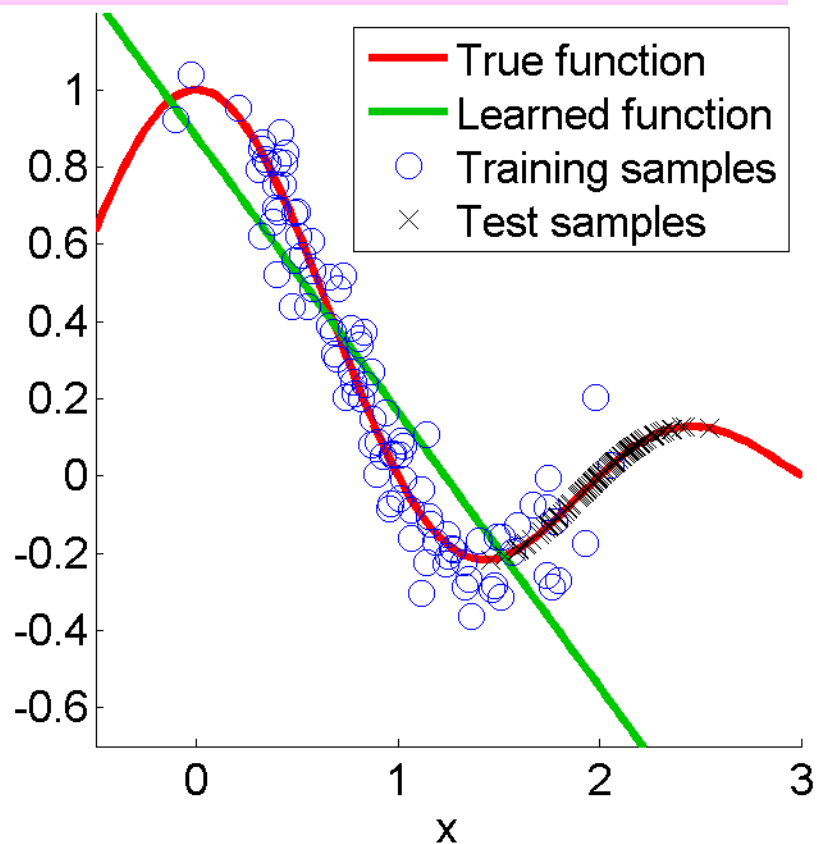
$$w(\boldsymbol{x}) = \frac{p_{\text{test}}(\boldsymbol{x})}{p_{\text{train}}(\boldsymbol{x})}$$

■ 重要度が
精度良く
推定できている

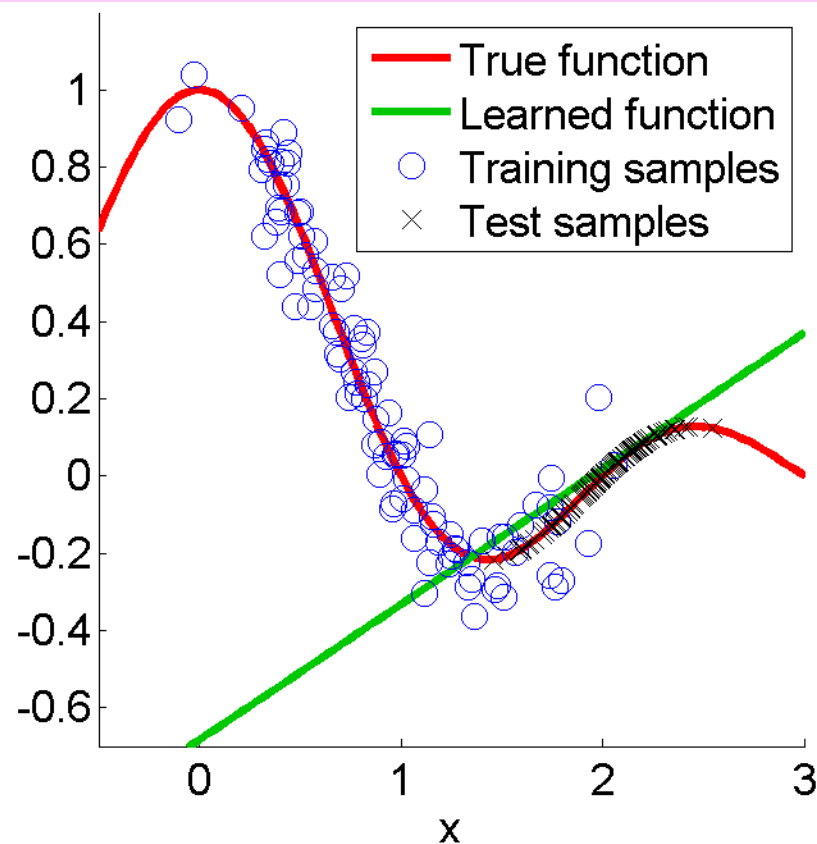
重要度重み付き最小二乗法 の実行例

32

$$\min_{\theta} \left[\sum_{i=1}^n \left(f_{\theta}(\mathbf{x}_i) - y_i \right)^2 \right]$$



$$\min_{\theta} \left[\sum_{i=1}^n \frac{p_{\text{test}}(\mathbf{x}_i)}{p_{\text{train}}(\mathbf{x}_i)} \left(f_{\theta}(\mathbf{x}_i) - y_i \right)^2 \right]$$



■ 転移学習:

- 異なる確率分布に従う訓練標本を用いて, テスト標本の出力を精度良く予測したい

■ 共変量シフト:

- 入力 of 確率分布のみが変化する

■ 重要度 (確率密度の比) で重みを付けて学習すれば, 確率分布の違いを補正できる

■ 重要度は, 個々の確率密度を推定することなく, 直接推定できる

講義の流れ



1. 半教師付き学習(16章)

2. 転移学習(18章)

A) 共変量シフト

B) クラスバランス変化

クラス比変化

■ 顔画像からの性別予測

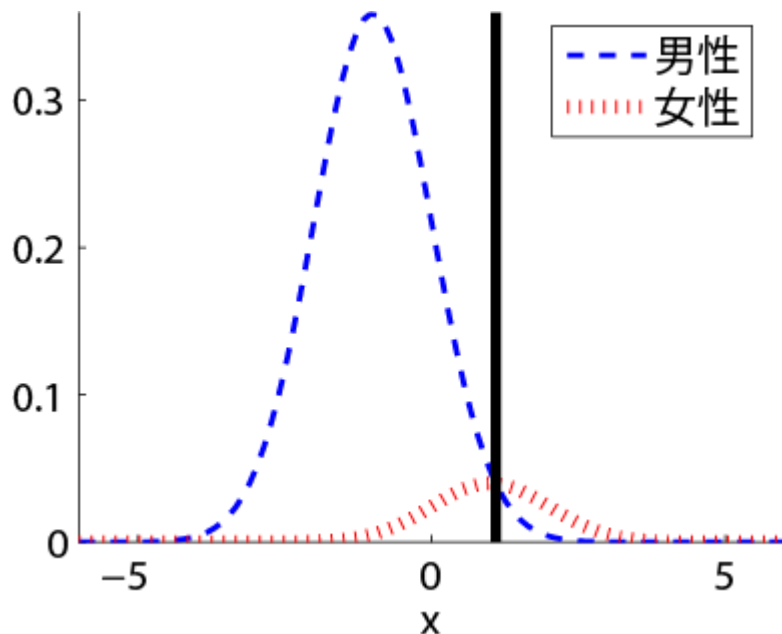
- 訓練標本: 大学には**女性が少ない**
- テスト標本: 世間での男女比は**一対一**



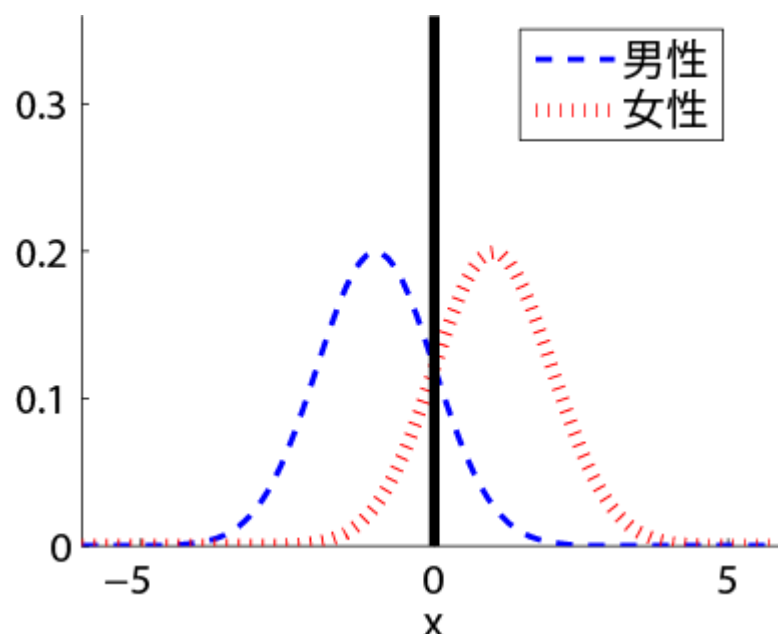
The Yale Face Database B

■ 識別境界がずれる

訓練標本



テスト標本



- テスト標本と比べて，訓練標本中で
 - 比率の小さいクラスの標本には大きな重みを与える
 - 比率の大きいクラスの標本には小さな重みを与える

$$\min_{\theta} \sum_{i=1}^n \frac{p_{\text{test}}(y_i)}{p_{\text{train}}(y_i)} \left(f_{\theta}(\mathbf{x}_i) - y_i \right)^2$$

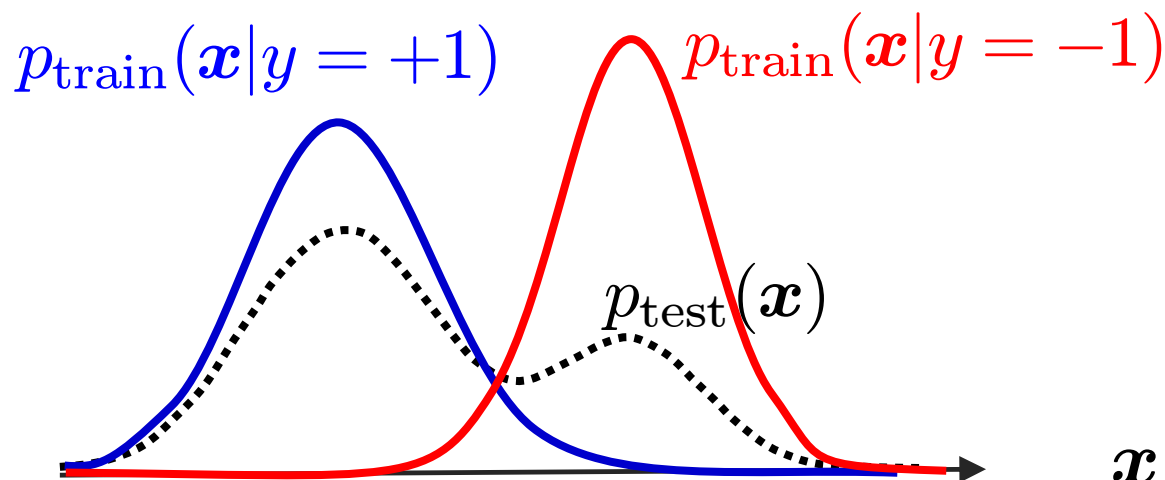
- テスト標本のクラスバランスがわからない場合はどうすればよいか？

クラス比の推定

- テスト入力の分布 p_{test} が各クラスの訓練入力の分布の混合分布 q_{π} で表されると仮定すると、真の混合比 π は(何らかの)距離尺度を最小化:

$$\min_{\pi \in [0,1]} \text{dist} \left(p_{\text{test}} \parallel q_{\pi} \right)$$

$$q_{\pi}(\mathbf{x}) = \pi p_{\text{train}}(\mathbf{x}|y = +1) + (1 - \pi) p_{\text{train}}(\mathbf{x}|y = -1)$$



確率分布間の距離

■ エネルギー距離 (の二乗) :

$$D_E^2(p_{\text{test}}, q_\pi) = 2\mathbb{E}_{\mathbf{x}' \sim p_{\text{test}}, \mathbf{x} \sim q_\pi} \|\mathbf{x}' - \mathbf{x}\|$$

$$- \mathbb{E}_{\mathbf{x}', \tilde{\mathbf{x}}' \sim p_{\text{test}}} \|\mathbf{x}' - \tilde{\mathbf{x}}'\| - \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim q_\pi} \|\mathbf{x} - \tilde{\mathbf{x}}\|$$

- $D_E(p_{\text{test}}, q_\pi) \geq 0$
- $D_E(p_{\text{test}}, q_\pi) = 0 \iff p_{\text{test}} = q_\pi$

$$q_\pi(\mathbf{x}) = \pi p_{\text{train}}(\mathbf{x}|y = +1) + (1 - \pi)p_{\text{train}}(\mathbf{x}|y = -1)$$

■ 期待値で定義されるため, 標本平均で簡単に近似できる

$$D_E^2(p_{\text{test}}, q_\pi) = 2\mathbb{E}_{\mathbf{x}' \sim p_{\text{test}}, \mathbf{x} \sim q_\pi} \|\mathbf{x}' - \mathbf{x}\|$$

$$- \mathbb{E}_{\mathbf{x}', \tilde{\mathbf{x}}' \sim p_{\text{test}}} \|\mathbf{x}' - \tilde{\mathbf{x}}'\| - \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim q_\pi} \|\mathbf{x} - \tilde{\mathbf{x}}\|$$

$$q_\pi(\mathbf{x}) = \pi p_{\text{train}}(\mathbf{x}|y = +1) + (1 - \pi)p_{\text{train}}(\mathbf{x}|y = -1)$$

■ エネルギー距離(の二乗)は, π の関数として次のように表現できる:

証明は宿題

$$J(\pi) = (2A_{+1,-1} - A_{+1,+1} - A_{-1,-1})\pi^2$$

$$-2(A_{+1,-1} - A_{-1,-1} - b_{+1} + b_{-1})\pi + \text{Const.}$$

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

$$b_y = \mathbb{E}_{\mathbf{x}' \sim p_{\text{test}}, \mathbf{x} \sim p_{\text{train}}(\mathbf{x}|y)} \|\mathbf{x}' - \mathbf{x}\|$$

エネルギー距離の標本近似

40

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

$$b_y = \mathbb{E}_{\mathbf{x}' \sim p_{\text{test}}, \mathbf{x} \sim p_{\text{train}}(\mathbf{x}|y)} \|\mathbf{x}' - \mathbf{x}\|$$

■ 期待値を標本平均で近似すれば,

$$\hat{J}(\pi) = (2\hat{A}_{+1,-1} - \hat{A}_{+1,+1} - \hat{A}_{-1,-1})\pi^2$$

$$-2(\hat{A}_{+1,-1} - \hat{A}_{-1,-1} - \hat{b}_{+1} + \hat{b}_{-1})\pi + \text{const.}$$

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

$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{\text{i.i.d.}}{\sim} p_{\text{train}}(\mathbf{x})$$

$$\{\mathbf{x}'_{i'}\}_{i'=1}^{n'} \stackrel{\text{i.i.d.}}{\sim} p_{\text{test}}(\mathbf{x})$$

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

n_y : クラス y の標本数

エネルギー距離に基づく クラス比の推定

$$\begin{aligned}\hat{J}(\pi) = & (2\hat{A}_{+1,-1} - \hat{A}_{+1,+1} - \hat{A}_{-1,-1})\pi^2 \\ & - 2(\hat{A}_{+1,-1} - \hat{A}_{-1,-1} - \hat{b}_{+1} + \hat{b}_{-1})\pi + \text{const.}\end{aligned}$$

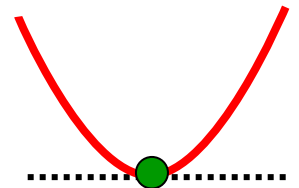
■ クラス比の推定値

$$\hat{\pi} = \underset{0 \leq \pi \leq 1}{\operatorname{argmin}} J(\pi)$$

は解析的に求められる:

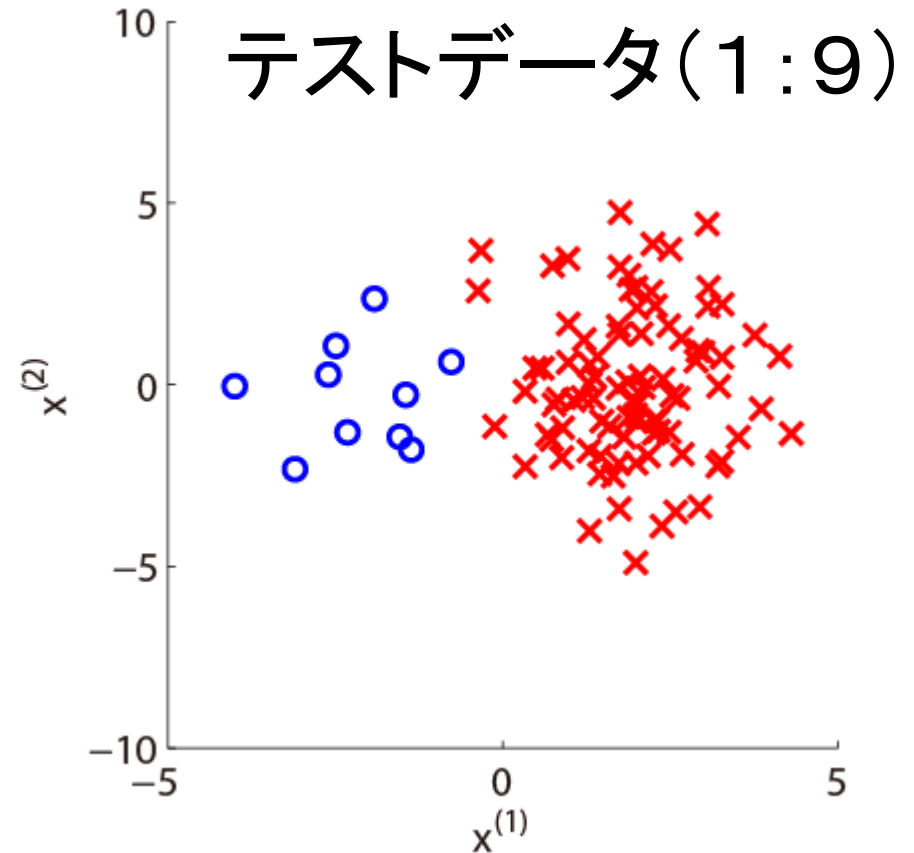
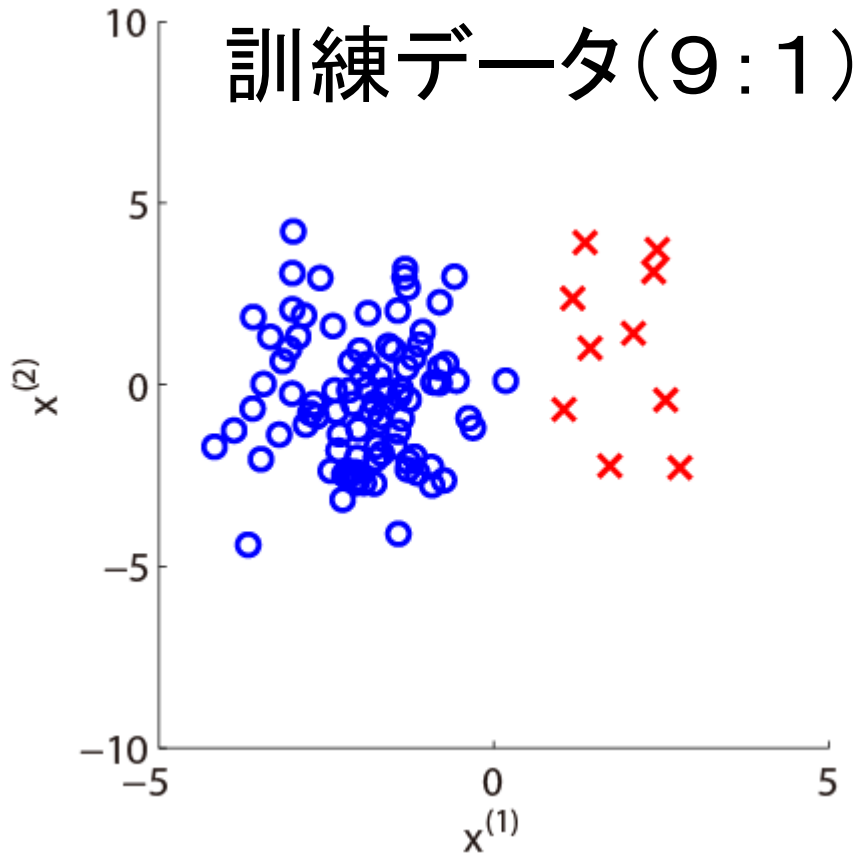
$$\tilde{\pi} = \frac{\hat{A}_{+1,-1} - \hat{A}_{-1,-1} - \hat{b}_{+1} + \hat{b}_{-1}}{2\hat{A}_{+1,-1} - \hat{A}_{+1,+1} - \hat{A}_{-1,-1}}$$

$$\hat{\pi} = \min(1, \max(0, \tilde{\pi}))$$



クラス比の推定例

42



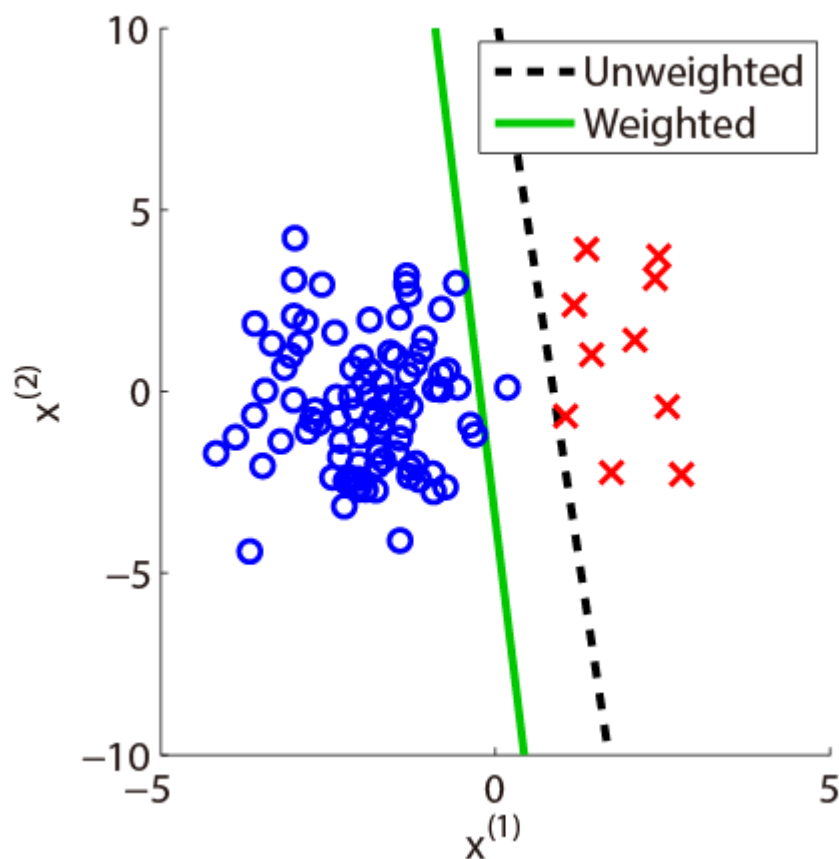
■ テストデータのクラスバランス推定結果

(0.18:0.82)

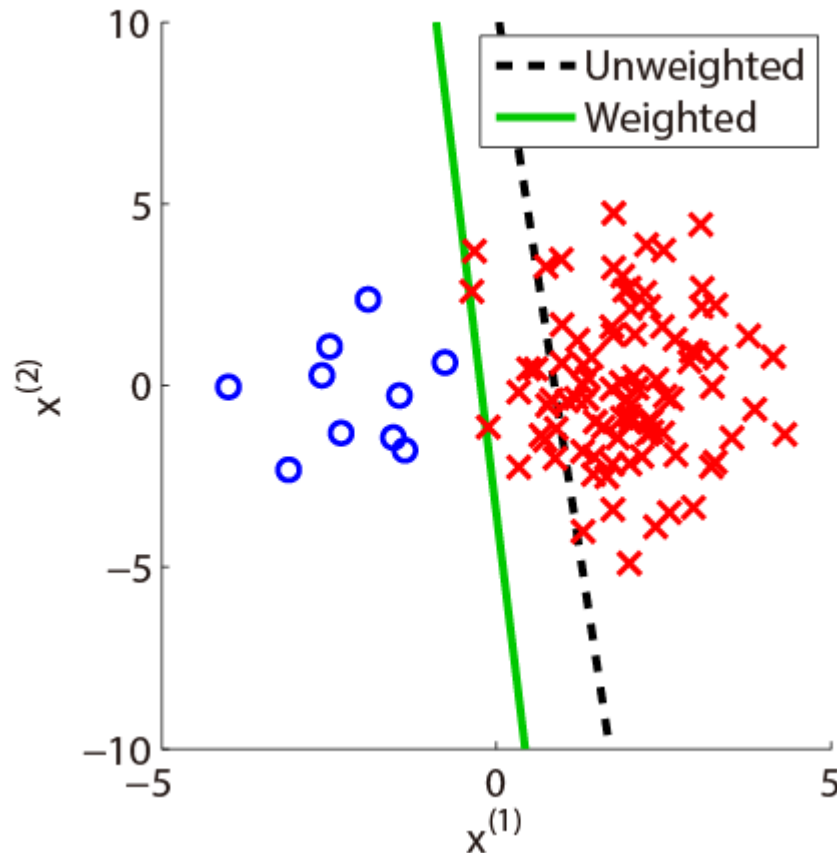
クラス比重み付き 最小二乗法の実行例

- テストデータに対する分類精度が向上している

訓練データ



テストデータ



■ 転移学習：

- 異なる確率分布に従う訓練標本を用いて、テスト標本の出力を精度良く予測したい

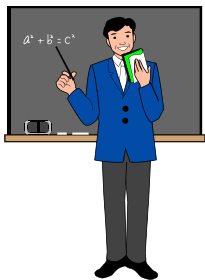
■ クラス比変化：

- クラスのバランスのみが変化する

■ 重要度（クラス比）で重みを付けて学習すれば確率分布の違いを補正できる

■ 重要度は、適当な距離尺度のもとでの適合によって推定できる

講義の流れ

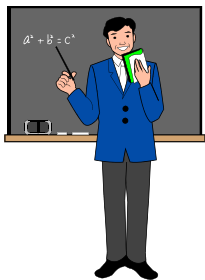


1. 半教師付き学習(16章)
2. 転移学習(18章)

- **半教師付き学習** : ラベルなしデータを活用する
 - 同じかたまりに属するデータは同じクラスに属すると仮定
 - **ラプラス正則化**によりラベルを伝播
- **転移学習** : 訓練時とテスト時でデータの生成分布が変化
 - **共変量シフト** : 入力分布が変化
 - **クラス比変化** : 各クラスの標本のバランスが変化
 - **重要度重み付き学習**により適応
 - 確率分布を推定せず, 重要度を直接推定する

次回の予告

- 線形次元削減(13章, 17章)

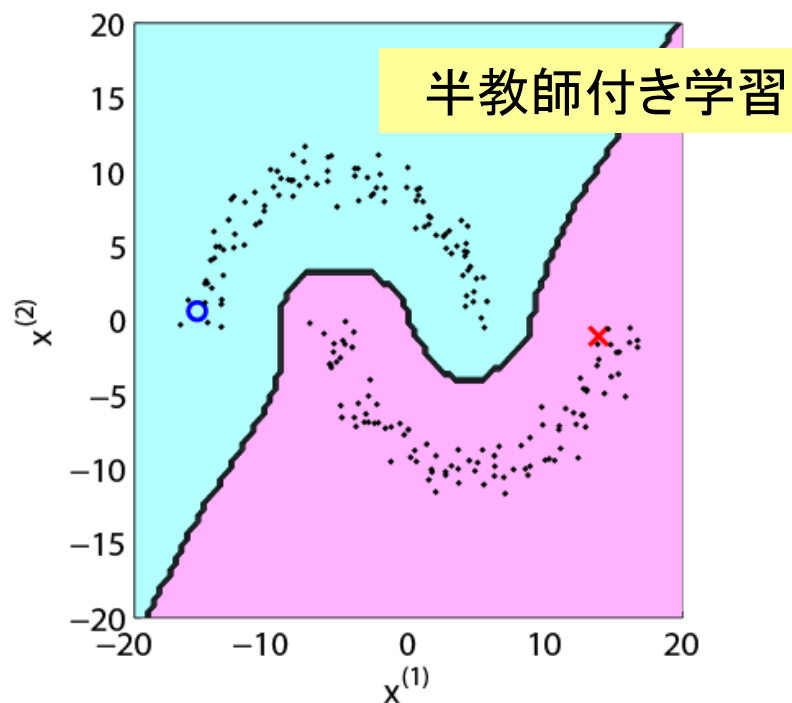
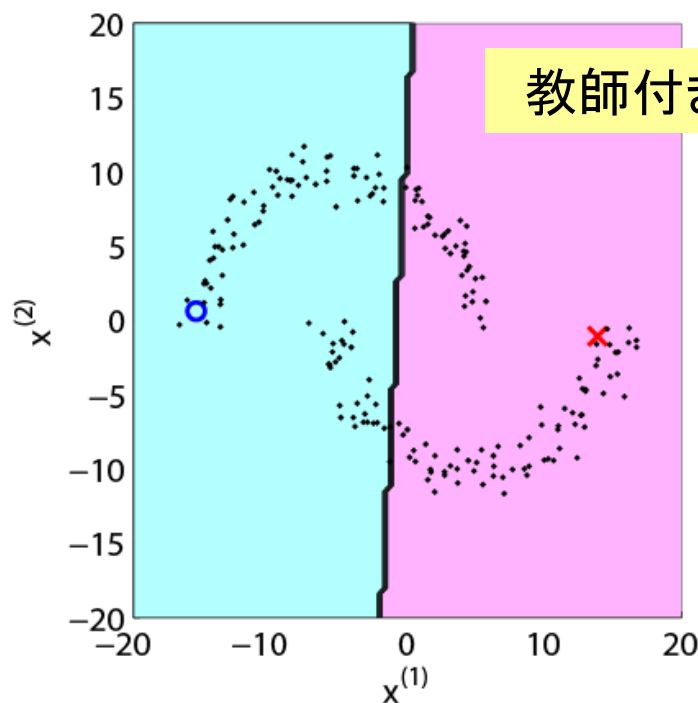


■ ガウスカーネルモデル

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

$$K(\boldsymbol{x}, \boldsymbol{c}) = \exp \left(-\frac{\|\boldsymbol{x} - \boldsymbol{c}\|^2}{2h^2} \right)$$

に対してラプラス正則化最小二乗分類を実装せよ



$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left[\sum_{i=1}^n \left(f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i \right)^2 + \lambda \|\boldsymbol{\theta}\|^2 + \nu \sum_{i,i'=1}^{n+n'} W_{i,i'} \left(f_{\boldsymbol{\theta}}(\mathbf{x}_i) - f_{\boldsymbol{\theta}}(\mathbf{x}_{i'}) \right)^2 \right]$$

- 正則化項には例えば $\lambda = \nu = 1$ を用いてよい
- 近傍グラフの重みにはガウスカーネル

$$W_{i,i'} = \exp \left(-\frac{\|\mathbf{x}_i - \mathbf{x}_{i'}\|^2}{2h^2} \right)$$

を用いてよい

宿題1(続き)

50

```
clear all; rand('state',0); randn('state',0);
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]; % Training input
y=zeros(n,1); y(1)=-1; y(n)=1; % Training output
hh=2*1^2; % Gauss width (2h^2)
```

LRLS.m
を実装せよ

```
t=LRLS(x,y,hh); % Parameter learning
```

```
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.');
```

$$D_{\text{E}}^2(p_{\text{test}}, q_{\pi}) = 2\mathbb{E}_{\mathbf{x}' \sim p_{\text{test}}, \mathbf{x} \sim q_{\pi}} \|\mathbf{x}' - \mathbf{x}\| \\ - \mathbb{E}_{\mathbf{x}', \tilde{\mathbf{x}}' \sim p_{\text{test}}} \|\mathbf{x}' - \tilde{\mathbf{x}}'\| - \mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim q_{\pi}} \|\mathbf{x} - \tilde{\mathbf{x}}\|$$

$$q_{\pi}(\mathbf{x}) = \pi p_{\text{train}}(\mathbf{x}|y = +1) + (1 - \pi)p_{\text{train}}(\mathbf{x}|y = -1)$$

■ $D_{\text{E}}^2(p_{\text{test}}, q_{\pi})$ の π に関する以下の表現を導け

$$J(\pi) = (2A_{+1,-1} - A_{+1,+1} - A_{-1,-1})\pi^2 \\ - 2(A_{+1,-1} - A_{-1,-1} - b_{+1} + b_{-1})\pi + \text{Const.}$$

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

$$b_y = \mathbb{E}_{\mathbf{x}' \sim p_{\text{test}}, \mathbf{x} \sim p_{\text{train}}(\mathbf{x}|y)} \|\mathbf{x}' - \mathbf{x}\|$$

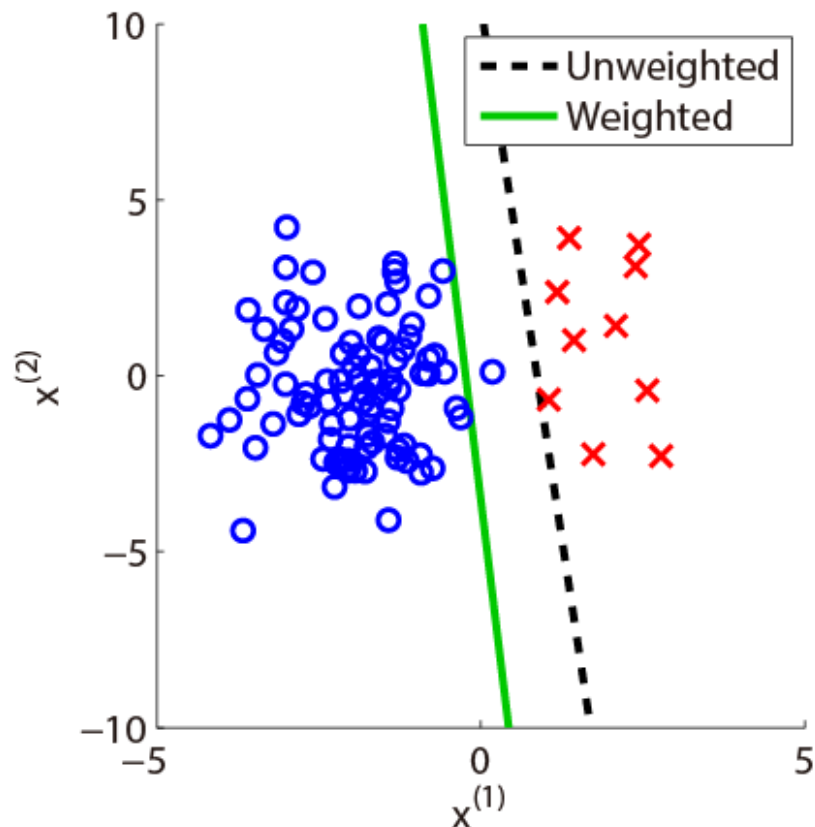
■ ヒント: $\mathbb{E}_{\tilde{\mathbf{x}} \sim q_{\pi}} [f(\tilde{\mathbf{x}})] = \pi \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{\text{train}}(\tilde{\mathbf{x}}|+1)} [f(\tilde{\mathbf{x}})] \\ + (1 - \pi) \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{\text{train}}(\tilde{\mathbf{x}}|-1)} [f(\tilde{\mathbf{x}})]$

宿題3

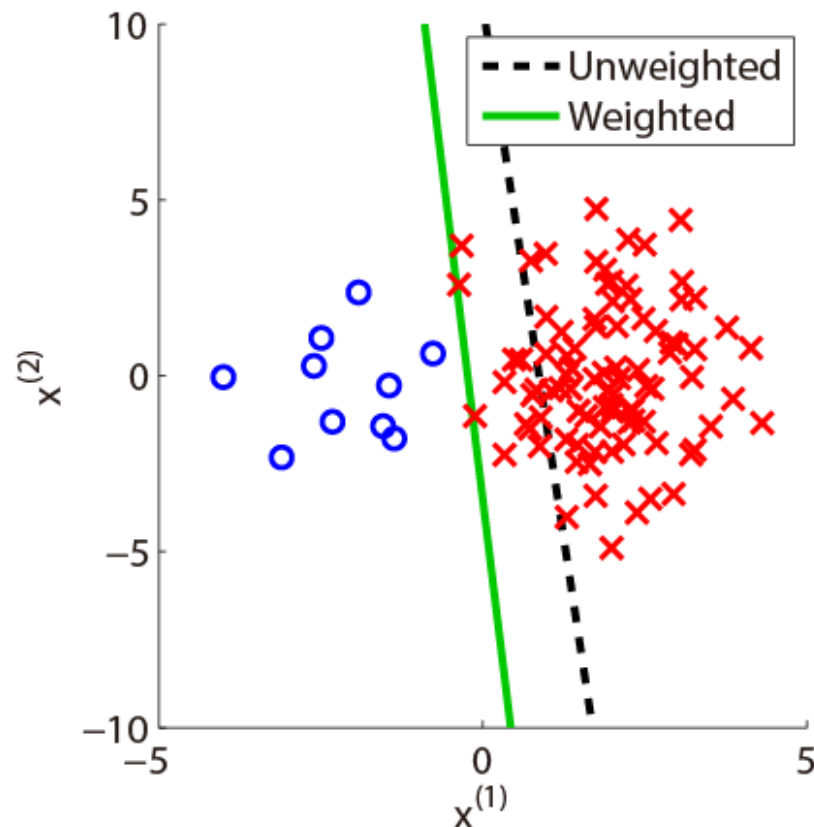
52

- 線形モデル $f_{\theta}(x) = \theta^{\top} x + \theta_0$ に対して, クラス比重み付き最小二乗法を実装せよ

訓練データ



テストデータ



宿題3(続き)

```
clear all; rand('state',0); randn('state',0);
x=[[randn(90,1)-2; randn(10,1)+2] 2*randn(100,1)];
y=[ones(90,1); 2*ones(10,1)];
X=[[randn(10,1)-2; randn(90,1)+2] 2*randn(100,1)];
% x, y, X: Training input, training output, test input

t=CWLS(x,y,X); % Parameter learning

Y=[ones(10,1); 2*ones(90,1)]; % Training output
figure(1); clf; hold on
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('Weighted');
axis([-5 5 -10 10])
```

CWLS.mを
実装せよ