

機械学習特論

～ 理論とアルゴリズム ～

第 9 回

(Regularization for Logistic Regression)

講師：西郷浩人

Previous exercise 1: Implement Logistic Regression by Newton's method

- Download today's data
- Executing ex1.m loads and plots data
- Complete logistic.m
- Logistic function $h(x)$ is already declared.

logistic.m

```
function [b,J]=logistic(X,y)
[n p]=size(X); b = zeros(p);
J=0; itr=0;

while 1
    itr = itr + 1;
    disp(["iteration: ",num2str(itr)]);
    prev_J = J;

    % Compute J here
    % J = -(y'*log(h(X*b))+(ones(n,1)-y)'*log(ones(n,1)-h(X*b)))
    J = -(y'*X*b-sum(log(1+exp(X*b))))
    if abs(prev_J-J) < 1/n, break, end

    % Compute the first gradient here
    grad = X'*(h(X*b)-y);

    % Compute the second gradient here
    W=diag(h(X*b).*(ones(n,1)-h(X*b)));
    H=X'*W*X;

    % Update beta here
    b = b - inv(H)*grad; % update weights
end
endfunction

function [y] = h(x)
y = 1.0 ./ (1.0 + exp(-x));
endfunction
```

$$J(\boldsymbol{\beta}) = -(\mathbf{y}' \log h(\mathbf{X}\boldsymbol{\beta}) + (\mathbf{1} - \mathbf{y})' \log (\mathbf{1} - h(\mathbf{X}\boldsymbol{\beta})))$$

$$= -(\mathbf{y}' \mathbf{X}\boldsymbol{\beta} - \sum (\log(1 + \exp(\mathbf{X}\boldsymbol{\beta}))))$$

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

$$\mathbf{H} = \mathbf{X}' \mathbf{W} \mathbf{X}$$

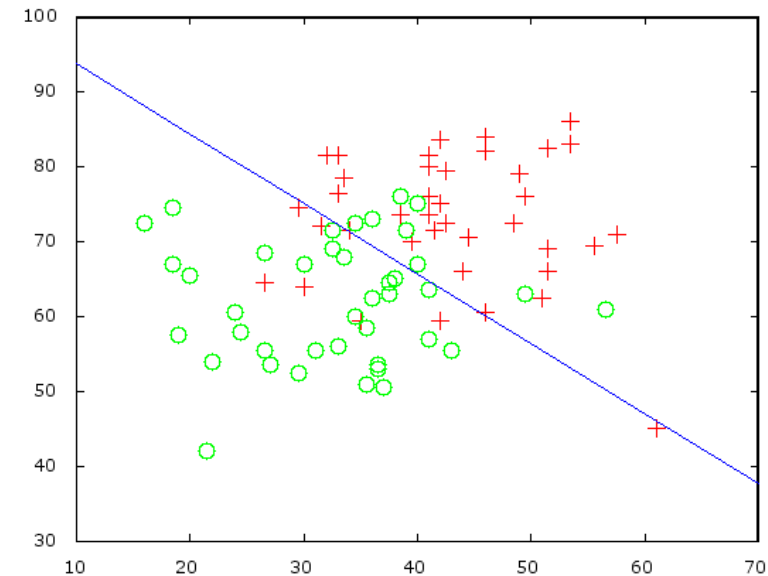
$$\boldsymbol{\beta} \leftarrow \boldsymbol{\beta} - \mathbf{H}^{-1} \nabla_{\boldsymbol{\beta}} J$$

Decision boundary on test data

- Coefficients and decision boundary

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

beta =
-16.37874
0.14834
0.15891
J = 32.436



- Equation for the decision boundary

$$\begin{aligned}-16.4 + 0.15x_1 + 0.16x_2 &= 0 \\ \longleftrightarrow x_2 &= \frac{16.4}{0.16} - \frac{0.15}{0.16}x_1\end{aligned}$$

Solving Logistic Regression by Iteratively Reweighted Least Squares (IRLS)

Algorithm: 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 $J(\beta)$
 - If $\text{prev_J} - J < \text{SMALL_NUMBER}$
 - Break
 - Compute nabla_J $\nabla_{\beta} J = X'(h(X\beta) - y)$
 - Compute H $W = \text{diag}(h(X\beta)(1 - h(X\beta)))$ $H = X'WX$
 - Update β $\beta \leftarrow \beta - H^{-1} \nabla_{\beta} J$

Drawbacks

- It requires to calculate the inverse of a Hessian matrix.
 - This step can be numerically unstable, and/or slow.
 - Therefore IRLS (Iteratively Re-weighted Least Squares) is often used instead for solving logistic regression.

Rewriting the update rule

$$\nabla_{\beta} J = X' (h(X\beta) - y) \quad H = X' W X$$

- Rewrite the update rule

$$\begin{aligned} \beta &\leftarrow \beta - H^{-1} \nabla_{\beta} J \\ &= \beta + (X' W X)^{-1} X' (y - h(X\beta)) \\ &= (X' W X)^{-1} X' W (X\beta + W^{-1} (y - h(X\beta))) \\ &= (X' W X)^{-1} X' W z \end{aligned}$$

- where $z = X\beta + W^{-1} (y - h(X\beta))$
- Remember that the OLS solution for $X\beta = y$ is $\beta = (X' X)^{-1} X' y$
- New update rule iteratively solves weighted least squares problem $(X' W X)\beta = X' W z$
 - Whose solution is $\beta = (X' W X)^{-1} X' W z$
 - However, z and W are functions of β , so update is repeated until convergence.

Algorithm: Logistic Regression by IRLS

- 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 $J(\beta)$
 - If $\text{prev_J} - J < \text{SMALL_NUMBER}$
 - Break
 - Compute w $W = \text{diag}(h(X\beta)(1 - h(X\beta)))$
 - Compute z $z = (X\beta + W^{-1}(y - h(X\beta)))$
 - Solve $(X'WX)\beta = X'Wz$ for β

Ex. 1: implement logistic regression by IRLS

- Implement logistic regression by IRLS and name it as `logisticIRLS.m`
- You can load data by typing "init"
- Remember to use "back slash operator" for solving normal equation.
- You can check time spent for calculation using "tic" and "toc".

```
>tic, [beta J]=logistic(X,y), toc
```

- Compare the time and the solution of `logistic.m` with that of `logisticIRLS.m`.

logisticIRLS.m

```
function [beta, J] = logisticIRLS(X,y)
[n p] = size(X);
beta = zeros(p,1);
itr = 0;
J = 0;
while 1
    itr = itr + 1;
    prev_J = J;

    % compute J here
    %

    if abs(prev_J-J) < 1/n break, end

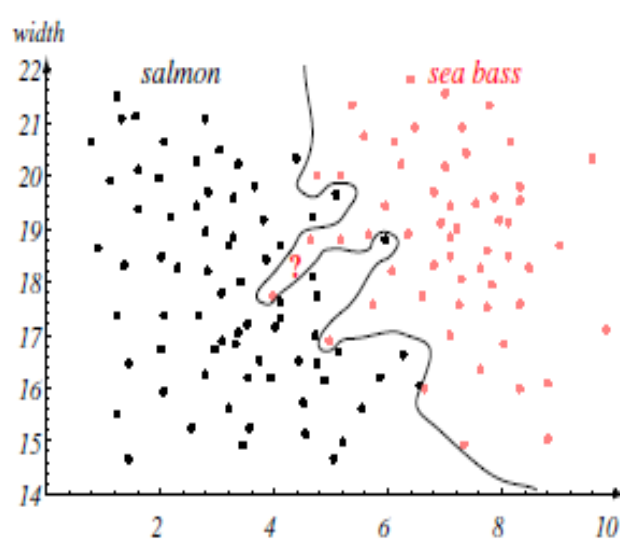
    % update beta here
    %

end
endfunction

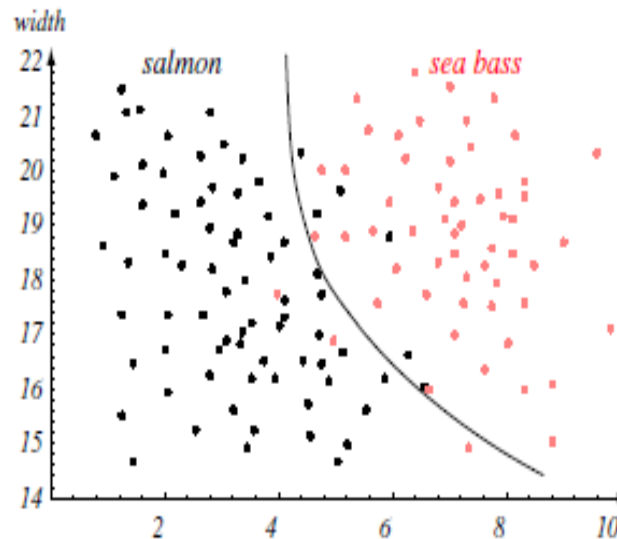
function [y] = h(x)
y = 1.0 ./ (1.0 + exp(-x));
endfunction
```

Regularization in classification

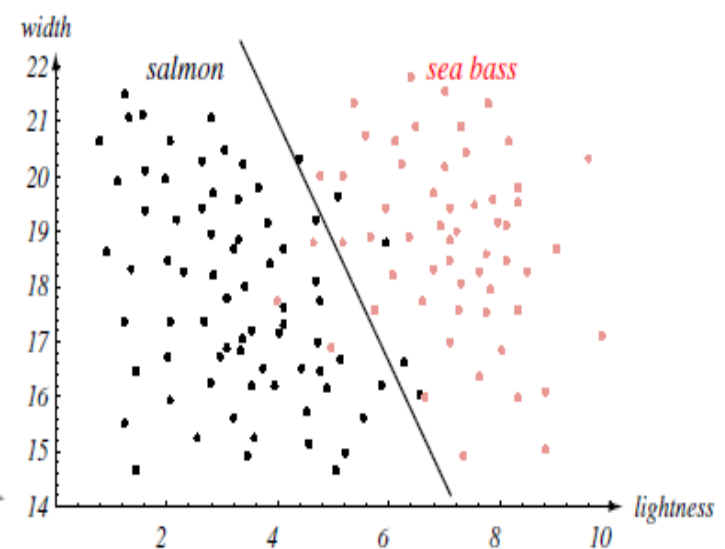
- An optimal classification rule is deemed to exist between overfitting and underfitting.
- How can we adjust between the two ?
 - > via regularization parameter



Overfitting



Optimal



Underfitting

Regularization for regression

$$RSS = \| \mathbf{X} \boldsymbol{\beta} - \mathbf{y} \|_2^2$$

- OLS

$$\min_{\boldsymbol{\beta}} RSS$$

- Ridge regression (L2-regularization)

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

- LASSO regression (L1-regularization)

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

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 by Iteratively Reweighted Ridge Regression (IRRR)

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

Solving L2-Logistic Regression

- 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

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

- Same way as the previous lecture.

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

- Vector form

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

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.

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

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

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

Rewriting the update rule

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

$$H = X' W X + \lambda I$$

- Rewrite the update rule

$$\begin{aligned}\beta &\leftarrow \beta - H^{-1} \nabla_{\beta} J \\ &= \beta - (X' W X + \lambda I)^{-1} (X' (h(X\beta) - y) + \lambda \beta) \\ &= (I - (X' W X + \lambda I)^{-1} \lambda) \beta + (X' W X + \lambda I)^{-1} X' (h(X\beta) - y) \\ &= (X' W X + \lambda I)^{-1} X' W X \beta + (X' W X + \lambda I)^{-1} X' (h(X\beta) - y) \\ &= (X' W X + \lambda I)^{-1} X' W (X\beta + W^{-1} (y - h(X\beta))) \\ &= (X' W X + \lambda I)^{-1} X' W z\end{aligned}$$

- where $z = X\beta + W^{-1} (y - h(X\beta))$
- New update rule iteratively solves weighted ridge regression problem $(X' W X + \lambda I) \beta = X' W z$

Algorithm: L2-Logistic Regression by IRLS

- 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
$$J(\beta) = -\sum_{i=1}^n \{y_i \beta' x_i - \log(1 + \exp(\beta' x_i))\} + \frac{\lambda}{2} \|\beta\|_2^2$$
- If $\text{prev_J} - J < \text{SMALL_NUMBER}$
 - Break
- Compute w
$$W = \text{diag}(h(X\beta)(1 - h(X\beta)))$$
- Compute z
$$z = (X\beta + W^{-1}(y - h(X\beta)))$$
- Solve $(X'WX + \lambda I)\beta = X'Wz$ for β

Exercise. 2: Implement L2-Logistic Regression by IRRR

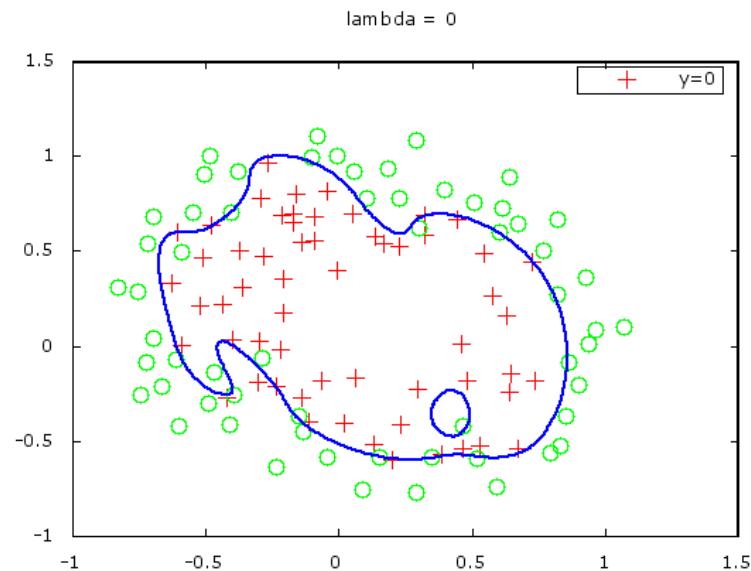
- Implement Algorithm 2, and name it as “logisticIRLS2.m”.
- After completion, you can plot the decision boundary by using “plotLR.m” function.

- Usage example when $\lambda=0$

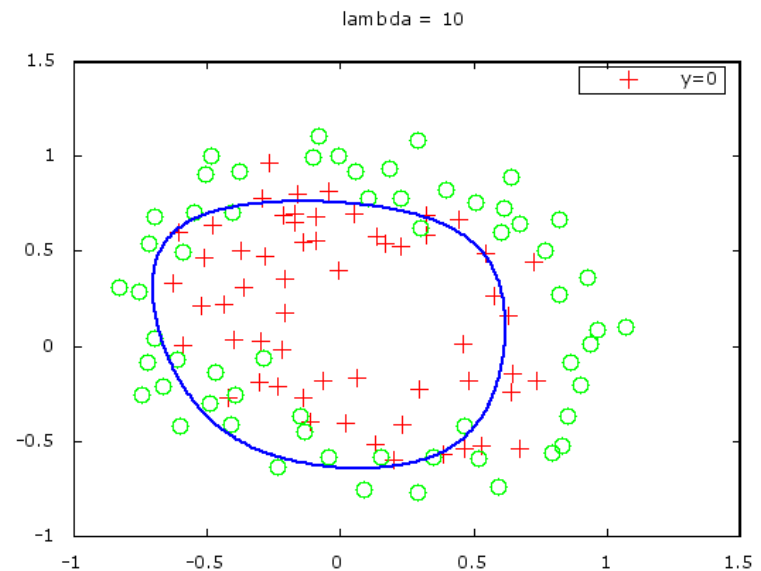
```
>plotLR(X,y,0)
```

- It uses fifth order polynomial function via map_feature.m
- Observe how decision boundary changes as you increase λ .

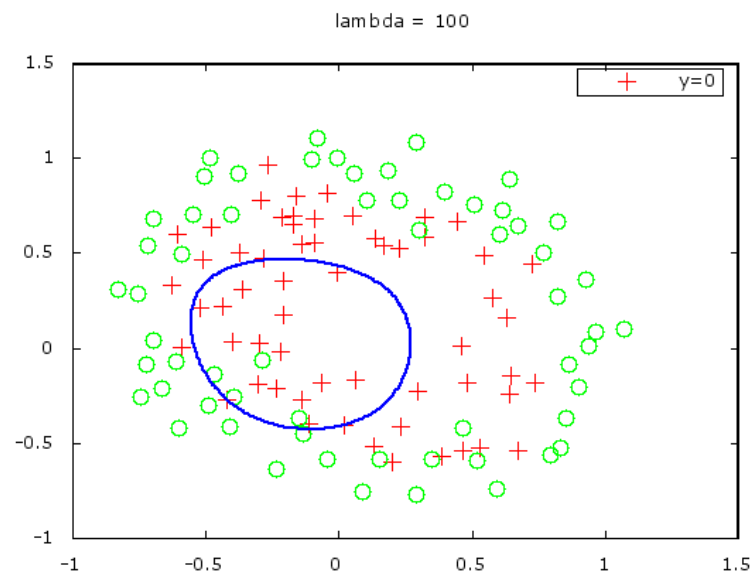
Demonstration: L2



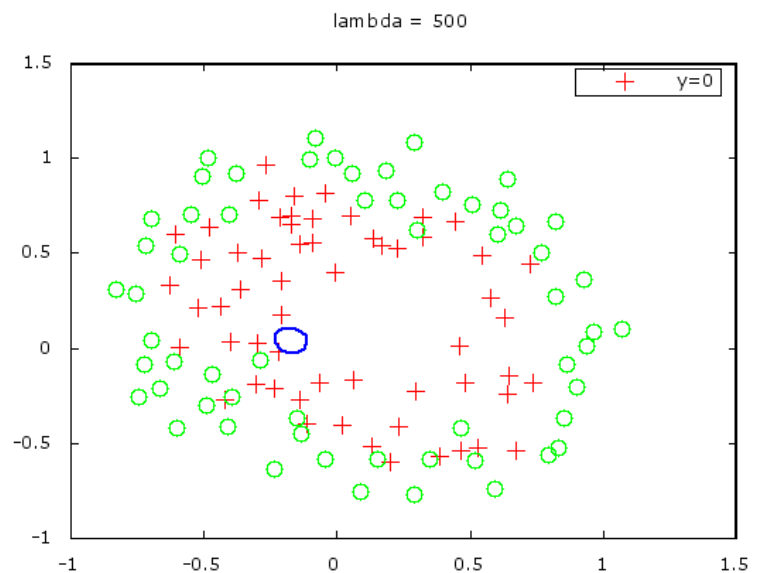
Sparsity: 0%, Likelihood: -23.4



Sparsity: 0%, Likelihood: -76.8



Sparsity: 4%, Likelihood: -85.5



Sparsity: 28.6%, Likelihood: -87.1

Solving L1-Logistic Regression by Iteratively Reweighted Ridge Regression (IRRR)

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

Solving L1-logistic regression

$$\min_{\beta} J + \lambda \|\beta\|_1 \quad \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$$

Remember that we have solved L1-regularized least squares problem (LASSO) by using the following update rule.

$$\beta^{lasso} \leftarrow (\mathbf{X}' \mathbf{X} + \lambda \mathbf{B})^{-1} \mathbf{X}' \mathbf{y}$$

$$\mathbf{B} = \text{diag}(1/|\beta_1|, 1/|\beta_2|, \dots, 1/|\beta_p|)$$

where

In logistic regression, L1-norm solution is easily obtained by solving the following normal equation in IRLS.

$$\beta^{L1LR} \leftarrow (\mathbf{X}' \mathbf{W} \mathbf{X} + \lambda \mathbf{B})^{-1} \mathbf{X}' \mathbf{W} \mathbf{z}$$

Algorithm: L1-Logistic Regression by IRLS

- 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
$$J(\beta) = -\sum_{i=1}^n \{y_i \beta' x_i - \log(1 + \exp(\beta' x_i))\} + \lambda \|\beta\|_1$$
 - If $\text{prev_J} - J < \text{SMALL_NUMBER}$
 - Break
 - Compute w
$$W = \text{diag}(h(X\beta)(1 - h(X\beta)))$$
 - Compute z
$$z = (X\beta + W^{-1}(y - h(X\beta)))$$
 - Solve $(X'WX + \lambda B)\beta = X'Wz$ for β

Hint

- You can build $\mathbf{B} = \text{diag}(1/|\beta_1|, 1/|\beta_2|, \dots, 1/|\beta_p|)$ by

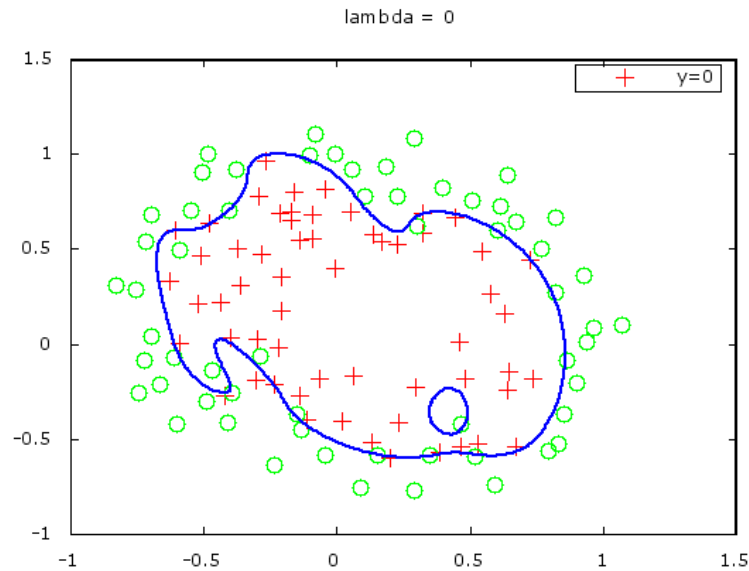
```
pinv(diag(abs(beta)))
```

- where pinv means "psuedo inverse", which is useful even when inverse of a matrix cannot be defined because of singularity.

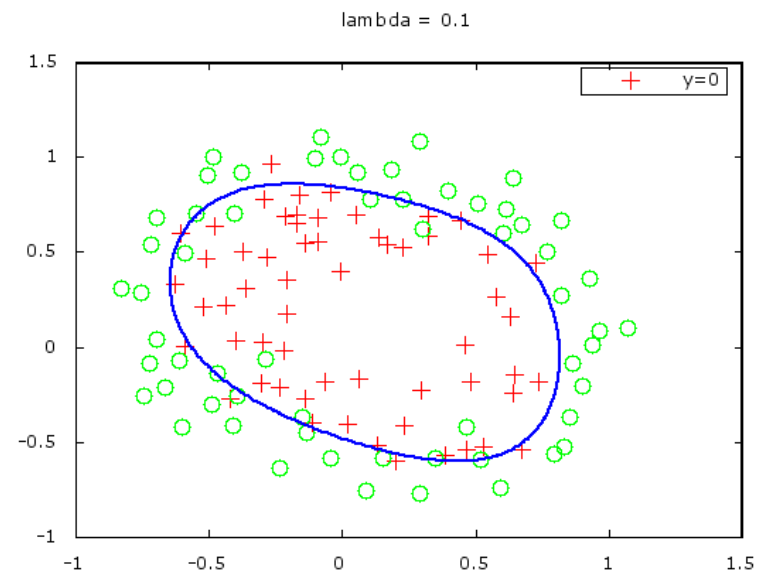
Ex. 3: Implement L1-Logistic Regression by IRLS

- Modify two lines in “logisticIRLSL2.m”.
 - Objective function
 - Least squares fit
- Compare the obtained likelihood and sparsity with those obtained by L2-Logistic Regression.
- You can plot data with
 - `plotLR(X,y,0.3,1)`, where the last option specifies norm, and takes either '1' or '2'.

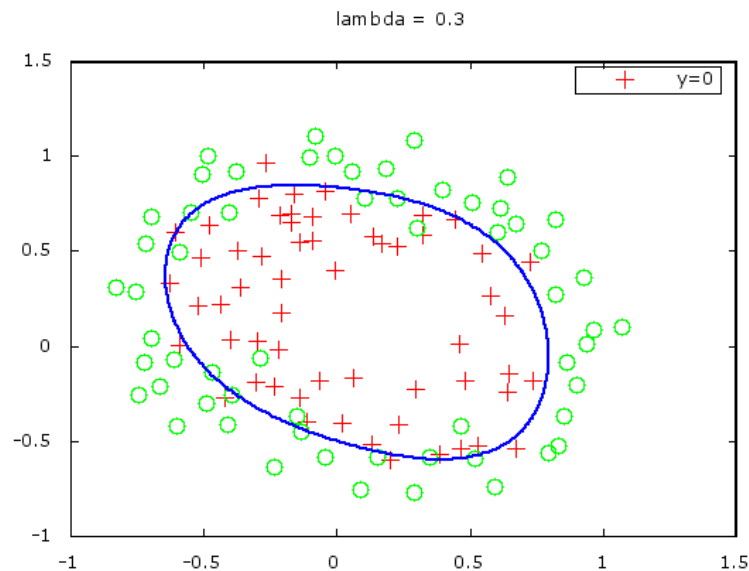
Demonstration: L1



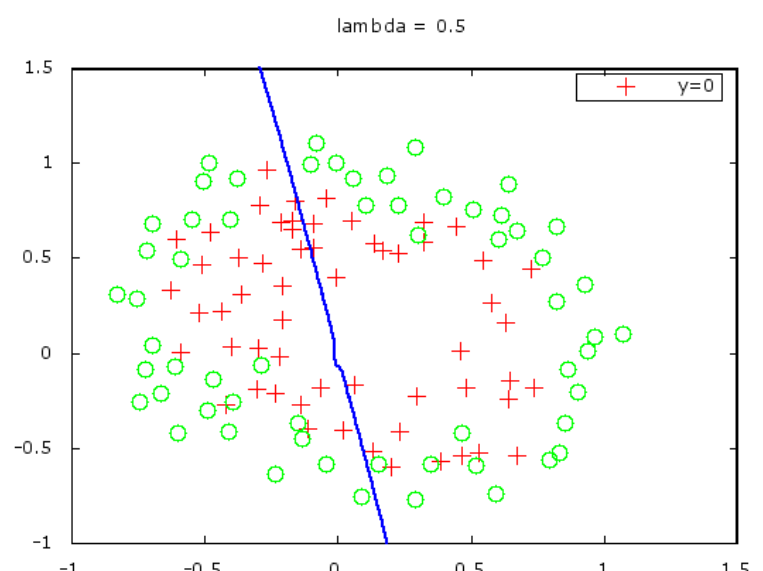
Sparsity: 0%, Likelihood: -0.2



Sparsity: 32%, Likelihood: -88



Sparsity: 60%, Likelihood: -73



Sparsity: 100%, Likelihood: -81

Summary

- Iteratively reweighted least squares (IRLS) can solve logistic regression in a robust and efficient manner.
- OLS
 - L1-norm (Ridge)
 - L2-norm (LASSO)
- Logistic Regression
 - L1-norm
 - L2-norm