

Statistics 101C - Binary Classification: Surrogate Loss functions

Shirong Xu

University of California, Los Angeles

shirong@stat.ucla.edu

October 10, 2023

Question

- **Question:** If you are asked to minimize or maximize a function $f : \mathbb{R} \rightarrow \mathbb{R}$, what would you do?

Question

- **Question:** If you are asked to minimize or maximize a function $f : \mathbb{R} \rightarrow \mathbb{R}$, what would you do?
- **Answer:** Take the derivative with respect to x and solve the equation

$$\frac{d}{dx}f(x) = 0$$

Question

- **Question:** If you are asked to minimize or maximize a function $f : \mathbb{R} \rightarrow \mathbb{R}$, what would you do?
- **Answer:** Take the derivative with respect to x and solve the equation

$$\frac{d}{dx}f(x) = 0$$

- **Question:** What is the equation $\frac{d}{dx}f(x) = 0$ is hard to solve?

Question

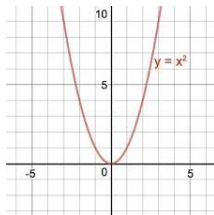
- **Question:** If you are asked to minimize or maximize a function $f : \mathbb{R} \rightarrow \mathbb{R}$, what would you do?
- **Answer:** Take the derivative with respect to x and solve the equation

$$\frac{d}{dx}f(x) = 0$$

- **Question:** What is the equation $\frac{d}{dx}f(x) = 0$ is hard to solve?
- **Answer:** Gradient ascent (For maximization) or Gradient descent (For minimization)

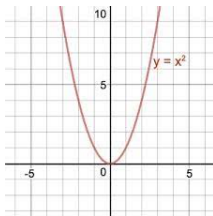
Gradient Descent: An Example

- Suppose we want to minimize $f(x) = x^2$



Gradient Descent: An Example

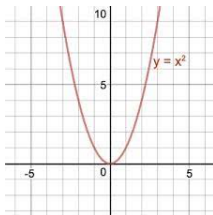
- Suppose we want to minimize $f(x) = x^2$



- First, we randomly choose an initial point $x^{(0)} = 3$

Gradient Descent: An Example

- Suppose we want to minimize $f(x) = x^2$

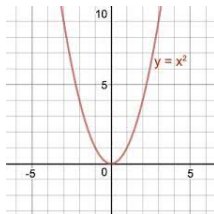


- First, we randomly choose an initial point $x^{(0)} = 3$
- Second, we calculate the derivative at the point $x^{(0)} = 3$

$$f'(3) = 6$$

Gradient Descent: An Example

- Suppose we want to minimize $f(x) = x^2$



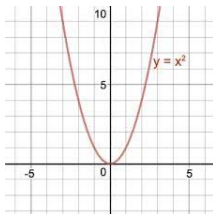
- First, we randomly choose an initial point $x^{(0)} = 3$
- Second, we calculate the derivative at the point $x^{(0)} = 3$

$$f'(3) = 6$$

- Update the new value $x^{(1)} = x^{(0)} - \lambda \cdot f'(3)$, λ is the step size.

Gradient Descent: An Example

- Suppose we want to minimize $f(x) = x^2$



- First, we randomly choose an initial point $x^{(0)} = 3$
- Second, we calculate the derivative at the point $x^{(0)} = 3$

$$f'(3) = 6$$

- Update the new value $x^{(1)} = x^{(0)} - \lambda \cdot f'(3)$, λ is the step size.
- Repeat the above process until convergence.

Gradient Descent

- **Motivation of Gradient descent:** gradient provides information to minimize the objective function.
- **General form of gradient descent:** Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a p -variate function. Then gradient descent has the form

$$\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)} - \lambda \nabla f(\mathbf{x}),$$

where $\nabla f(\mathbf{x}) = \left(\frac{\partial}{\partial x_1} f(\mathbf{x}), \frac{\partial}{\partial x_2} f(\mathbf{x}), \dots, \frac{\partial}{\partial x_p} f(\mathbf{x}) \right)$

Gradient Descent

- **Motivation of Gradient descent:** gradient provides information to minimize the objective function.
- **General form of gradient descent:** Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a p -variate function. Then gradient descent has the form

$$\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)} - \lambda \nabla f(\mathbf{x}),$$

where $\nabla f(\mathbf{x}) = \left(\frac{\partial}{\partial x_1} f(\mathbf{x}), \frac{\partial}{\partial x_2} f(\mathbf{x}), \dots, \frac{\partial}{\partial x_p} f(\mathbf{x}) \right)$

- **Question:** When does gradient descent stop?

Gradient Descent

- **Motivation of Gradient descent:** gradient provides information to minimize the objective function.
- **General form of gradient descent:** Let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be a p -variate function. Then gradient descent has the form

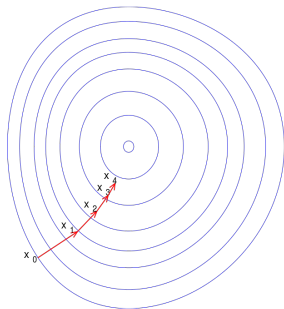
$$\mathbf{x}^{(t)} = \mathbf{x}^{(t-1)} - \lambda \nabla f(\mathbf{x}),$$

where $\nabla f(\mathbf{x}) = \left(\frac{\partial}{\partial x_1} f(\mathbf{x}), \frac{\partial}{\partial x_2} f(\mathbf{x}), \dots, \frac{\partial}{\partial x_p} f(\mathbf{x}) \right)$

- **Question:** When does gradient descent stop?
- **Answer:** When $\nabla f(\mathbf{x}) \approx \mathbf{0}$

Application of Gradient Descent

- Gradient descent is usually employed, when optimization problem is non-convex or does not have **analytic solution** or high-dimensional \mathbf{x} , for example Deep neural networks (Non-convex)



- Convex Optimization Problem + Gradient Descent \Rightarrow Optimal Solution

Classification

- A typical dataset in classification $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
 - \mathbf{x}_i : the covariate vector of i -th instance
 - $y_i \in \{-1, 1\}$: binary label of i -th instance

Classification

- A typical dataset in classification $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
 - \mathbf{x}_i : the covariate vector of i -th instance
 - $y_i \in \{-1, 1\}$: binary label of i -th instance
- **Question:** Can we directly minimize the averaged 0-1 loss?

$$\text{Training Error} : \frac{1}{n} \sum_{i=1}^n I(f(\mathbf{x}_i) \neq y_i)$$

Classification

- A typical dataset in classification $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
 - \mathbf{x}_i : the covariate vector of i -th instance
 - $y_i \in \{-1, 1\}$: binary label of i -th instance
- **Question:** Can we directly minimize the averaged 0-1 loss?

$$\text{Training Error} : \frac{1}{n} \sum_{i=1}^n I(f(\mathbf{x}_i) \neq y_i)$$

- **Answer:** No, the 0-1 loss function is non-convex and discontinuous, so (sub)gradient methods cannot be applied.

Classification - Surrogate Loss

- We can replace the 0-1 loss by other loss functions

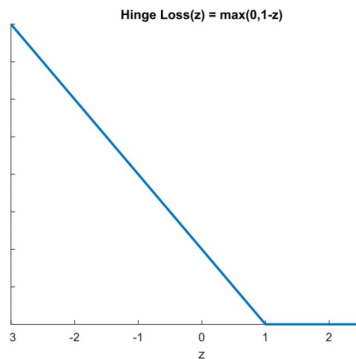
$$\frac{1}{n} \sum_{i=1}^n I(f(\mathbf{x}_i) \neq y_i) \Rightarrow \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i), y_i) = \frac{1}{n} \sum_{i=1}^n \phi(f(\mathbf{x}_i)y_i)$$

- Hinge loss: $\phi(x) = \max\{0, 1 - x\}$
- Logistic loss $\phi(x) = \log(1 + \exp(-x))$

Classification - Hinge Loss

- Definition of Hinge loss:

$$L_{\text{hinge}}(f(\mathbf{x}_i), y_i) = \begin{cases} 1 - f(\mathbf{x}_i)y_i & \text{if } f(\mathbf{x}_i)y_i \leq 1 \\ 0, & \text{if } f(\mathbf{x}_i)y_i > 1 \end{cases}$$



Classification - Hinge Loss

- Let $\eta(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$. The expected hinge loss: hinge risk.

$$\begin{aligned} R_{\text{hinge}}(f) &= \mathbb{E}_{\mathbf{X}, Y} [L_{\text{hinge}}(f(\mathbf{X}), Y)] \\ &= \mathbb{E}_{\mathbf{X}} \left[\eta(\mathbf{X})(1 - f(\mathbf{X}))_+ + (1 - \eta(\mathbf{X}))(1 + f(\mathbf{X}))_+ \right] \end{aligned}$$

Classification - Hinge Loss

- Let $\eta(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$. The expected hinge loss: hinge risk.

$$\begin{aligned} R_{\text{hinge}}(f) &= \mathbb{E}_{\mathbf{X}, Y} [L_{\text{hinge}}(f(\mathbf{X}), Y)] \\ &= \mathbb{E}_{\mathbf{X}} \left[\eta(\mathbf{X})(1 - f(\mathbf{X}))_+ + (1 - \eta(\mathbf{X}))(1 + f(\mathbf{X}))_+ \right] \end{aligned}$$

- Suppose that $f(\mathbf{X}) \in [-1, 1]$, for any \mathbf{X} , we have

$$\begin{aligned} &\eta(\mathbf{X})(1 - f(\mathbf{X})) + (1 - \eta(\mathbf{X}))(1 + f(\mathbf{X})) \\ &= \eta(\mathbf{X}) - 2\eta(\mathbf{X})f(\mathbf{X}) + 1 + f(\mathbf{X}) - \eta(\mathbf{X}) \\ &= f(\mathbf{X})(1 - 2\eta(\mathbf{X})) + 1. \end{aligned}$$

Classification - Hinge Loss

- Let $\eta(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$. The expected hinge loss: hinge risk.

$$\begin{aligned} R_{\text{hinge}}(f) &= \mathbb{E}_{\mathbf{X}, Y} [L_{\text{hinge}}(f(\mathbf{X}), Y)] \\ &= \mathbb{E}_{\mathbf{X}} \left[\eta(\mathbf{X})(1 - f(\mathbf{X}))_+ + (1 - \eta(\mathbf{X}))(1 + f(\mathbf{X}))_+ \right] \end{aligned}$$

- Suppose that $f(\mathbf{X}) \in [-1, 1]$, for any \mathbf{X} , we have

$$\begin{aligned} &\eta(\mathbf{X})(1 - f(\mathbf{X})) + (1 - \eta(\mathbf{X}))(1 + f(\mathbf{X})) \\ &= \eta(\mathbf{X}) - 2\eta(\mathbf{X})f(\mathbf{X}) + 1 + f(\mathbf{X}) - \eta(\mathbf{X}) \\ &= f(\mathbf{X})(1 - 2\eta(\mathbf{X})) + 1. \end{aligned}$$

- The optimal function f_{hinge}^* minimizing $R_{\text{hinge}}(f)$
 - If $1 - 2\eta(\mathbf{X}) > 0$, hinge loss is minimized at $f(\mathbf{X}) = -1$
 - If $1 - 2\eta(\mathbf{X}) < 0$, hinge loss is minimized at $f(\mathbf{X}) = 1$

Classification - Hinge Loss

- Recall that the optimal classifier (Bayes classifier) of Binary loss is defined as

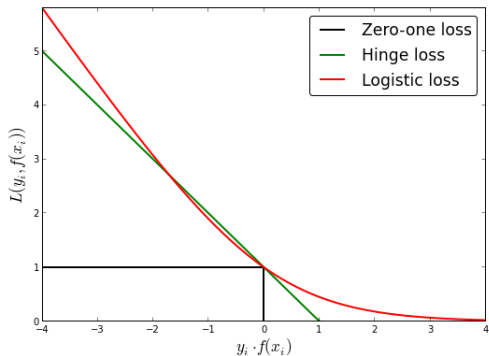
$$f^*(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}) - 1/2) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) > 1/2 \\ 0 & \text{if } \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) < 1/2 \end{cases}$$

- Observation:** f_{hinge}^* is exactly the Bayes classifier and the hinge loss is a **convex function**, which makes it possible to minimize the training error.

Classification - Logistic Loss

- Definition of Logistic loss:

$$L_{\log}(f(\mathbf{x}_i), y_i) = \log(1 + \exp(-f(\mathbf{x}_i)y_i))$$



Classification - Logistic Loss

- The logistic risk:

$$\begin{aligned} R_{\log}(f) &= \mathbb{E}_{\mathbf{X}, Y} \left[\log \left(1 + \exp(-f(\mathbf{X})Y) \right) \right] \\ &= \mathbb{E}_{\mathbf{X}} \left[\eta(\mathbf{X}) \log \left(1 + \exp(-f(\mathbf{X})) \right) + (1 - \eta(\mathbf{X})) \log \left(1 + \exp(f(\mathbf{X})) \right) \right] \end{aligned}$$

Classification - Logistic Loss

- The logistic risk:

$$\begin{aligned} R_{\log}(f) &= \mathbb{E}_{\mathbf{X}, Y} \left[\log \left(1 + \exp(-f(\mathbf{X})Y) \right) \right] \\ &= \mathbb{E}_{\mathbf{X}} \left[\eta(\mathbf{X}) \log \left(1 + \exp(-f(\mathbf{X})) \right) + (1 - \eta(\mathbf{X})) \log \left(1 + \exp(f(\mathbf{X})) \right) \right] \end{aligned}$$

- Take the derivative with respect to f , it follows that

$$\begin{aligned} & -\eta(\mathbf{X}) \frac{\exp(-f(\mathbf{X}))}{1 + \exp(-f(\mathbf{X}))} + (1 - \eta(\mathbf{X})) \frac{\exp(f(\mathbf{X}))}{1 + \exp(f(\mathbf{X}))} \\ &= -\eta(\mathbf{X}) \frac{1}{1 + \exp(f(\mathbf{X}))} + (1 - \eta(\mathbf{X})) \frac{\exp(f(\mathbf{X}))}{1 + \exp(f(\mathbf{X}))} \\ &= \frac{\exp(f(\mathbf{X}))}{1 + \exp(f(\mathbf{X}))} - \eta(\mathbf{X}) = 0 \iff f_{\log}^*(\mathbf{X}) = \log \frac{\eta(\mathbf{X})}{1 - \eta(\mathbf{X})} \end{aligned}$$

Connection between binary loss and surrogate losses

- The Bayes classifier $f^*(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}) - 1/2)$
- The optimal classifier of Hinge risk $f_{\text{hinge}}^*(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}) - 1/2)$
- The optimal classifier of Logistic risk $f_{\text{log}}^*(\mathbf{x}) = \log \frac{\eta(\mathbf{x})}{1-\eta(\mathbf{x})}$

Connection between binary loss and surrogate losses

- The Bayes classifier $f^*(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}) - 1/2)$
- The optimal classifier of Hinge risk $f_{\text{hinge}}^*(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}) - 1/2)$
- The optimal classifier of Logistic risk $f_{\text{log}}^*(\mathbf{x}) = \log \frac{\eta(\mathbf{x})}{1-\eta(\mathbf{x})}$
- **Question:** what is the connection between these optimal classifiers?

Connection between binary loss and surrogate losses

- The Bayes classifier $f^*(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}) - 1/2)$
- The optimal classifier of Hinge risk $f_{\text{hinge}}^*(\mathbf{x}) = \text{sign}(\eta(\mathbf{x}) - 1/2)$
- The optimal classifier of Logistic risk $f_{\text{log}}^*(\mathbf{x}) = \log \frac{\eta(\mathbf{x})}{1-\eta(\mathbf{x})}$
- **Question:** what is the connection between these optimal classifiers?
- **Answer:** They are consistent in sign.

Classification

- A typical dataset in classification $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
 - \mathbf{x}_i : the covariate vector of i -th instance
 - $y_i \in \{0, 1\}$: binary label of i -th instance
- Bayes classifier f^* :

$$f^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) > 1/2 \\ 0 & \text{if } \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) < 1/2 \end{cases}$$

- Minimal risk $R(f^*)$:

$$R(f^*) = \mathbb{E}[f^*(\mathbf{X}) \neq Y] = \mathbb{E}[\min(\eta(\mathbf{X}), 1 - \eta(\mathbf{X}))],$$

where $\eta(\mathbf{x}) = \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})$.

How can we construct classifier?

- **Discriminative models**

- Discriminative modeling studies the $P(Y|\mathbf{X})$
- Examples: Logistic regression (LR) and Support Vector Machine (SVM)

How can we construct classifier?

- **Discriminative models**

- Discriminative modeling studies the $P(Y|\mathbf{X})$
- Examples: Logistic regression (LR) and Support Vector Machine (SVM)

- **Generative models**

- Generative models studies the joint probability distribution $\mathbb{P}(\mathbf{X}, Y)$
- Examples: linear discriminant analysis and quadratic discriminant analysis

Logistic Regression

- **Logistic regression** estimates the conditional probability probability:

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

- In logistic regression, it is assumed that

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}^T \mathbf{x})},$$

where

- $\mathbf{x} = (x_1, \dots, x_p)^T$ is a p -dimensional predictor
- β_0 and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$ are unknown parameters
- $\boldsymbol{\beta}^T \mathbf{x} = \sum_{i=1}^p \beta_i x_i$

Why Logistic Regression

- The **odds ratio**: the probability that $Y = 1$ divided by the probability that $Y = 0$ conditional on $\mathbf{X} = \mathbf{x}$.

$$\exp(\beta_0 + \beta^T \mathbf{x}) = \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{1 - \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})} = \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})}$$

Why Logistic Regression

- The **odds ratio**: the probability that $Y = 1$ divided by the probability that $Y = 0$ conditional on $\mathbf{X} = \mathbf{x}$.

$$\exp(\beta_0 + \beta^T \mathbf{x}) = \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{1 - \mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})} = \frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})}$$

- The **log-odds**: linear with respect to β

$$\beta_0 + \beta^T \mathbf{x} = \log \left(\frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x})} \right)$$

- The log-odds take any value in \mathbb{R}
- The log-odds equals a linear combination of the predictors.
- β_i can then be interpreted as the average change in the log-odds ratio given by a one-unit increase in x_i

Estimation in Logistic Regression

- Likelihood function $L(\beta_0, \boldsymbol{\beta})$:

$$L(\beta_0, \boldsymbol{\beta}) = \prod_{i=1}^n \left(\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) \right)^{y_i} \left(\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x}) \right)^{1-y_i}$$

Estimation in Logistic Regression

- Likelihood function $L(\beta_0, \beta)$:

$$L(\beta_0, \beta) = \prod_{i=1}^n \left(\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) \right)^{y_i} \left(\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x}) \right)^{1-y_i}$$

- Logarithm of $L(\beta_0, \beta)$:

$$\begin{aligned} \log L(\beta_0, \beta) &= \sum_{i=1}^n \left[y_i \log \left(\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) \right) + (1 - y_i) \log \left(\mathbb{P}(Y = 0 | \mathbf{X} = \mathbf{x}) \right) \right] \\ &= \sum_{i=1}^n \left[y_i (\beta_0 + \beta^T \mathbf{x}) - \log \left(1 + \exp(\beta_0 + \beta^T \mathbf{x}) \right) \right] \end{aligned}$$

Estimation in Logistic Regression

Estimate β_0 and β (Gradient Ascent):

$$\begin{aligned}\beta_0^{(t+1)} &\leftarrow \beta_0^{(t)} + \lambda \sum_{i=1}^n \left[y_i - \frac{\exp(\beta_0^{(t)} + \beta^{(t)T} \mathbf{x})}{1 + \exp(\beta_0^{(t)} + \beta^{(t)T} \mathbf{x})} \right] \\ \beta^{(t+1)} &\leftarrow \beta^{(t)} + \lambda \sum_{i=1}^n \left[y_i - \frac{\exp(\beta_0^{(t)} + \beta^{(t)T} \mathbf{x})}{1 + \exp(\beta_0^{(t)} + \beta^{(t)T} \mathbf{x})} \right] \mathbf{x}_i\end{aligned}$$

Classification

Once we get estimated parameters, we have

$$\hat{\eta}(\mathbf{x}) = \frac{\exp(\hat{\beta}_0 + \hat{\beta}^T \mathbf{x})}{1 + \exp(\hat{\beta}_0 + \hat{\beta}^T \mathbf{x})}$$

Then we make predictions by

$$\hat{f}(\mathbf{x}) = \begin{cases} 1, & \text{if } \hat{\eta}(\mathbf{x}) > 1/2 \\ 0, & \text{if } \hat{\eta}(\mathbf{x}) < 1/2 \end{cases}$$

If $\hat{\eta}(\mathbf{x}) = 1/2$, then just randomly assign a label to it.

Example: Data Generation

Dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{5000}$, where $\mathbf{x}_i = (x_{i1}, x_{i2}, x_{i3}, x_{i4})$.

- Features are generated from uniform distribution $x_{il} \sim \text{Unif}(0, 2)$, $l = 1, 2, 3, 4$.
- $\beta_0 = 0.5$ and $\beta = (\beta_1, \dots, \beta_2)$ with $\beta_i \sim \text{Unif}(-1, 1)$ for $i = 1, 2, 3, 4$.
- Model:

$$y_i \sim \text{Bernoulli}\left(\frac{\exp(\beta_0 + \beta^T \mathbf{x})}{1 + \exp(\beta_0 + \beta^T \mathbf{x})}\right)$$

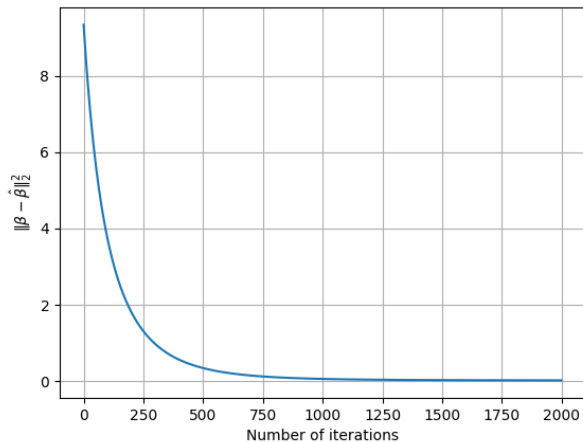
Python Codes

```
import numpy as np
np.random.seed(2)
n,p = 5000,4 # Set training datasize and dimension of features
X = np.random.uniform(-1,1,[n,p]) # Generation of features
beta = np.random.uniform(0,2,4) # Generation of parameters
beta_0 = 0.5 # Set the intercept term to 0.5
logOdd = (X * beta).sum(axis=1)+beta_0 # Log-odds
Prob = np.exp(logOdd)/(1+np.exp(logOdd)) # Probability
Y = np.array(Prob - np.random.uniform(0,1,n)>0,dtype=int) # Generate labels
```

Python Codes

```
Beta_0_hat = 0, # Initialization of intercept term
Beta_hat = np.zeros(p) # Initialization of beta
lamb = 0.1 # Learning rate
Error = [] # Error set
for i in range(2000): # Iterations of gradient ascent
    logOdd_hat = (X * Beta_hat).sum(axis=1)+Beta_0_hat
    Beta_0_hat = Beta_0_hat + lamb * np.mean(Y - np.exp(logOdd_hat)/(1+np.exp(logOdd_hat)))
    Beta_hat = Beta_hat + lamb * ((Y - np.exp(logOdd_hat)/(1+np.exp(logOdd_hat))) * X.T).mean(axis=1)
    Error.append(np.linalg.norm(Beta_hat-beta)**2)
import matplotlib.pyplot as plt
plt.plot(np.arange(0,2000),Error)
plt.xlabel('Number of iterations')
plt.ylabel('$\Vert \hat{\beta} - \beta \Vert_2^2$')
plt.grid()
```

Example: gradient ascent for logistic regression



Application of Logistic Regression to Banknote dataset

```
1 library(mclust)
2 data(banknote)
3 set.seed(123)
4 i <- 1:dim(banknote)[1]
5 i.train <- sample(i, 130, replace = FALSE)
6 bn.train <- banknote[i.train,]
7 bn.test <- banknote[-i.train,]
8 ml <- glm(Status Length + Right + Left + Top, data =
  bn.train,family = "binomial")
9 summary(ml)
10 pred.log.odds <- predict(ml)
11 pred.probs <- predict(ml, type = 'response')
12 my.thres <- 0.5
13 pred.log.odds.test <- predict(ml, bn.test[,-1])
14 pred.probs.test <- predict(ml, bn.test[,-1], type = 'response')
15 predicted.counterfeit <- pred.probs.test > my.thres
```

- `table('Reference' = bn.test[,1] == 'counterfeit', "Predicted" = predicted.counterfeit)`

	Predicted	
Reference	FALSE	TRUE
FALSE	2	31
TRUE	31	6