

MACHINE LEARNING

Electrical Summer Workshops (ESW) 2022

Electrical Engineering Department - Sharif University of Technology

Instructors: *Alireza Gargoori Motlagh – Amir Mirrashid – Ali Nourian*



LOGISTIC REGRESSION

Linear Regression as Classifier

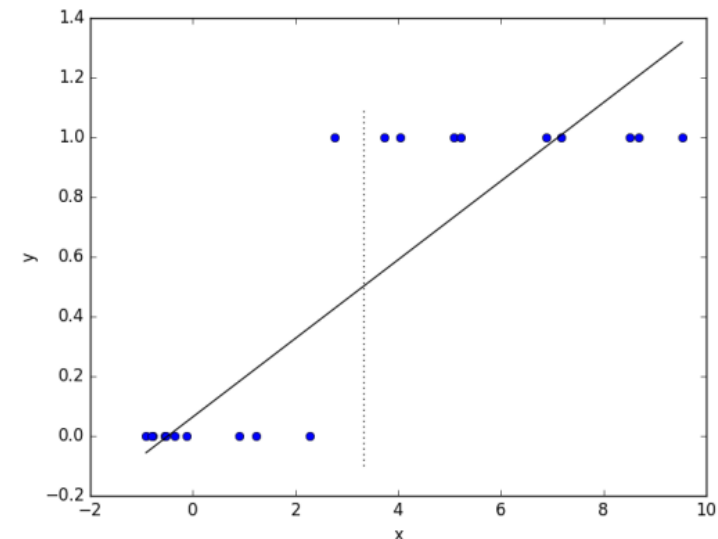
A basic approach to solve classification problems could use linear regression as follows: ($y^{(i)} \in \{0,1\}$)

$$RSS = \sum_{i=1}^n (y^{(i)} - \boldsymbol{\beta}^T \mathbf{x}^{(i)})^2$$

$$\rightarrow y^{(i)} = \begin{cases} 1 & \boldsymbol{\beta}^T \mathbf{x}^{(i)} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

- Arbitrary threshold
- Optimal coefficients try to make data points hug the line as closely as possible, no matter the effect on classification performance

- Our loss function needs to be modified.



Logit Transform

Predicting probabilities : $\mathbb{R}^n \rightarrow (0,1)$

- Dot product to \mathbb{R}_+^n through exponentiation
- Divide the result to 1 + *itself*

→ *sigmoid* function:

$$\sigma(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$$

where: $z = \boldsymbol{\beta}^T \mathbf{x}$ and $\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{bmatrix}$

$$\mathbb{P}(Y = 1|\mathbf{X}) = \frac{1}{1 + e^{-\boldsymbol{\beta}^T \mathbf{X}}}$$

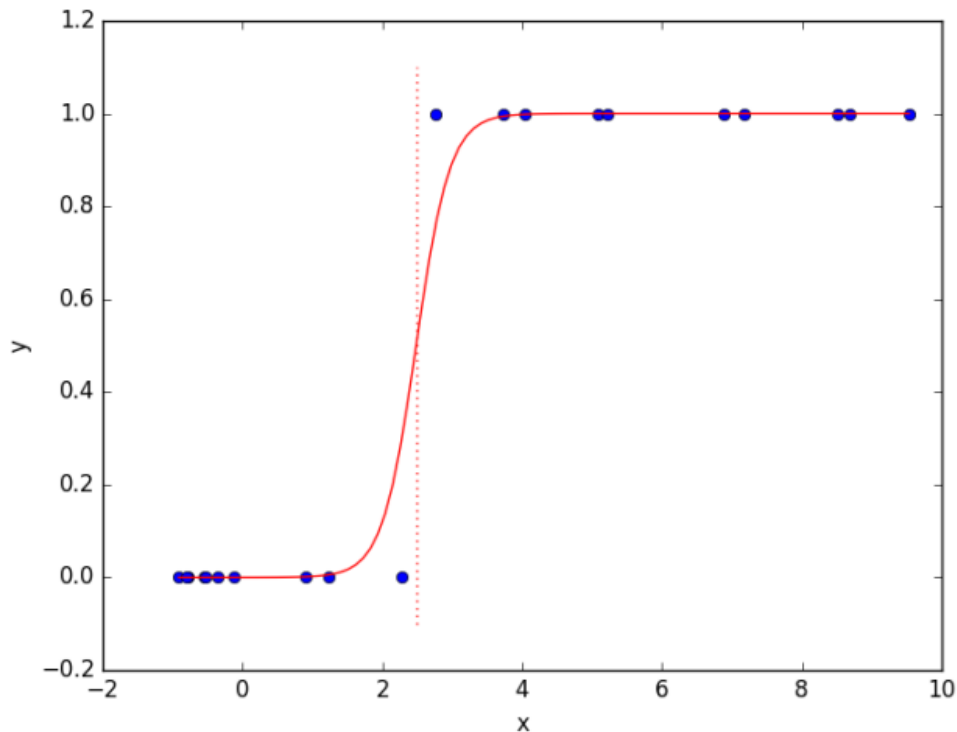
$$\mathbb{P}(Y = 0|\mathbf{X}) = \frac{1}{1 + e^{\boldsymbol{\beta}^T \mathbf{X}}}$$

$$\xrightarrow{\log\text{-odds}} \log\left(\frac{\mathbb{P}(Y = 1|\mathbf{X})}{\mathbb{P}(Y = 0|\mathbf{X})}\right) = \boldsymbol{\beta}^T \mathbf{X} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p$$

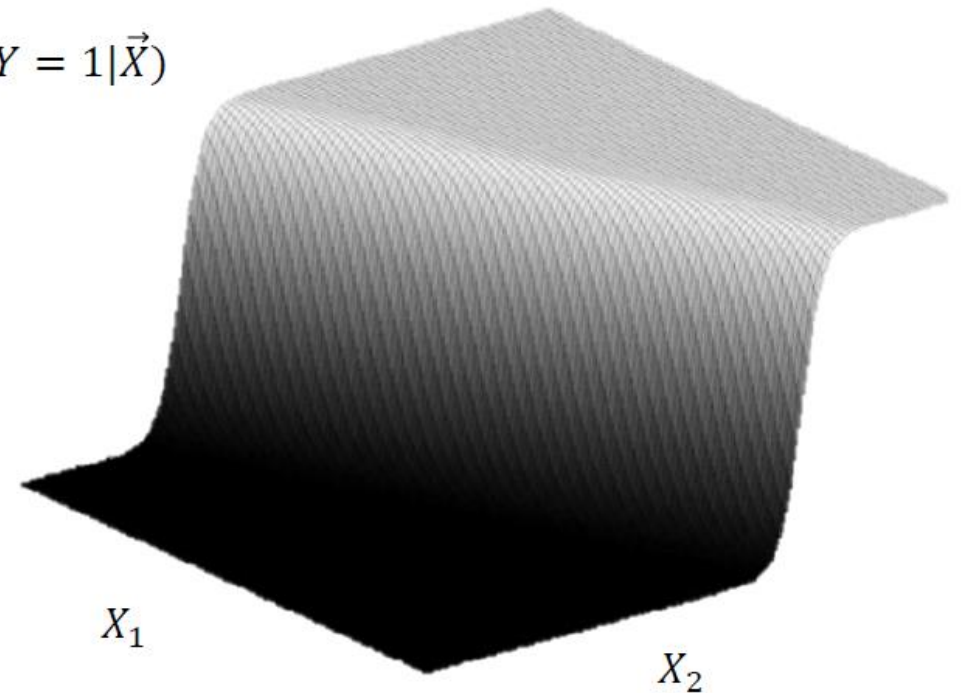
Logistic Regression Decision Boundary

$$y^{(i)} = \begin{cases} 1 & \sigma(\boldsymbol{\beta}^T \mathbf{x}^{(i)}) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

- Other thresholds might be needed in some cases!



$$P(Y = 1 | \vec{X})$$



Maximum Likelihood Estimation (MLE)

The problem of finding two probabilities reduces to estimation of regression parameters:

- *Maximum likelihood estimation (MLE)*: Finding the parameters that maximize the joint probability of occurring all observed samples

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n \mathbb{P}(Y = y^{(i)} | \mathbf{X} = \mathbf{x}^{(i)}) = \prod_{i=1}^n \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}^{(i)})^{y^{(i)}} \mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x}^{(i)})^{1-y^{(i)}}$$

$$\rightarrow L(\boldsymbol{\beta}) = \prod_{i=1}^n \sigma(\boldsymbol{\beta}^T \mathbf{x}^{(i)})^{y^{(i)}} (1 - \sigma(\boldsymbol{\beta}^T \mathbf{x}^{(i)}))^{1-y^{(i)}}$$

$$\begin{aligned} l(\boldsymbol{\beta}) = \log L(\boldsymbol{\beta}) &= \sum_{i=1}^n y^{(i)} \log(\sigma(\boldsymbol{\beta}^T \mathbf{x}^{(i)})) + (1 - y^{(i)}) (\log(1 - \sigma(\boldsymbol{\beta}^T \mathbf{x}^{(i)}))) = \\ &= \sum_{i=1}^n (y^{(i)} (\boldsymbol{\beta}^T \mathbf{x}^{(i)}) - \log(1 + e^{\boldsymbol{\beta}^T \mathbf{x}^{(i)}})) \end{aligned}$$

Cost Function

Since sigmoid outputs probability, we use negative log likelihood to represent the error:

$$J(\boldsymbol{\beta}) = -\frac{1}{n}l(\boldsymbol{\beta}) = -\frac{1}{n}\sum_{i=1}^n \left(y^{(i)}(\boldsymbol{\beta}^T \mathbf{x}^{(i)}) - \log(1 + e^{\boldsymbol{\beta}^T \mathbf{x}^{(i)}}) \right)$$

We can update the parameters by minimizing the cost function (or equivalently, maximizing likelihood function)

Estimation of the Parameters

$$J(\boldsymbol{\beta}) = -\frac{1}{n} \sum_{i=1}^n \left(y^{(i)} (\boldsymbol{\beta}^T \mathbf{x}^{(i)}) - \log(1 + e^{\boldsymbol{\beta}^T \mathbf{x}^{(i)}}) \right)$$

$$\begin{aligned} \frac{\partial J(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= -\frac{1}{n} \sum_{i=1}^n \left(y^{(i)} \mathbf{x}^{(i)} - \mathbf{x}^{(i)} \frac{e^{\boldsymbol{\beta}^T \mathbf{x}^{(i)}}}{1 + e^{\boldsymbol{\beta}^T \mathbf{x}^{(i)}}} \right) = -\frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \left(y^{(i)} - \frac{e^{\boldsymbol{\beta}^T \mathbf{x}^{(i)}}}{1 + e^{\boldsymbol{\beta}^T \mathbf{x}^{(i)}}} \right) = \\ &\quad \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} (\sigma(\boldsymbol{\beta}^T \mathbf{x}^{(i)}) - y^{(i)}) \end{aligned}$$

$$\rightarrow \frac{\partial J(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} (\sigma(\boldsymbol{\beta}^T \mathbf{x}^{(i)}) - y^{(i)}) = \frac{1}{n} \mathbf{X}^T (\boldsymbol{\sigma}(\mathbf{X}\boldsymbol{\beta}) - \mathbf{y})$$

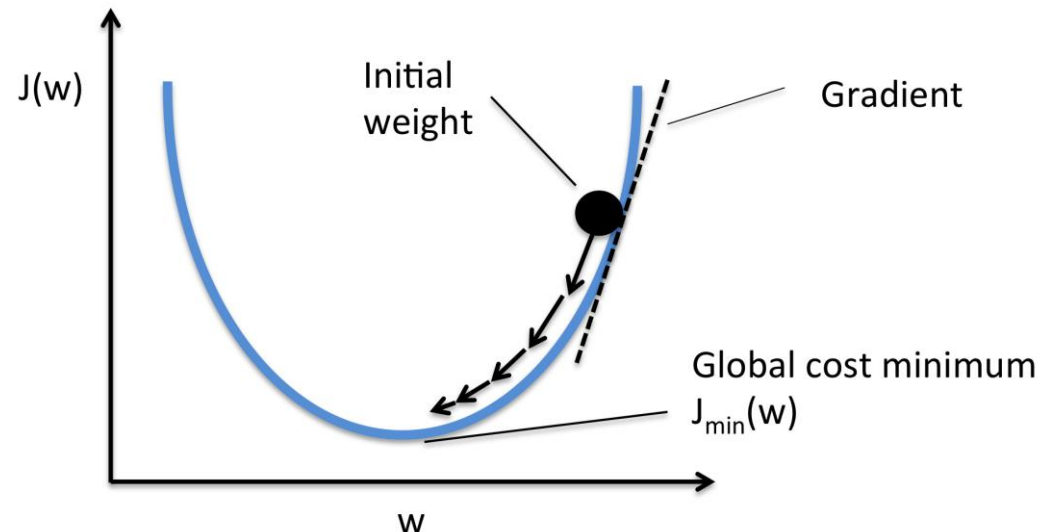
$$\text{where } \mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_p^{(1)} \\ \vdots & \vdots & \dots & \vdots \\ 1 & x_1^{(n)} & \dots & x_p^{(n)} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \boldsymbol{\sigma}(\mathbf{X}\boldsymbol{\beta}) = \begin{bmatrix} \sigma(\boldsymbol{\beta}^T \mathbf{x}^{(1)}) \\ \sigma(\boldsymbol{\beta}^T \mathbf{x}^{(2)}) \\ \vdots \\ \sigma(\boldsymbol{\beta}^T \mathbf{x}^{(n)}) \end{bmatrix}$$

Updating the Parameters using Gradient Descent

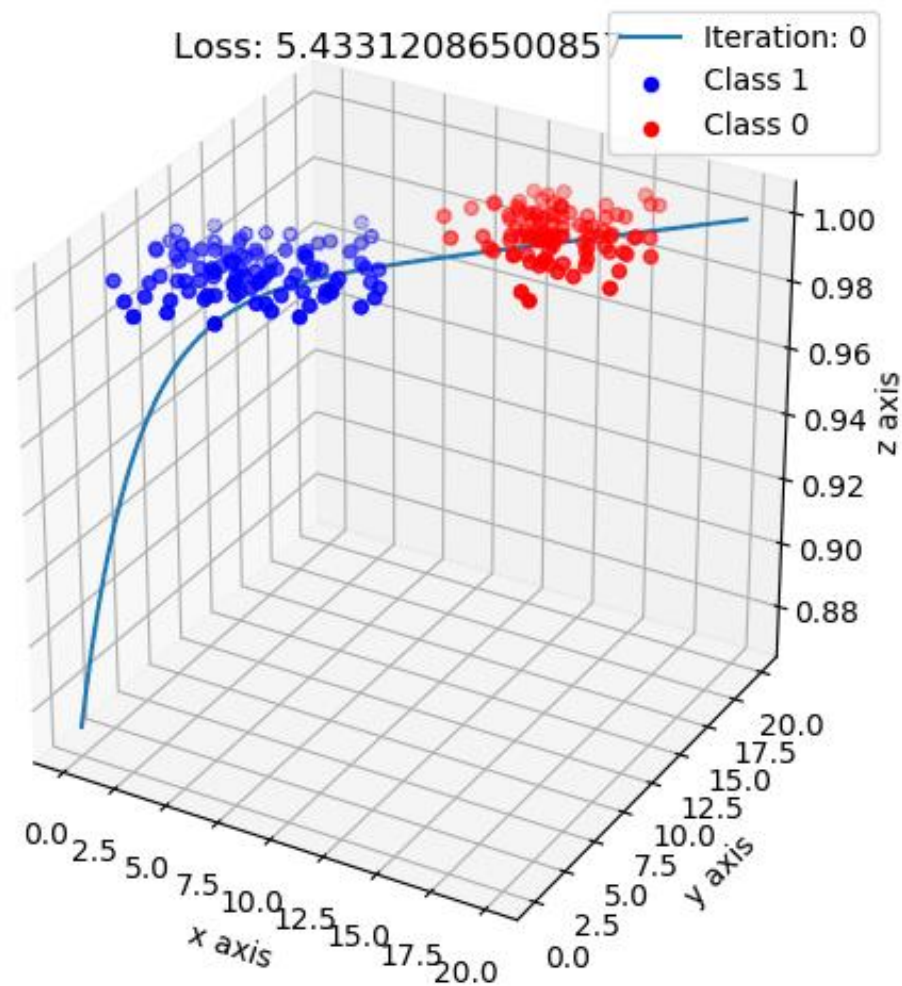
Now that we have the derivatives we can update the parameters using *Gradient Descent*:

$$\boldsymbol{\beta}^{new} = \boldsymbol{\beta}^{old} - \alpha \frac{\partial J(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \boldsymbol{\beta}^{old} - \alpha \frac{1}{n} \mathbf{X}^T (\boldsymbol{\sigma}(\mathbf{X}\boldsymbol{\beta}) - \mathbf{y})$$

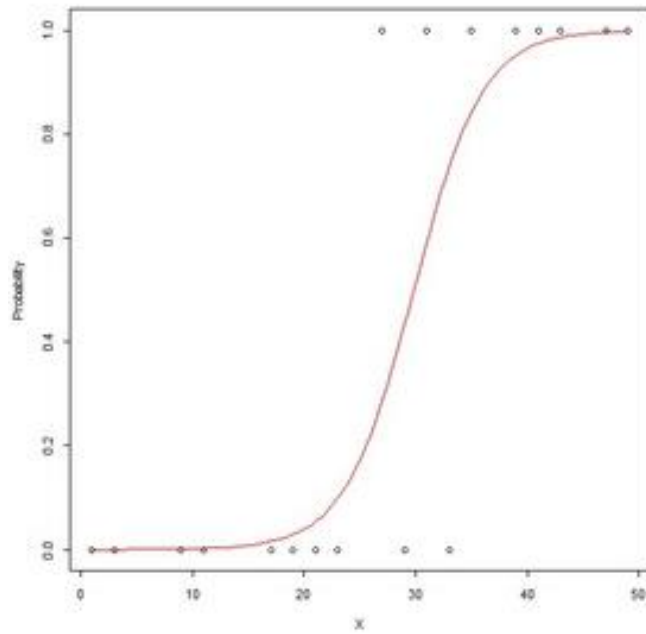
- α (*learning rate*): controls the rate of convergence (should be moderate).



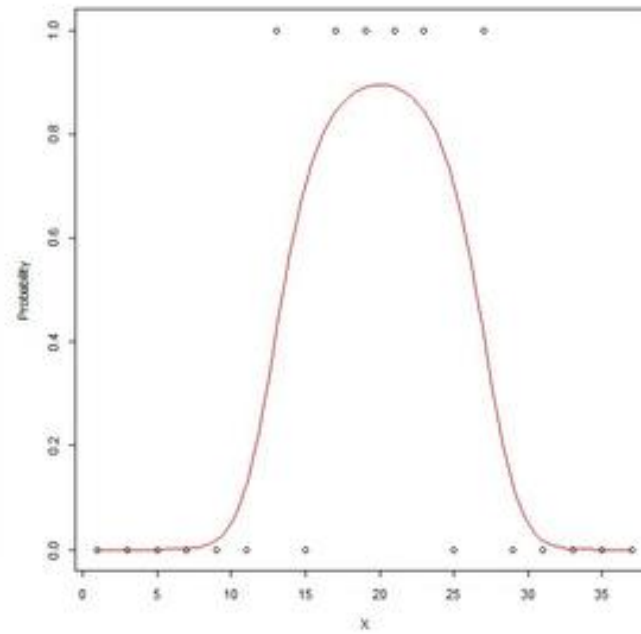
Example



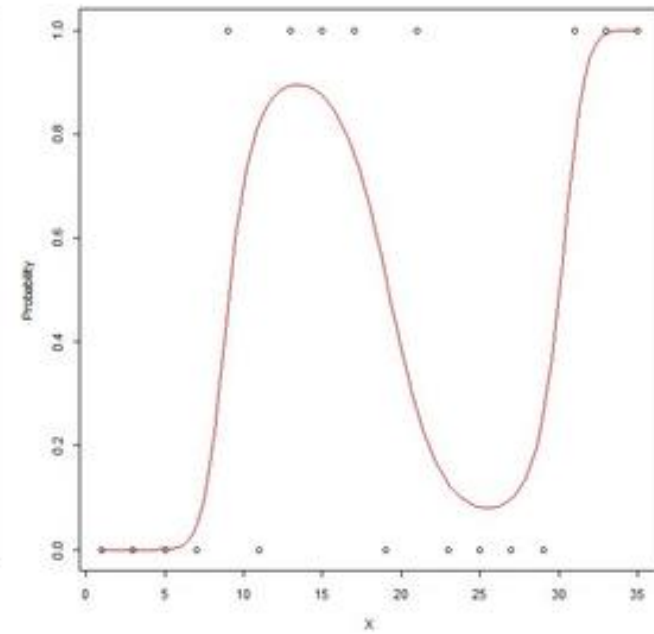
Non-Linear Logistic Regression



Linear



Quadratic



Cubic

*Multi-Class Logistic Regression

We can use logistic regression for multi-class classification as well. This is called *Multinomial Logistic Regression* or *Softmax Regression*.

- Each class has its own probability function with its own parameter vector, s.t. for any given x , they all sum to one.

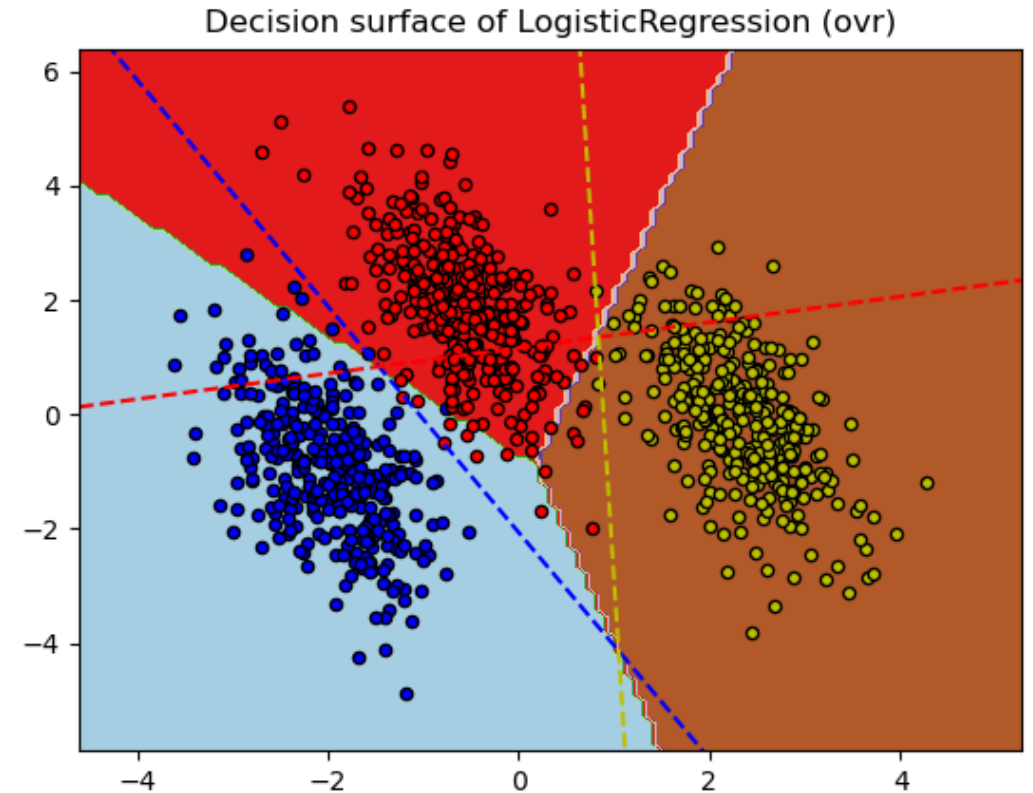
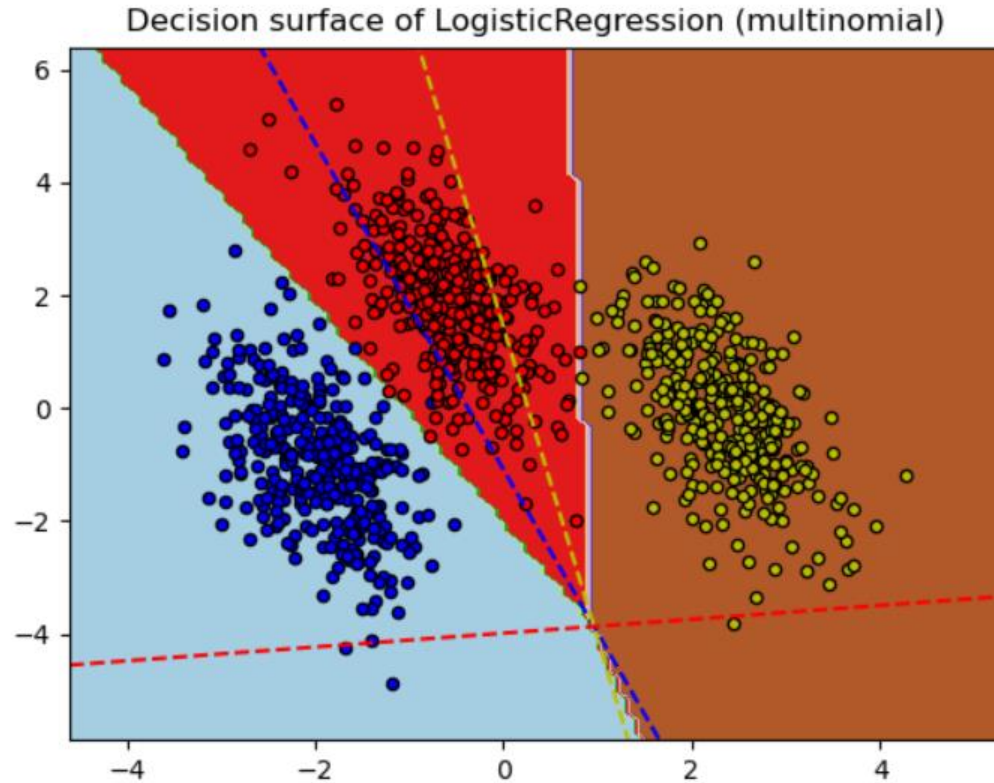
$$\mathbb{P}_j(\mathbf{x}, \boldsymbol{\beta}) = \mathbb{P}(Y = j | \mathbf{x}, \boldsymbol{\beta}) = \frac{e^{\boldsymbol{\beta}_j^T \mathbf{x}}}{\sum_{j=1}^K e^{\boldsymbol{\beta}_j^T \mathbf{x}}}, \quad j = 1, 2, \dots, K$$

Other approaches:

- One-vs-All
- One-vs-One

- In softmax model we should use penalty term to keep the parameters small. (L_2 or L_1 regularizers)

*Multi-Class Logistic Regression Decision Surface



*Regularized Logistic Regression

Like linear regression, in order to reduce the variance of the model, we can use regularization for the logistic regression as well:

- L_2 norm regularizer:

$$cost = -l(\boldsymbol{\beta}) + \lambda ||\boldsymbol{\beta}||_{L_2}^2 = -l(\boldsymbol{\beta}) + \lambda \sum_{j=1}^p \beta_j^2$$

- L_1 norm regularizer:

$$cost = -l(\boldsymbol{\beta}) + \lambda ||\boldsymbol{\beta}||_{L_1} = -l(\boldsymbol{\beta}) + \lambda \sum_{j=1}^p |\beta_j|$$

- Scikit-learn uses L_2 norm regularization by default.