

機械学習特論

～ 理論とアルゴリズム ～

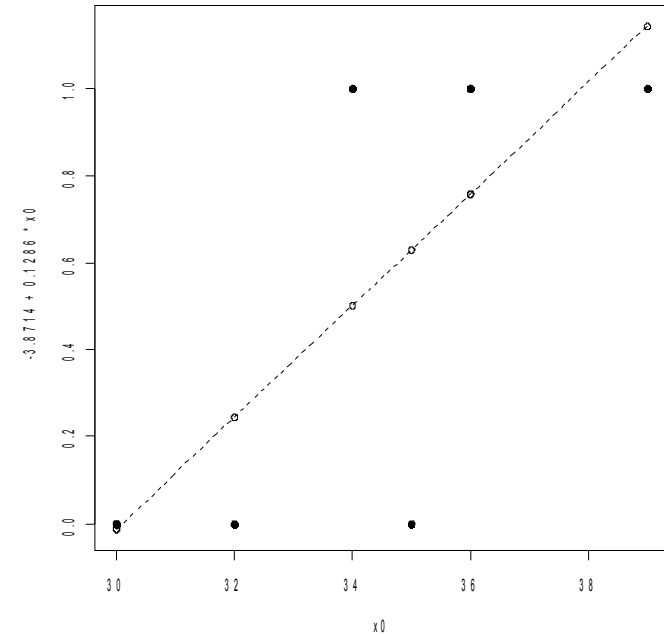
(Logistic Regression)

講師：西郷浩人

Ordinary Least Squares Regression vs Logistic Regression (LR)

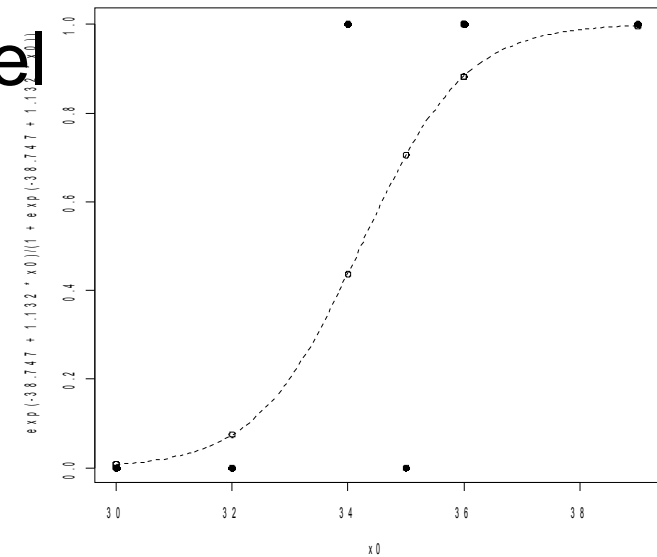
- OLS
 - a.k.a. Linear model

$$f(\mathbf{x}) = \beta_0 + \sum_{j=1}^p \beta_j x_j$$



- Logistic Regression
 - a.k.a. Generalized linear model

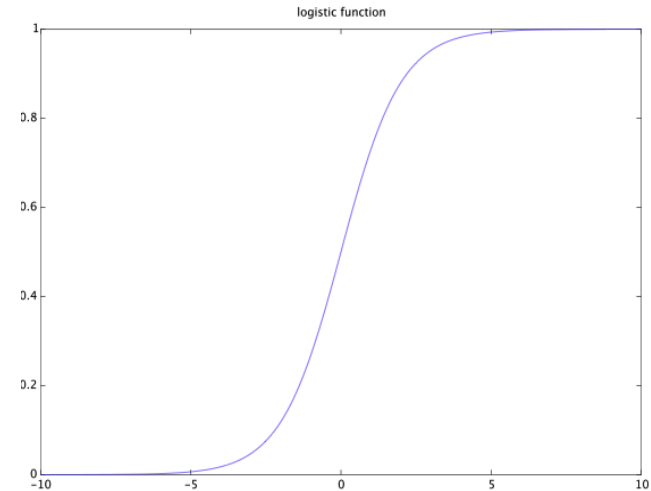
$$f(\mathbf{x}) = \frac{\exp(\beta_0 + \sum_{j=1}^p \beta_j x_j)}{1 + \exp(\beta_0 + \sum_{j=1}^p \beta_j x_j)}$$



Logistic Regression (LR)

- Logistic function maps linear function into a logistic form

$$h(\mathbf{x}) = \frac{\exp(-\mathbf{x})}{1 + \exp(-\mathbf{x})}$$



- Posterior probabilities

$$p(y_i=0|\mathbf{x}) = h(\boldsymbol{\beta}'\mathbf{x}) = \frac{\exp(-\beta_0 - \sum_{j=1}^p \beta_j x_j)}{1 + \exp(-\beta_0 - \sum_{j=1}^p \beta_j x_j)}$$

$$p(y_i=1|\mathbf{x}) = 1 - h(\boldsymbol{\beta}'\mathbf{x}) = 1 - \frac{\exp(-\beta_0 - \sum_{j=1}^p \beta_j x_j)}{1 + \exp(-\beta_0 - \sum_{j=1}^p \beta_j x_j)} = \frac{1}{1 + \exp(-\beta_0 - \sum_{j=1}^p \beta_j x_j)}$$

Logistic Regression (LR)

$$p(y_i=0|\mathbf{x})=h(\boldsymbol{\beta}'\mathbf{x})=\frac{\exp(-\beta_0-\sum_{j=1}^p\beta_jx_j)}{1+\exp(-\beta_0-\sum_{j=1}^p\beta_jx_j)}$$

$$p(y_i=1|\mathbf{x})=1-h(\boldsymbol{\beta}'\mathbf{x})=1-\frac{\exp(-\beta_0-\sum_{j=1}^p\beta_jx_j)}{1+\exp(-\beta_0-\sum_{j=1}^p\beta_jx_j)}=\frac{1}{1+\exp(-\beta_0-\sum_{j=1}^p\beta_jx_j)}$$

- Decision function is built by taking the ratio

$$f(\mathbf{x})=\frac{p(y_i=1|\mathbf{x})}{p(y_i=0|\mathbf{x})}=\exp(\beta_0+\sum_{j=1}^p\beta_jx_j)$$

- which has a linear function in an exponential form. (so-called generalized linear form)
- With this, we can perform 2-class classification by either if $f(\mathbf{x}) > 1$ or $f(\mathbf{x}) < 1$.

Solving logistic regression

$$p(y_i=0|\mathbf{x}) = h(\boldsymbol{\beta}'\mathbf{x}) = \frac{\exp(-\sum_{j=1}^p \beta_j x_j)}{1 + \exp(-\sum_{j=1}^p \beta_j x_j)}$$

$$p(y_i=1|\mathbf{x}) = 1 - h(\boldsymbol{\beta}'\mathbf{x}) = \frac{1}{1 + \exp(-\sum_{j=1}^p \beta_j x_j)}$$

Forget β_0
for a moment

- Define likelihood function as

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n p(y_i|\mathbf{x}_i)$$

Below we denote by \mathbf{x}_i
a column vector containing
features of the i-th example.

- Consider maximizing log-likelihood $l(\theta)$

$$l(\boldsymbol{\beta}) = \log L(\boldsymbol{\beta})$$

$$= \sum_{i=1}^n \{y_i \log p(y_i=1|\mathbf{x}_i) + (1-y_i) \log p(y_i=0|\mathbf{x}_i)\} \quad \leftarrow y_i=0,1$$

$$= \sum_{i=1}^n \{y_i \log h(\boldsymbol{\beta}'\mathbf{x}_i) + (1-y_i) \log (1-h(\boldsymbol{\beta}'\mathbf{x}_i))\}$$

$$= \sum_{i=1}^n \{y_i (\log h(\boldsymbol{\beta}'\mathbf{x}_i) - \log (1-h(\boldsymbol{\beta}'\mathbf{x}_i))) + \log (1-h(\boldsymbol{\beta}'\mathbf{x}_i))\}$$

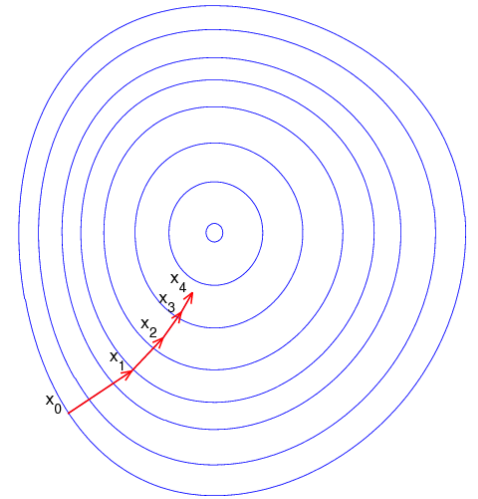
$$= \sum_{i=1}^n \{y_i \boldsymbol{\beta}'\mathbf{x}_i - \log (1 + \exp(\boldsymbol{\beta}'\mathbf{x}_i))\}$$

$$\log \left(\frac{h(\boldsymbol{\beta}'\mathbf{x})}{1-h(\boldsymbol{\beta}'\mathbf{x})} \right) = \boldsymbol{\beta}'\mathbf{x}_i$$

$$1-h(\boldsymbol{\beta}'\mathbf{x}) = \frac{1}{1 + \exp(-\sum_{j=1}^p \beta_j x_j)}$$

Searching for a minimum/maximum of a function: Gradient Descent

$$x \leftarrow x - \epsilon \frac{df(x)}{dx}$$



- A simple idea is to use the **first derivative** of the objective function.
- Step size ϵ needs to be determined via line search.
- Slow convergence.

Deriving the first derivative

$$l(\beta) = \sum_{i=1}^n \{y_i \beta' x_i - \log(1 + \exp(\beta' x_i))\}$$

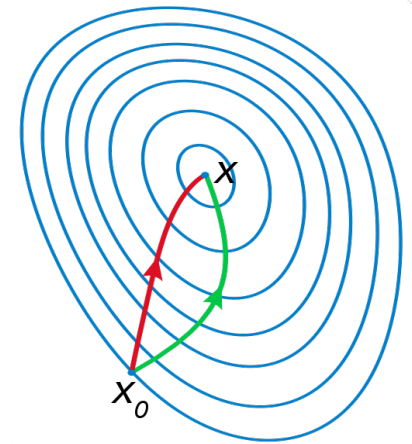
- The above objective function is convex, and the first derivative w.r.t. β is

$$\begin{aligned} \frac{\partial l}{\partial \beta} &= \frac{\partial}{\partial \beta} \sum_{i=1}^n \{y_i \beta' x_i - \log(1 + \exp(\beta' x_i))\} \\ &= \sum_{i=1}^n x_i (y_i - h(\beta' x_i)) \\ &= X' (y - h(X\beta)) \end{aligned}$$

- which again is a function of β . So we cannot solve it directly in an analytical form.
- Optimization such as Newton's method or gradient descent is required.

Searching for a minimum/maximum of a function: Newton's method

$$x \leftarrow x - \frac{\frac{df(x)}{dx}}{\frac{d^2f(x)}{dx^2}}$$



- **Both the first derivative and the second derivative** are used for searching the direction.
- Faster convergence.

Deriving the second derivative

$$\frac{\partial l}{\partial \boldsymbol{\beta}} = \sum_{i=1}^n \mathbf{x}_i (y_i - h(\boldsymbol{\beta}' \mathbf{x}_i))$$

- In order to use newton's method, we derive the second derivative from the first derivative.

$$\begin{aligned} \frac{\partial^2 l}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}} &= \frac{\partial}{\partial \boldsymbol{\beta}} \left\{ \sum_{i=1}^n -\mathbf{x}_i \left(y_i - \frac{\exp(\boldsymbol{\beta}' \mathbf{x}_i)}{1 + \exp(\boldsymbol{\beta}' \mathbf{x}_i)} \right) \right\} \\ &= \sum_{i=1}^n \left(\frac{-\mathbf{x}_i' \mathbf{x}_i \exp(\boldsymbol{\beta}' \mathbf{x}_i)}{(1 + \exp(\boldsymbol{\beta}' \mathbf{x}_i))^2} \right) \end{aligned}$$

- A matrix storing the second derivatives is called *Hessian*.

$$H = \sum_{i=1}^n h(\mathbf{x}_i' \boldsymbol{\beta}) (1 - h(\mathbf{x}_i' \boldsymbol{\beta})) \mathbf{x}_i \mathbf{x}_i' = \mathbf{X}' \mathbf{W} \mathbf{X}$$

– where \mathbf{W} is a diagonal matrix with entries

$$h(\mathbf{x}_i' \boldsymbol{\beta}) (1 - h(\mathbf{x}_i' \boldsymbol{\beta})) = \frac{\exp(\boldsymbol{\beta}' \mathbf{x}_i)}{(1 + \exp(\boldsymbol{\beta}' \mathbf{x}_i))^2}$$

Newton's method in our case

$$\beta \leftarrow \beta - H^{-1} \nabla_{\beta} J$$

- Let our objective function to minimize be $J = -\log L$, then

$$\begin{aligned} J(\beta) &= -\sum_{i=1}^n \{y_i \log h(\beta' x_i) + (1 - y_i) \log (1 - h(\beta' x_i))\} \\ &= -(\mathbf{y}' \log h(\mathbf{X} \beta) + (\mathbf{1} - \mathbf{y})' (\mathbf{1} - \log h(\mathbf{X} \beta))) \\ &= -\sum_{i=1}^n \{y_i \beta' x_i - \log (1 + \exp(\beta' x_i))\} \end{aligned}$$

- The first derivative is a vector:

$$\nabla_{\beta} J = \frac{-\partial l}{\partial \beta} = \mathbf{X}' (h(\mathbf{X} \beta) - \mathbf{y})$$

- The second derivative is a matrix

$$\mathbf{H} = \sum_{i=1}^n h(x_i' \beta) (1 - h(x_i' \beta)) x_i x_i' = \mathbf{X}' \mathbf{W} \mathbf{X}$$

where \mathbf{W} is a diagonal matrix with entries

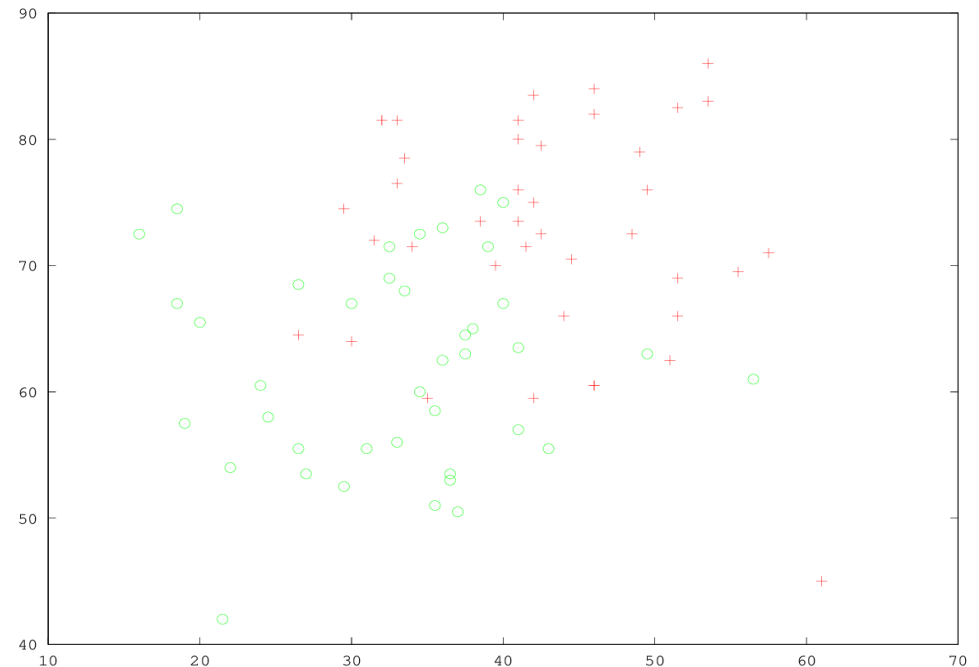
$$h(x_i' \beta) (1 - h(x_i' \beta))$$

Algorithm1: Logistic Regression

- Input:
 - X : $n \times p$ data matrix
 - y : $n \times 1$ binary response vector
- Output
 - β : $p \times 1$ coefficient vector
- Initialize
 - $J=0$; $\text{prev_J} = \text{LARGE_NUMBER}$;
- Repeat
 - Compute J
$$J(\beta) = \sum_{i=1}^n \{y_i \log p(y_i=1|\mathbf{x}_i) + (1-y_i) \log p(y_i=0|\mathbf{x}_i)\}$$
 - If $\text{prev_J} - J < \text{SMALL_NUMBER}$
 - Break
 - Compute nabla_J
$$\nabla_{\beta} J = X'(h(X\beta) - y)$$
 - Compute H
$$W = \text{diag}(h(X\beta) \cdot (1 - h(X\beta))) \quad H = X' W X$$
 - Update β
$$\beta \leftarrow \beta - H^{-1} \nabla_{\beta} J$$

Exercise 1: Implement Logistic Regression by Newton's method

- Download today's data
- Executing init.m loads and plots data
- Complete logistic.m
- Logistic sigmoid function $h(x)$ is already prepared.



Hint 1

- You can build a diagonal matrix with entries 1,2,3 by

```
diag([1,2,3])
```

- If correctly implemented, Newton's method will converge in 5 iterations with the following answer.
 - $\text{beta} = [-16.37874 \ 0.14834 \ 0.15891]$
 - $J = 32.436$

Hint 2

$$\begin{aligned} J(\boldsymbol{\beta}) &= -(\mathbf{y}' \log h(\mathbf{X}\boldsymbol{\beta}) + (\mathbf{1} - \mathbf{y})' \log (\mathbf{1} - h(\mathbf{X}\boldsymbol{\beta}))) \\ &= -(\mathbf{y}' \mathbf{X}\boldsymbol{\beta} - \sum (\log(1 + \exp(\mathbf{X}\boldsymbol{\beta})))) \end{aligned}$$

$$\begin{aligned} J &= -(y' * \log(h(X*b)) + (ones(n,1) - y)' * \log(ones(n,1) - h(X*b))) \\ &= -(y' * X*b - \text{sum}(\log(1 + \exp(X*b)))) \end{aligned}$$

Exercise 1 ~continued~

- At decision boundary, following relationship holds.

$$\begin{aligned}h(\beta'x) &= 1 - h(\beta'x) \\ \Leftrightarrow \frac{1}{1 + \exp(\beta'x)} &= \frac{\exp(\beta'x)}{1 + \exp(\beta'x)} \\ \Leftrightarrow \exp(\beta'x) &= 1 \\ \Leftrightarrow \beta_0 + \beta_1 x_1 + \beta_2 x_2 &= 0\end{aligned}$$

- Plot the decision boundary.
- Testing can be done with $h(X^*b)$, which returns a probability of one class.

$$\begin{aligned}p &= h(X^*b) \\ q &= 1 - h(X^*b)\end{aligned}$$

Exercise 1 ~continued2~

- Using p, you can plot decision boundary in 3D !

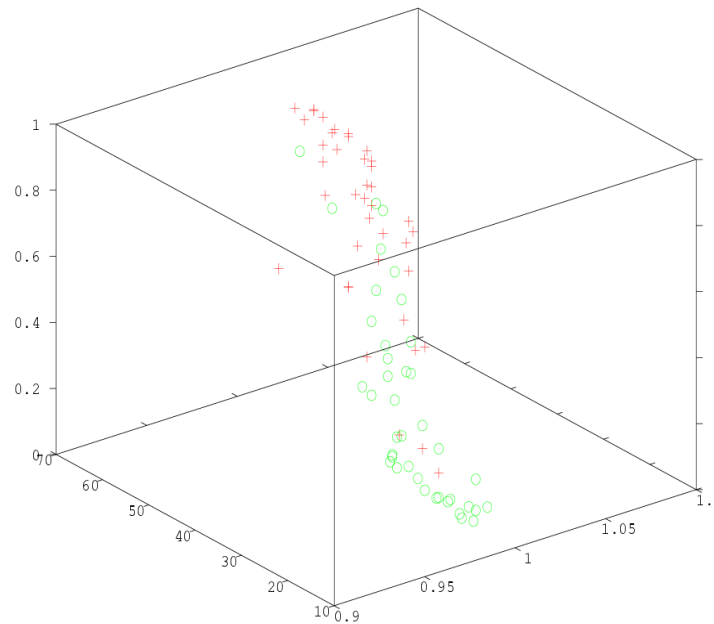
```
figure
```

```
plot3(X(pos,2),X(pos,3),p(pos),'go')
```

```
hold on
```

```
plot3(X(neg,2),X(neg,3),p(neg),'rx')
```

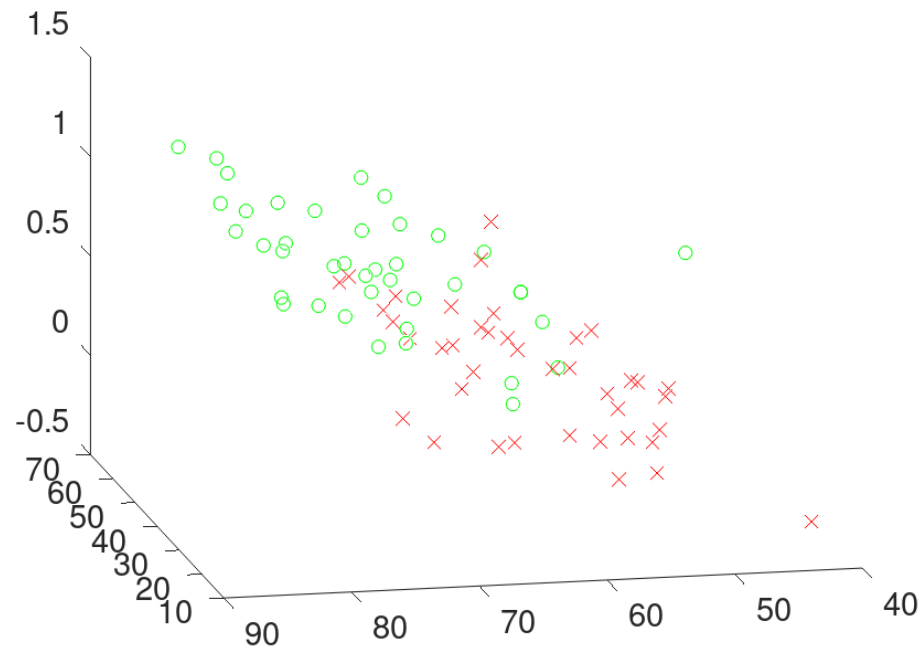
```
hold off
```



Exercise 1 ~continued3~

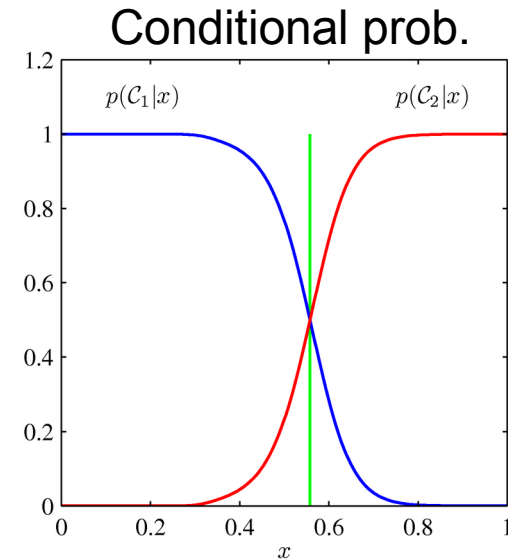
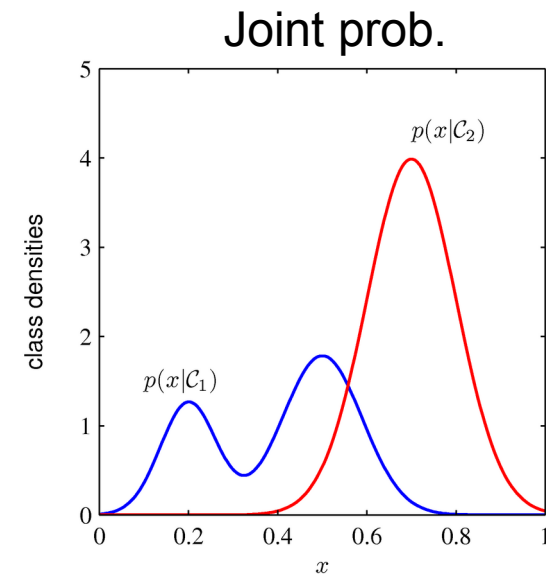
- Let's compare the results with that of OLS.

```
b_ols = X \ y
y_ols = X * b_ols
figure
plot3(X(pos,2),X(pos,3),y_ols(pos),'go')
hold on
plot3(X(neg,2),X(neg,3),y_ols(neg),'rx')
```



Generative vs Discriminative

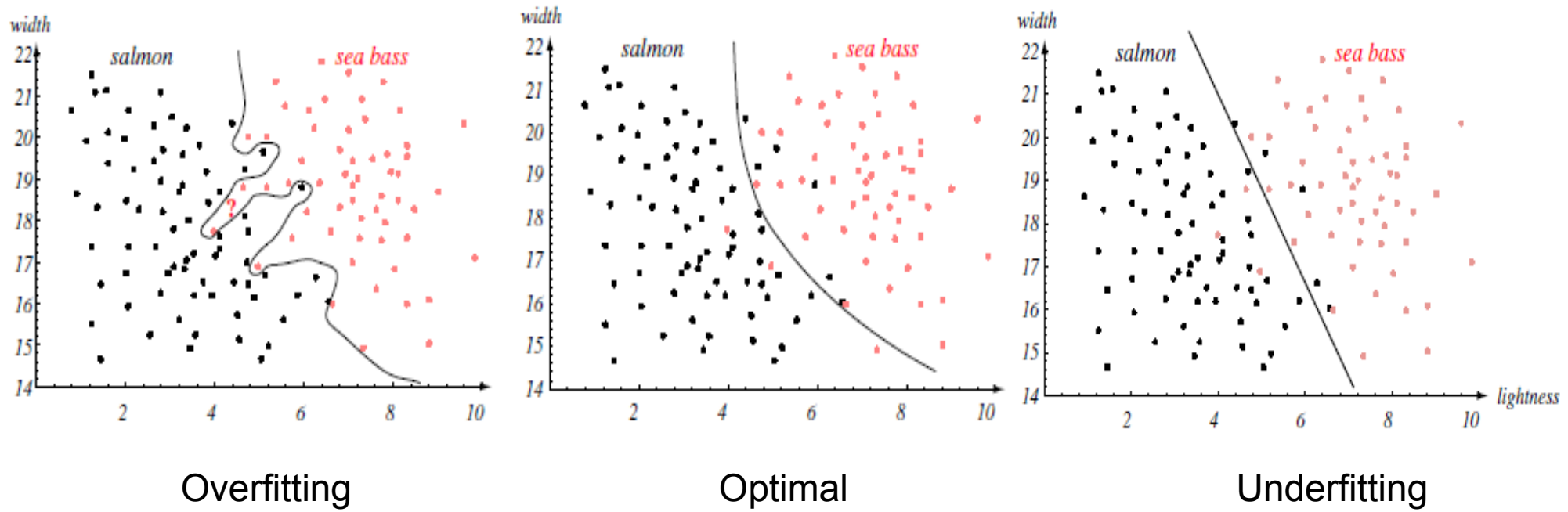
- The goal is to estimate $p(y|x)$
- generative approach(生成的手法)
first estimates joint probability $p(x, y)$
- then posterior
$$p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x, y)}{p(x)}$$
 - ex. LDA, QDA
- discriminative approach(識別的手法)
directly estimates $p(y|x)$
 - ex. Logistic Regression



LDA vs Logistic Regression(LR)

- LDA assumes that marginal probability $p(x,y)$ follows gaussian distribution, while LR does not have any assumption.
 - Since $p(y|x)$ is modeled directly, we do not need to model $p(x,y)$ in Logistic regression
- In situations where $p(x,y)$ follows gaussian distribution, LDA would work well, but in other cases, LR would perform better.
- In terms of computational complexity, LDA requires inversing a covariance once, while LR require it in every iteration, therefore much more expensive.

Regularization for classification



Regularization for regression

$$RSS = \|X\beta - y\|_2^2$$

- OLS

$$\min_{\beta} RSS$$

- Ridge regression (L2-regularization)

$$\min_{\beta} RSS + \lambda \|\beta\|_2$$

- LASSO regression (L1-regularization)

$$\min_{\beta} RSS + \lambda \|\beta\|_1$$

Regularization for classification

$$J = -l(\boldsymbol{\beta}) = -\sum_{i=1}^n \{y_i h(\mathbf{x}_i) + (1 - y_i) h(\mathbf{1} - \mathbf{x}_i)\}$$

- Logistic regression

$$\min_{\boldsymbol{\beta}} J$$

- L2-regularized Logistic regression

$$\min_{\boldsymbol{\beta}} J + \lambda \|\boldsymbol{\beta}\|_2$$

- L1-regularized Logistic regression

$$\min_{\boldsymbol{\beta}} J + \lambda \|\boldsymbol{\beta}\|_1$$

Solving L2-Logistic Regression

$$\min_{\boldsymbol{\beta}} J + \lambda \|\boldsymbol{\beta}\|_2$$

- Objective function

$$J = -l(\boldsymbol{\beta}) = -\sum_{i=1}^n \{y_i \boldsymbol{\beta}' \mathbf{x}_i - \log(1 + \exp(\boldsymbol{\beta}' \mathbf{x}_i))\} + \frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$

Deriving the first derivative

$$J = -l(\beta) = -\sum_{i=1}^n \{y_i \beta' \mathbf{x}_i - \log(1 + \exp(\beta' \mathbf{x}_i))\} + \frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$

- Same way as the ordinary logistic regression.

$$\begin{aligned} \nabla J &= \frac{-\partial l}{\partial \beta_j} = \frac{-\partial}{\partial \beta_j} \left(\sum_{i=1}^n \{y_i \beta' \mathbf{x}_i - \log(1 + \exp(\beta' \mathbf{x}_i))\} + \frac{\lambda}{2} \sum_{j=1}^p \beta_j^2 \right) \\ &= -\sum_{i=1}^n \mathbf{x}_i (y_i - h(\beta' \mathbf{x}_i)) + \lambda \beta_j \end{aligned}$$

- Vector form $\nabla_{\beta} J = \mathbf{X}'(h(\mathbf{X}\beta) - \mathbf{y}) + \lambda \beta$

Deriving the second derivative

$$\nabla J = \frac{-\partial l}{\partial \beta} = -X'(\mathbf{y} - h(X\beta)) + \lambda \beta$$

- Same way as the previous lecture.

$$\begin{aligned} \frac{\partial J}{\partial \beta^2} &= \frac{-\partial l^2}{\partial \beta \partial \beta} = \frac{-\partial}{\partial \beta} \left\{ \sum_{i=1}^n -\mathbf{x}_i \left(1 - \frac{\exp(\beta' \mathbf{x}_i)}{1 + \exp(\beta' \mathbf{x}_i)} \right) + \lambda \beta_j \right\} \\ &= -\sum_{i=1}^n \left(\frac{-\mathbf{x}_i' \mathbf{x}_i \exp(\beta' \mathbf{x}_i)}{(1 + \exp(\beta' \mathbf{x}_i))^2} \right) + \lambda \end{aligned}$$

$$H = \sum_{i=1}^n h(\mathbf{x}_i' \beta)(1 - h(\mathbf{x}_i' \beta)) \mathbf{x}_i \mathbf{x}_i' + \lambda \begin{bmatrix} 1 & 0 & 0 & \cdot \\ 0 & 1 & 0 & \cdot \\ 0 & 0 & 1 & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix} = X' W X + \lambda I$$

$$W = \text{diag}(h(X\beta)(1 - h(X\beta))) = \begin{bmatrix} h(\mathbf{x}_1' \beta)(1 - h(\mathbf{x}_1' \beta)) & 0 & 0 & \cdot \\ 0 & h(\mathbf{x}_2' \beta)(1 - h(\mathbf{x}_2' \beta)) & 0 & \cdot \\ 0 & 0 & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{bmatrix}$$

Algorithm: L2-Logistic Regression

- Input:
 - X : $n \times p$ data matrix
 - y : $n \times 1$ binary response vector
 - λ
- Output
 - β : $p \times 1$ coefficient vector
- Initialize
 - $J=0$; $\text{prev_J} = \text{LARGE_NUMBER}$;
- Repeat
 - Compute J $J(\beta) = \sum_{i=1}^n \{y_i \log p(y_i=1|\mathbf{x}_i) + (1-y_i) \log p(y_i=0|\mathbf{x}_i)\} + \frac{\lambda}{2} \|\beta\|_2^2$
 - If $\text{prev_J} - J < \text{SMALL_NUMBER}$
 - Break
 - Compute nabla_J $\nabla_{\beta} J = X' (h(X\beta) - y) + \lambda \beta$
 - Compute H $W = \text{diag}(h(X\beta) \cdot (1 - h(X\beta)))$ $H = X' W X + \lambda I$
 - Update β $\beta \leftarrow \beta - H^{-1} \nabla_{\beta} J$

Exercise 2

- Implement L2-logistic regression of the form
 - `[b] = function logistic2(x,y,lambda)`
- First verify that setting `lambda = 0` obtains the same solution as before.

- Observe

- How norm of `b` changes

`norm(b)`

- How accuracy changes

`sum((h(X*b)>0.5)==y)/length(y)`

with respect to change in `lambda`.