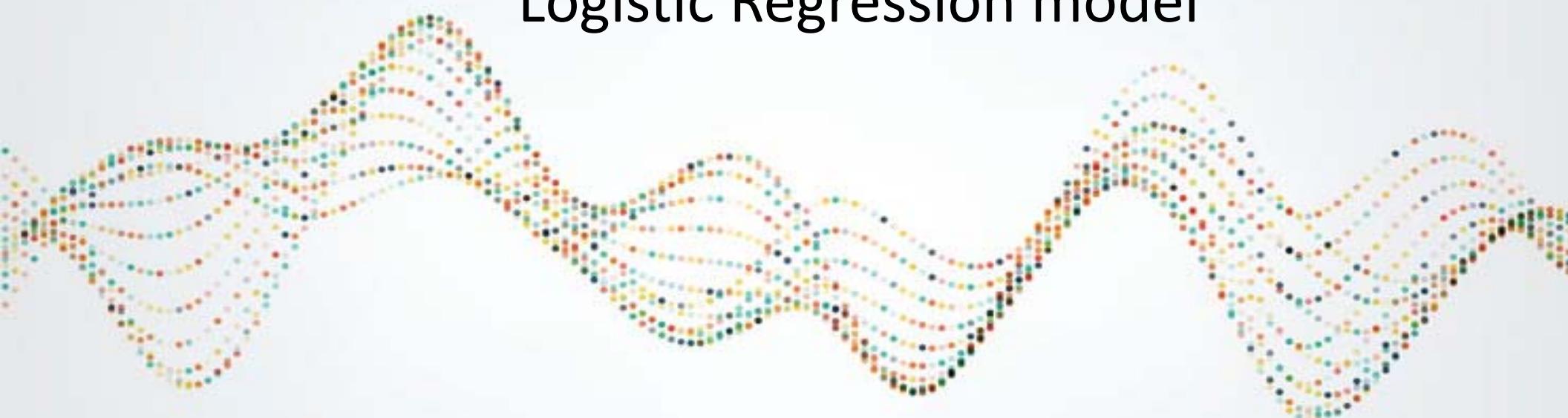


Credit card fraud detection using Logistic Regression model



Dr. Pritam Anand.
Assistant Professor,
DA-IICT, Gandhinagar.



Lottery Fraud



Juice Jacking



Card Skimming



Credit Card Fraud



Unknown mobile apps



Vishing Calls



Phishing Link

Federal Trade Commission (FTC) USA study [1]

- Nearly half of all American adults have had a fraudulent charge on their cards, amounting to around 127 million people. More than one in three card holders has experienced card fraud more than once^[2].
- The median fraudulent credit card charge was \$62 which approximate \$8 billion in fraudulent charges among all American consumers^[2].
- The credit card fraud has also become even more common since the start of the pandemic.
- Reports of credit card fraud increased by 44 percent between 2019 and 2020 according to the Federal Trade Commission (FTC) ^[1].

1. https://www.ftc.gov/system/files/documents/reports/consumer-sentinel-network-data-book-2020/csn_annual_data_book_2020.pdf
2. <https://www.security.org/digital-safety/credit-card-fraud-report/#references>

Scammers are growing smart..

They'll take notice of any new trends, like stimulus checks, unemployment payments, and general consumer financial trends, and capitalize upon them with phishing scams, identity theft, or other types of cybercrime.

Hari Ravichandran

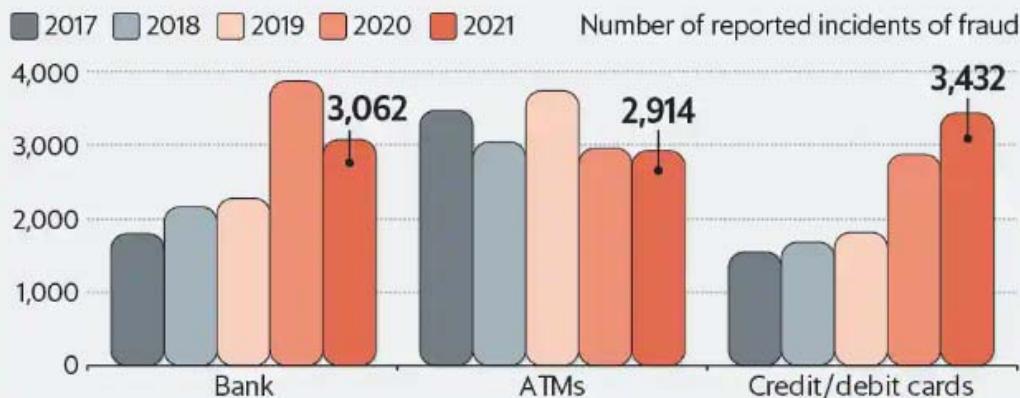
CEO, AURA
Identity Theft and fraud protection company



Indian Picture..

Fraud alert

NCRB's latest data shows that debit and credit card fraud is on the rise.

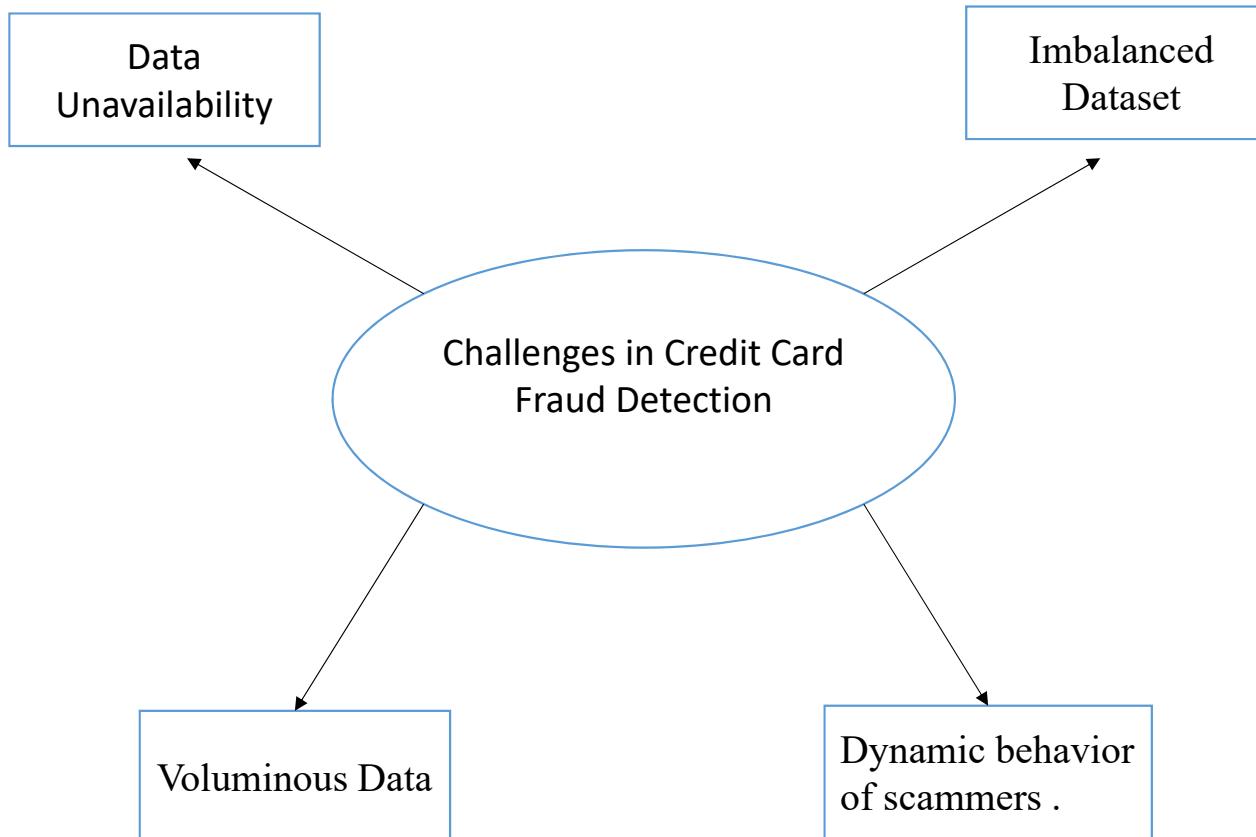


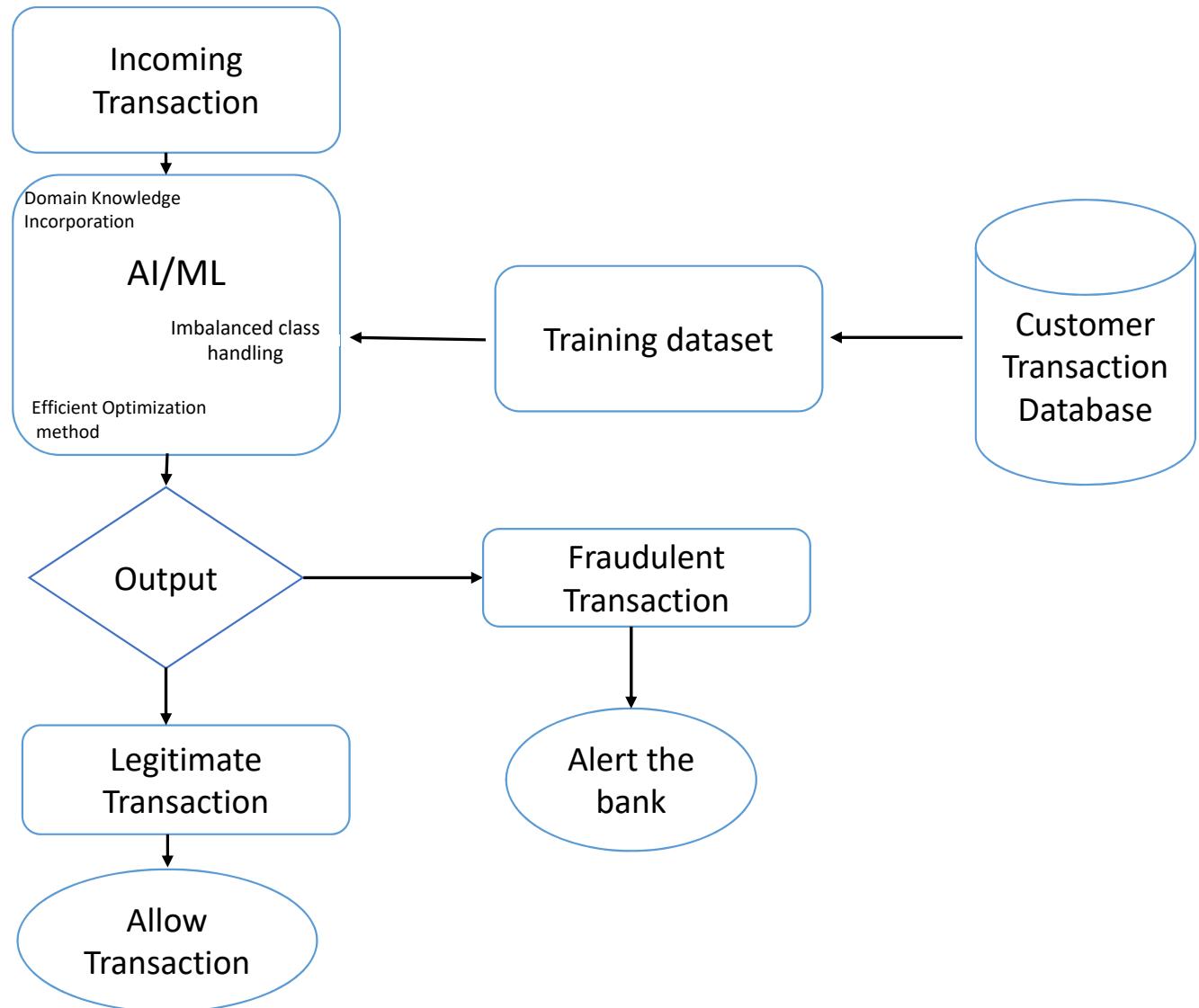
Report of National Crime Records Bureau, India [3]

According to the data, 3,432 cases of credit and debit card frauds were filed from across India in 2021, up nearly 20% from the year-earlier. In 2020, such frauds increased by over 70% [3].

[3] <https://www.livemint.com/industry/banking/debit-credit-card-frauds-on-rise-atm-scams-down-ncrb-11661885877307.html>

Challenges in Credit Card Fraud Detection

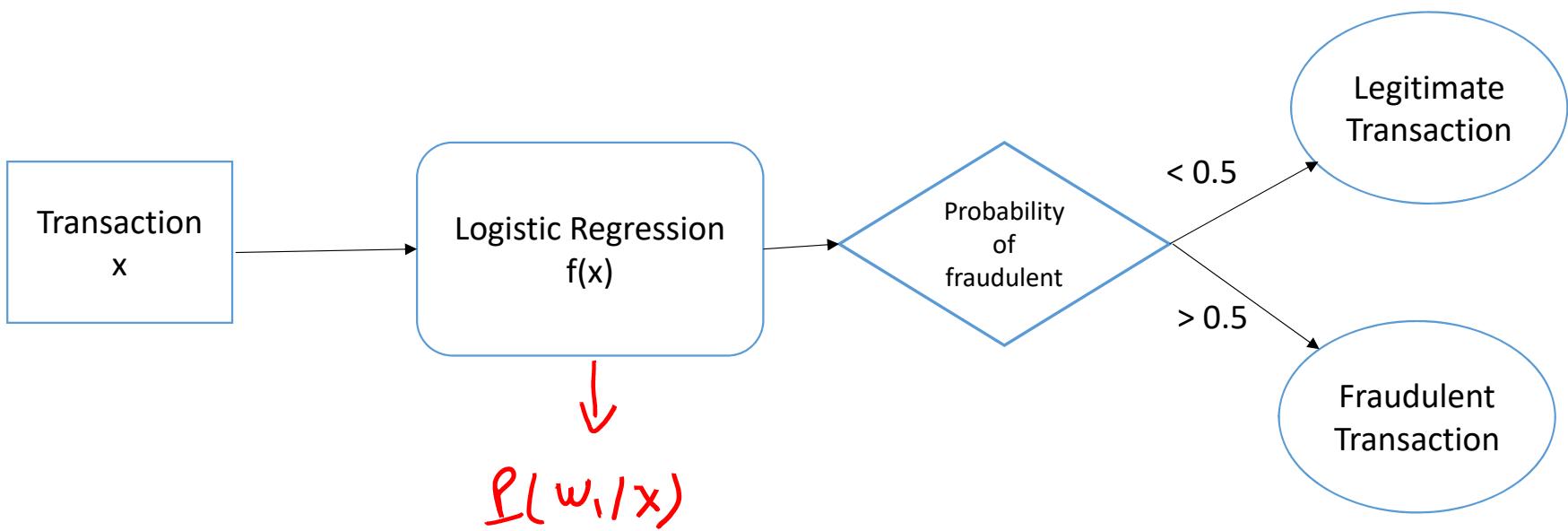


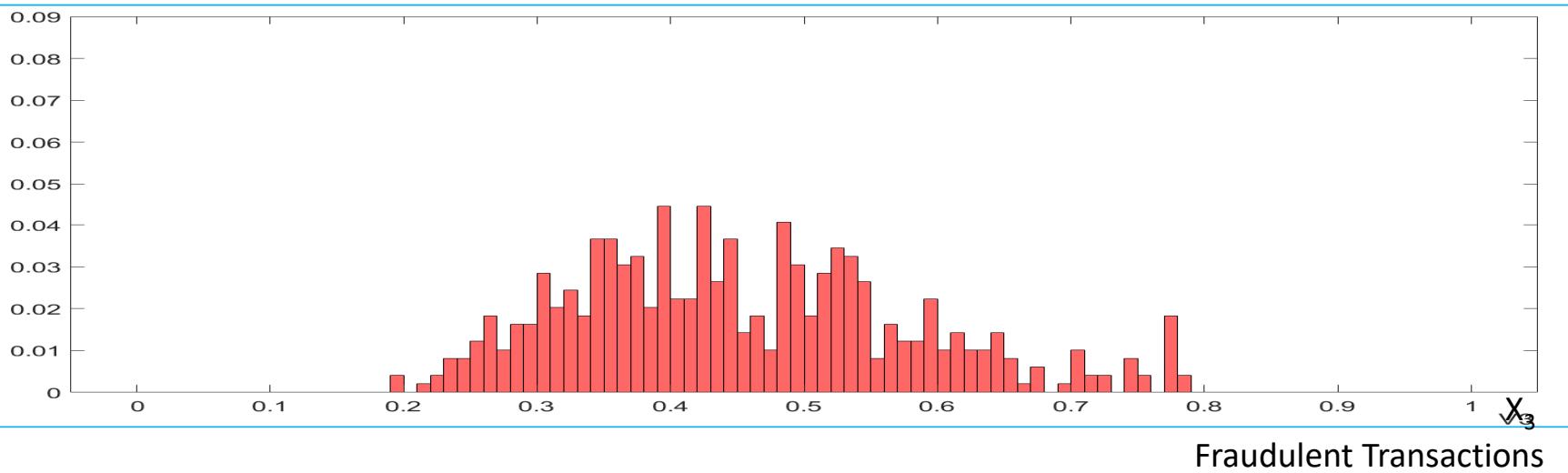


Real World Dataset

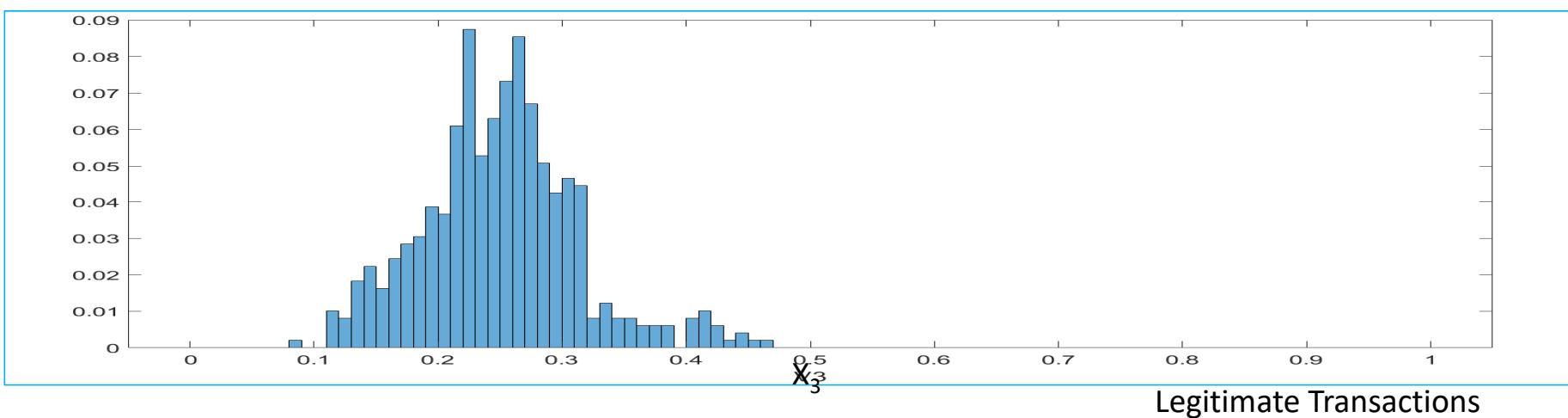
- The dataset is publicly available at <https://www.kaggle.com/datasets/mlg-ulb/creditcardfraud?select=creditcard.csv>
- It contains transactions made by credit cards in September 2013 by European cardholder in two days.
- Due to privacy issue, it does not provide original variables which provides the background information about the transaction.
- It contains 30 variables ,out of them X_1, X_2, \dots, X_{28} are transformed variables with Principal Component Analysis.
- Only Time and Amount is original variable which provides real information about the transaction.
- The dataset contains 492 fraudulent transaction out of 284807 transactions. The positive class account for 0.172% of all transaction.

Logistic Regression for credit card fraud detection

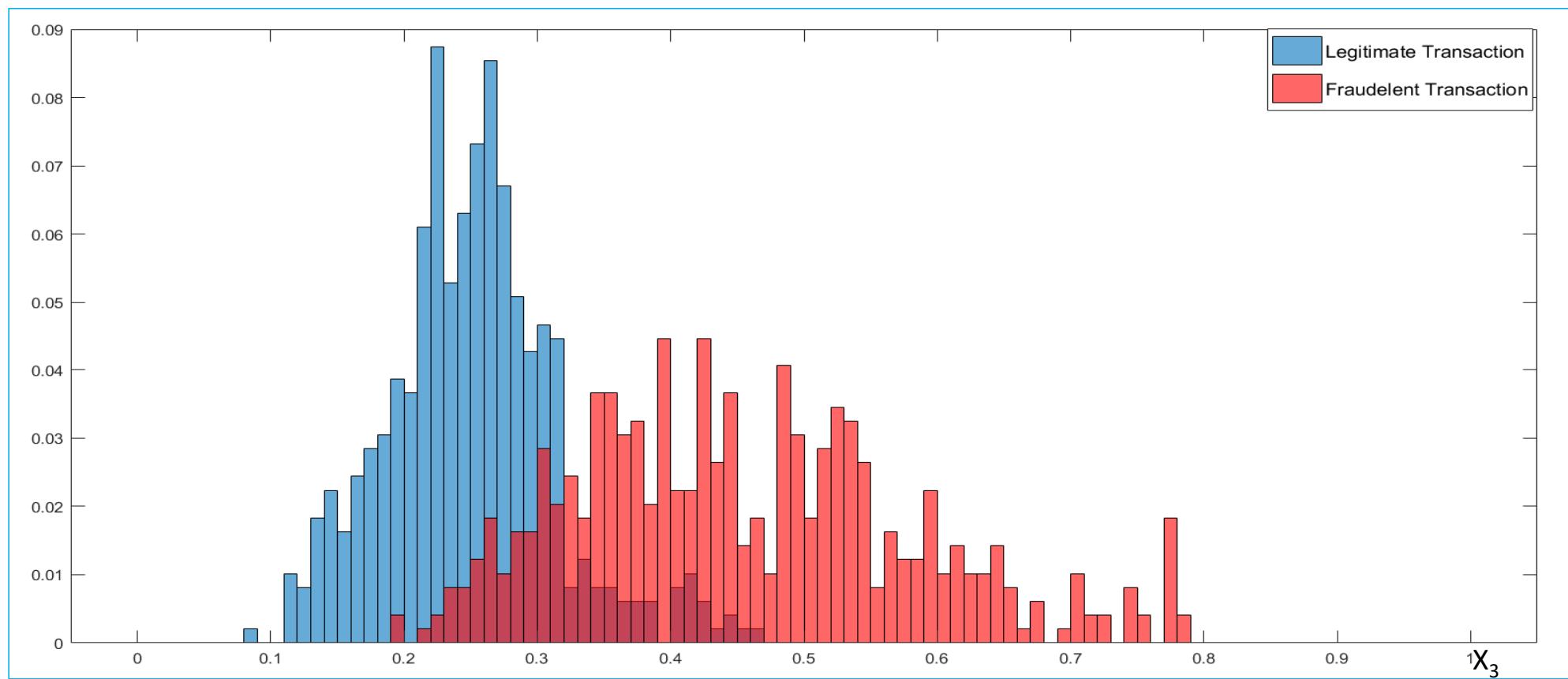


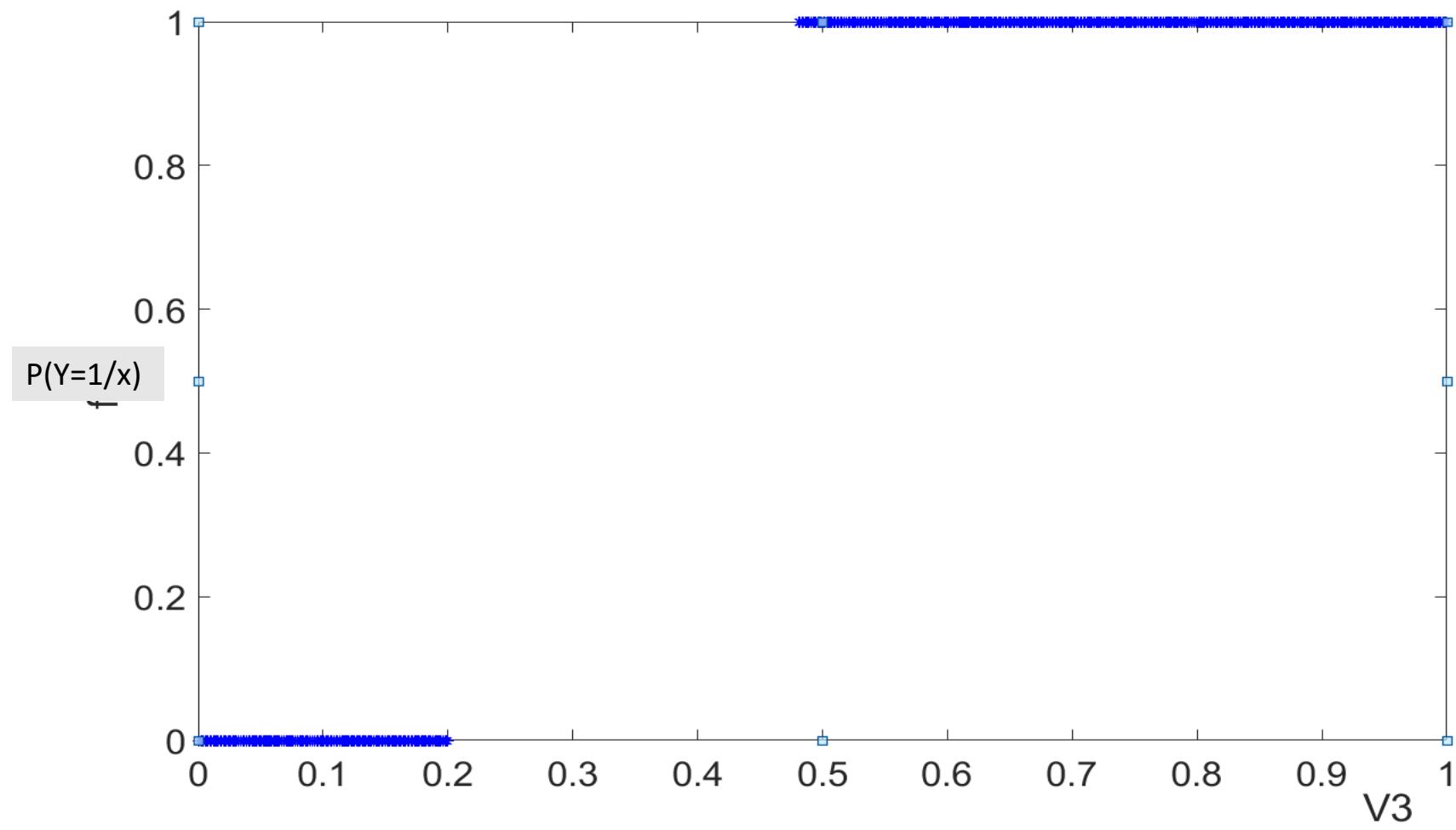


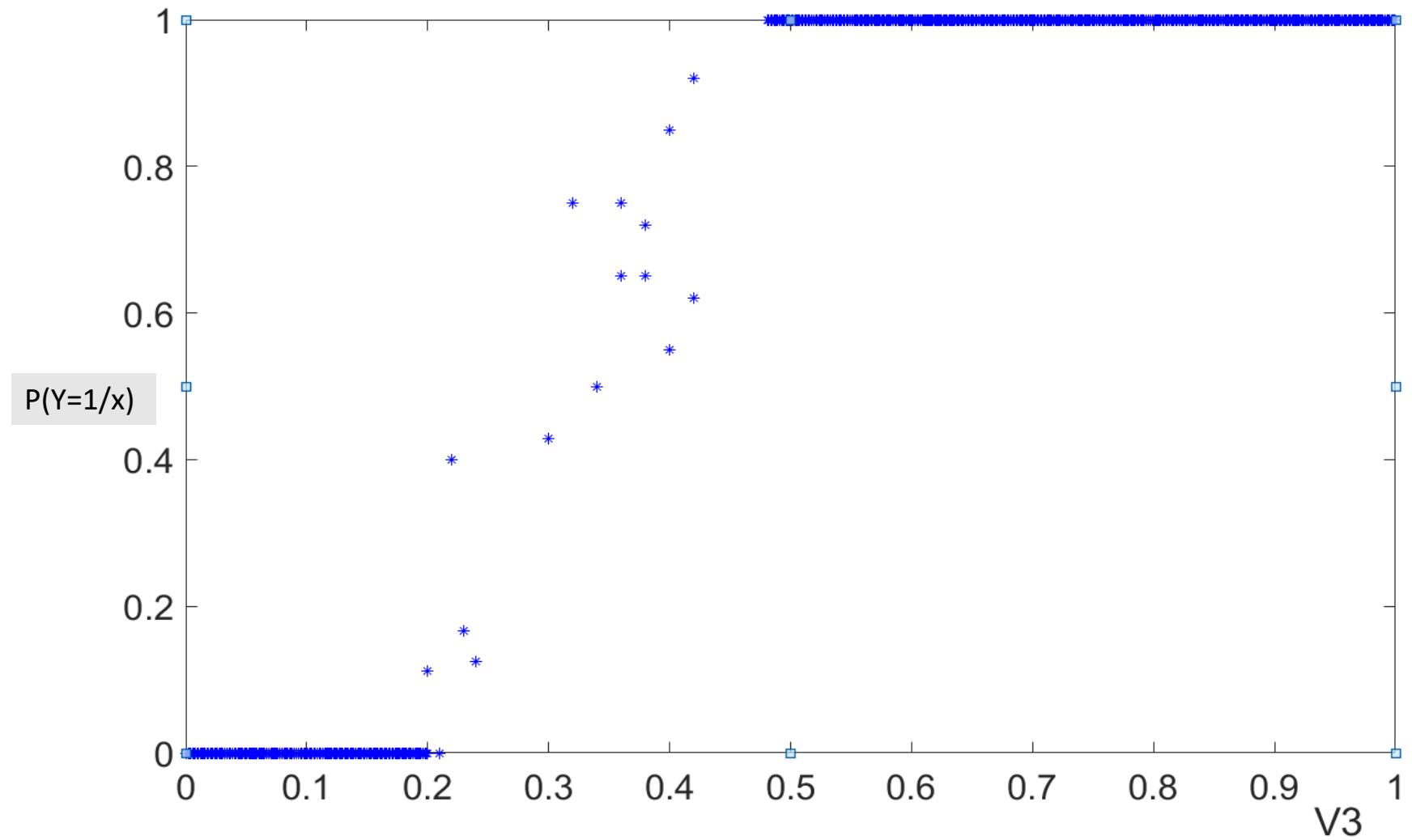
Fraudulent Transactions



Legitimate Transactions







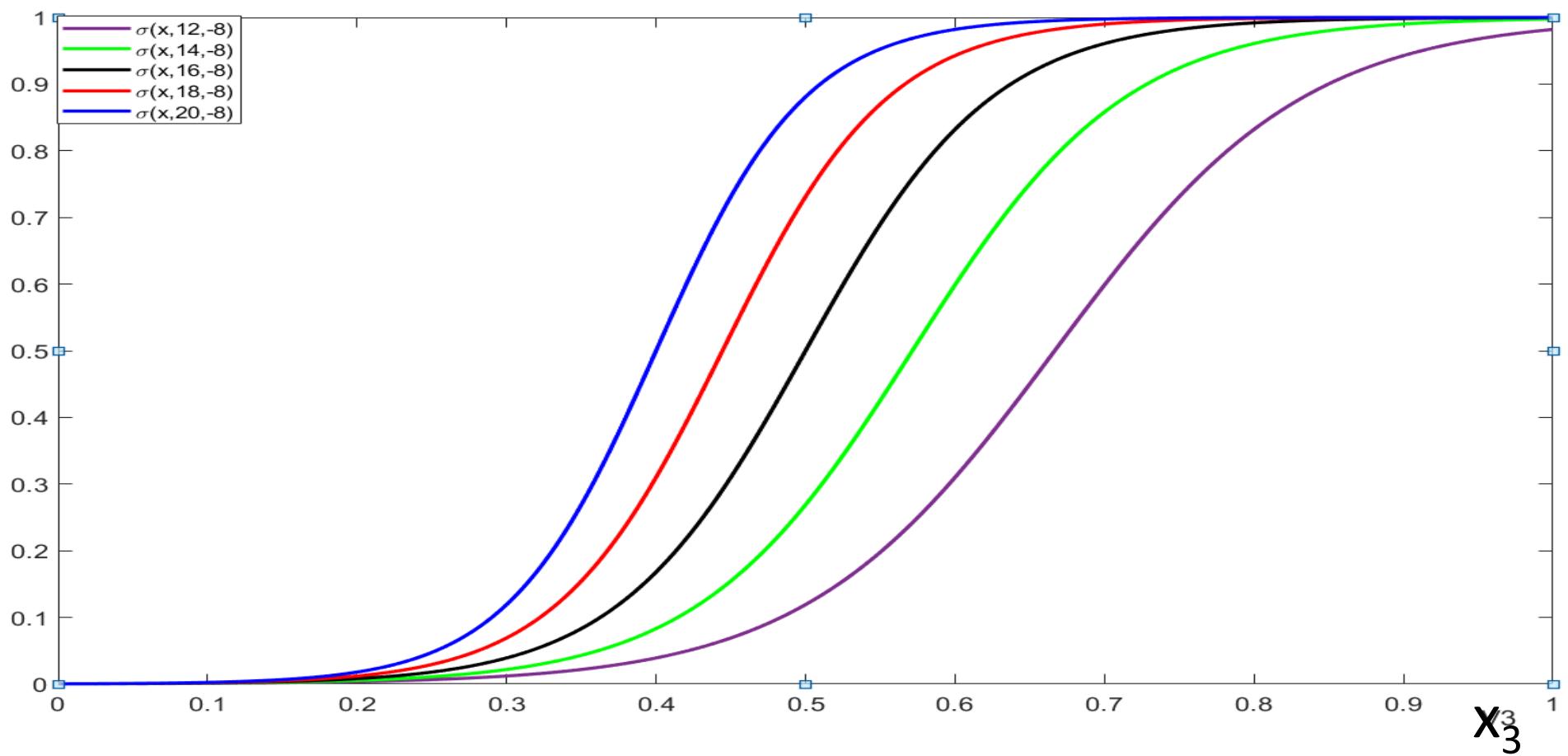
Sigmoidal Function

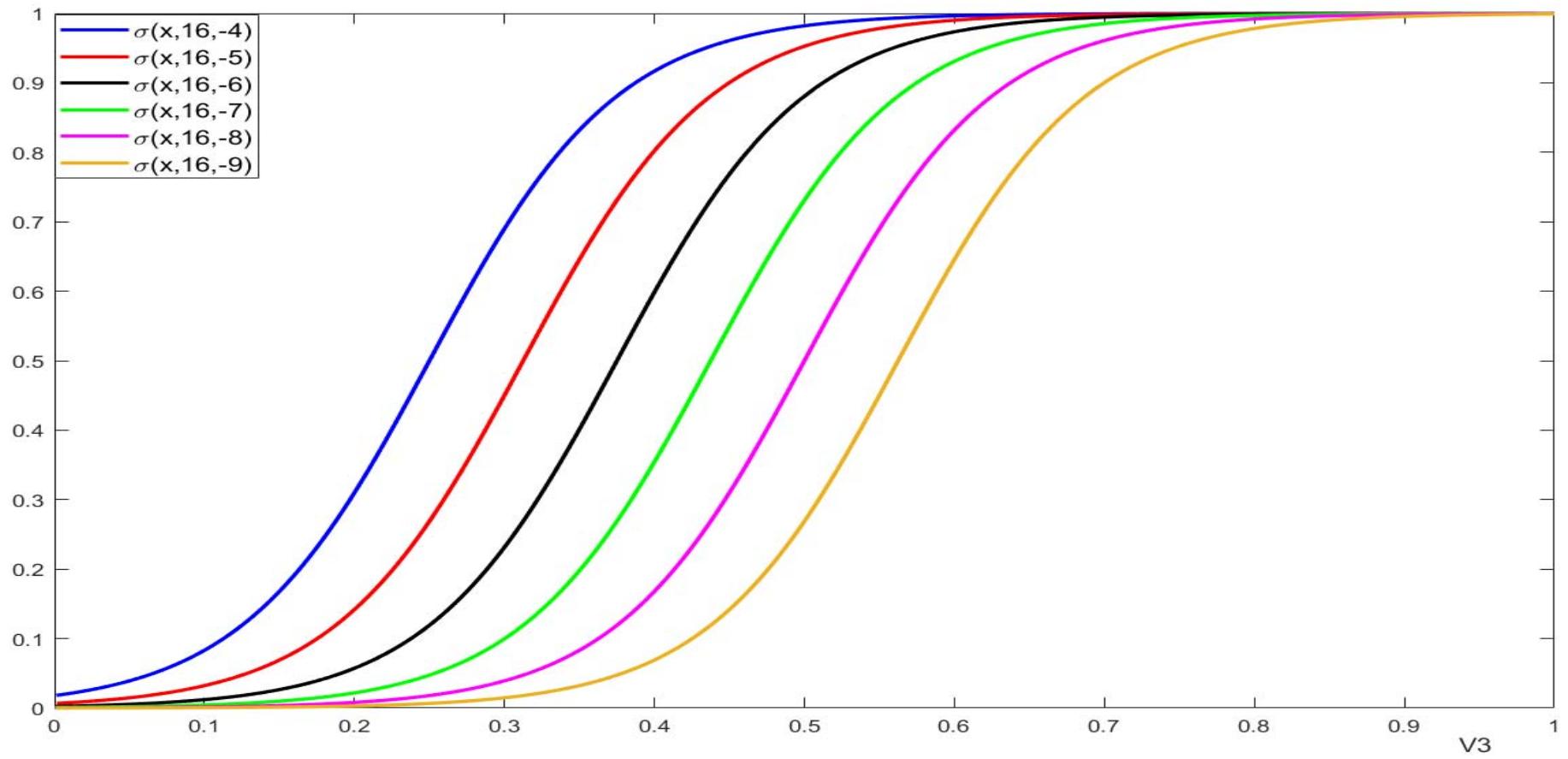
$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$

$$\hat{\sigma}(x) = \underline{\sigma(x)(1 - \sigma(x))}$$

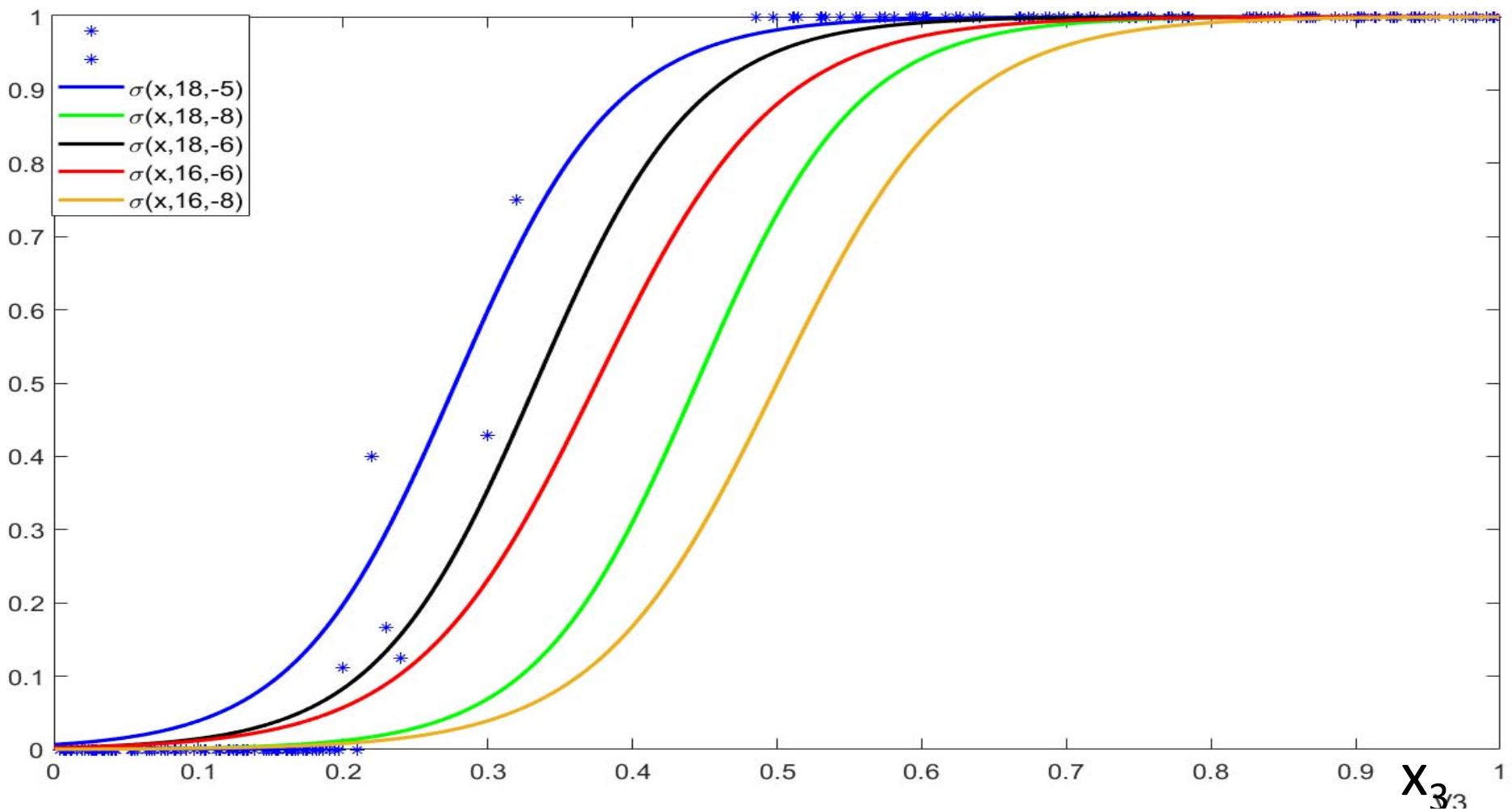
$$\frac{1}{1 + e^{-(\beta_1 x + \beta_0)}} = \frac{e^{(\beta_1 x + \beta_0)}}{1 + e^{(\beta_1 x + \beta_0)}}$$

<https://www.desmos.com/calculator/coknirwubg>

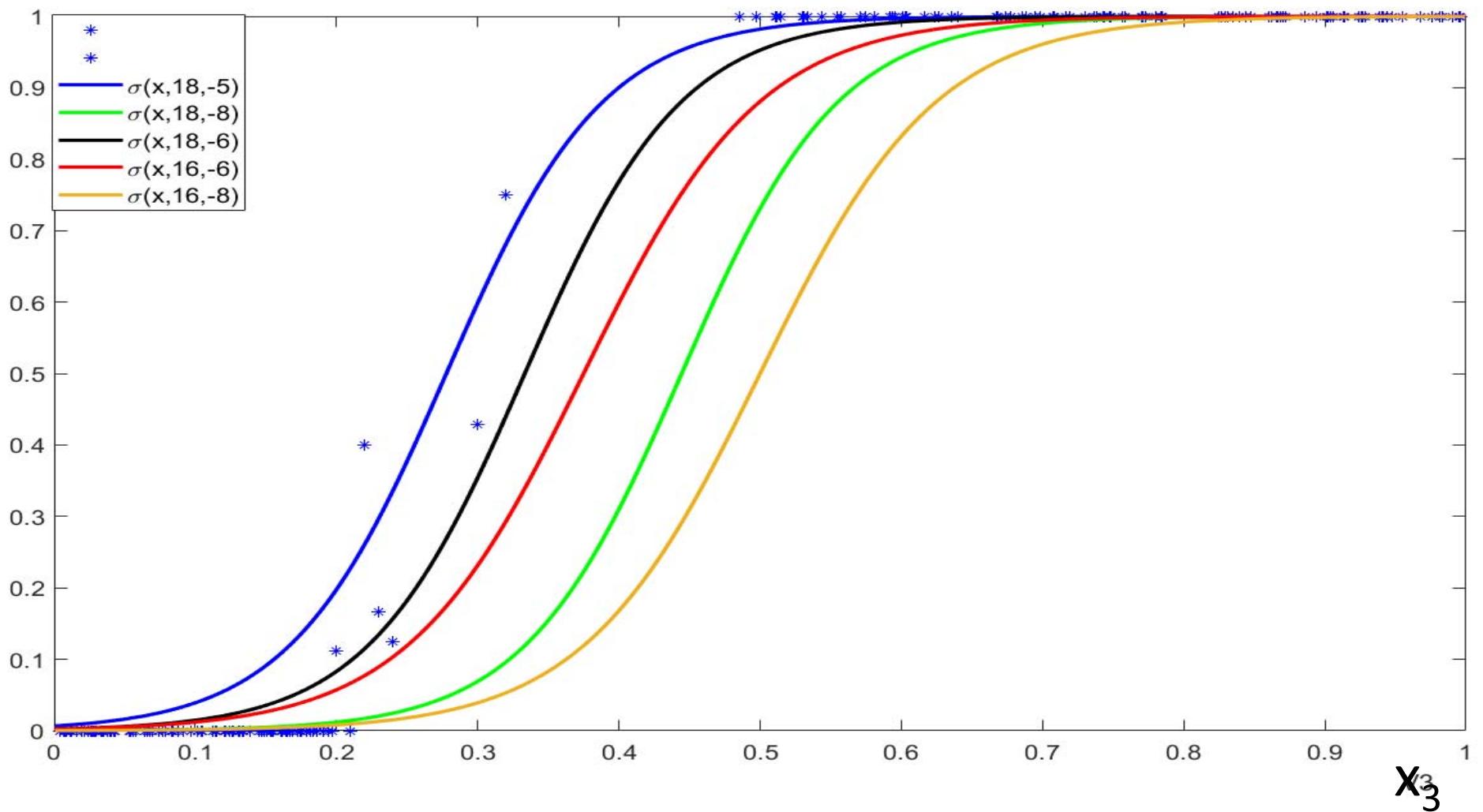




Logistic Function fitting



$\rho(x_1/x)$

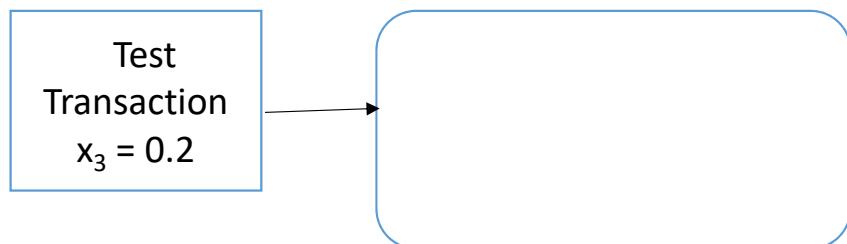


Logistic Regression on Test data

$$f(x) = \frac{p}{1 + e^{-(\beta_1 x + \beta_0)}}$$

Estimated

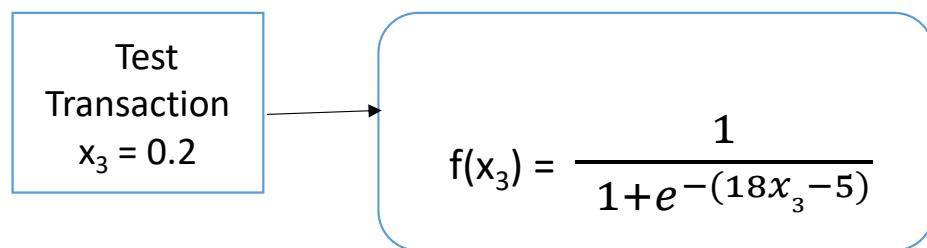
$$\begin{aligned}\beta_1 &= 18 \\ \beta_0 &= -5\end{aligned}$$



Logistic Regression on Test data

Estimated

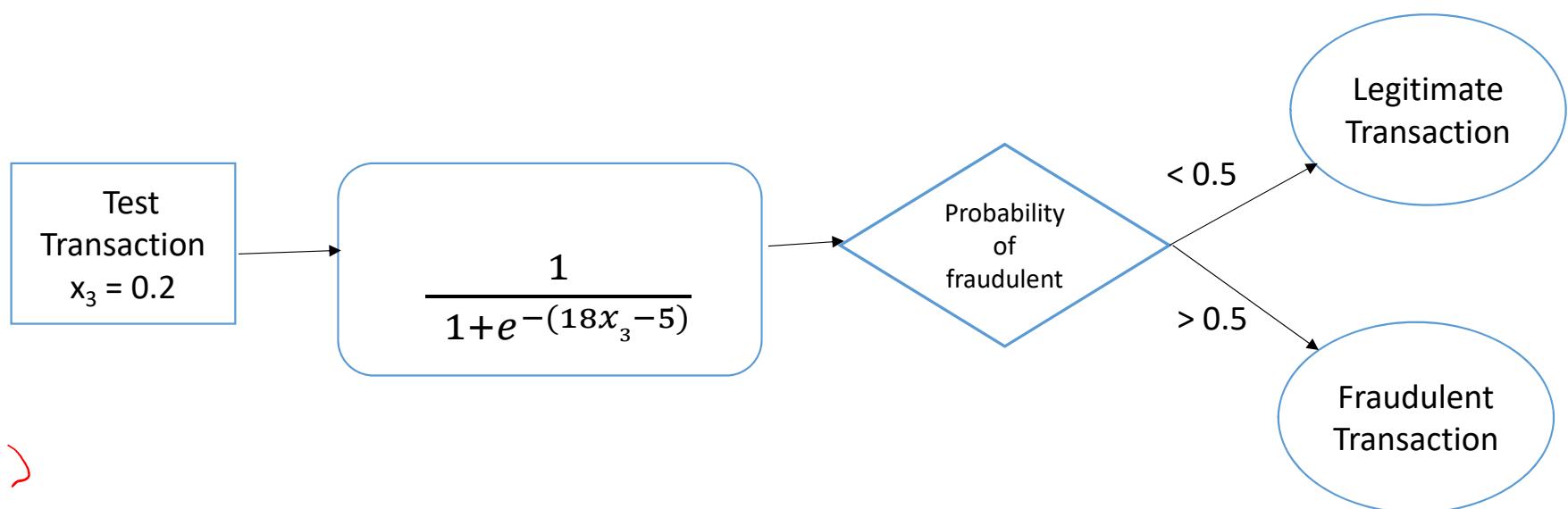
$$\begin{aligned}\beta_1 &= 18 \\ \beta_0 &= -5\end{aligned}$$



Logistic Regression on Test data

Estimated

$$\begin{aligned}\beta_1 &= 18 \\ \beta_0 &= -5\end{aligned}$$

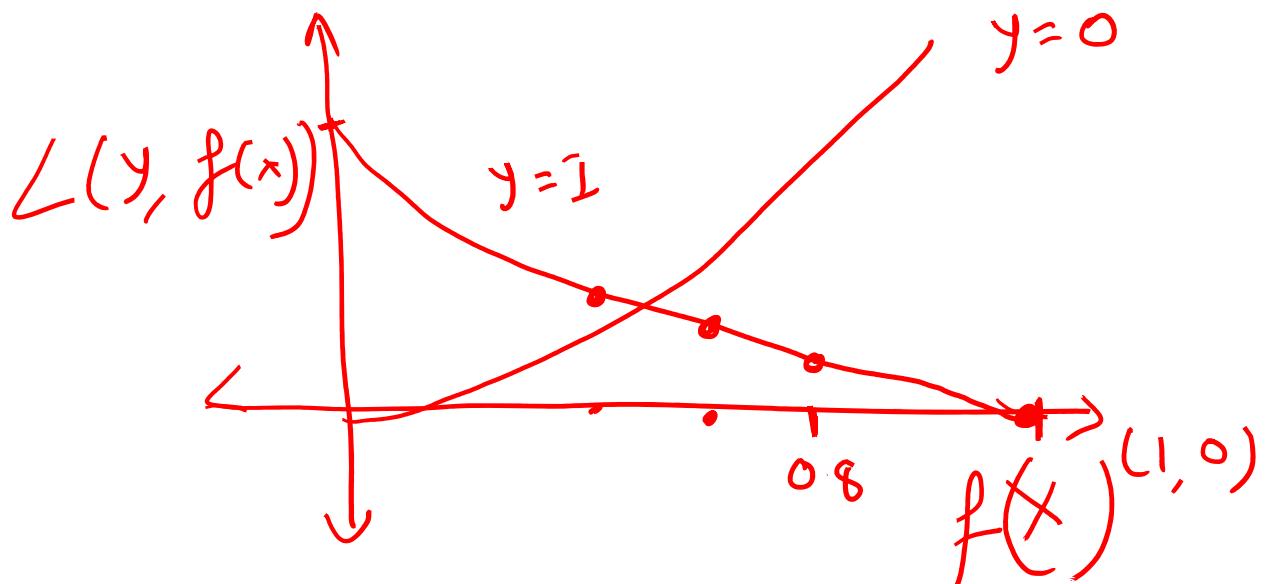


Logistic Regression and Cross entropy Loss

Gradient
Transcation

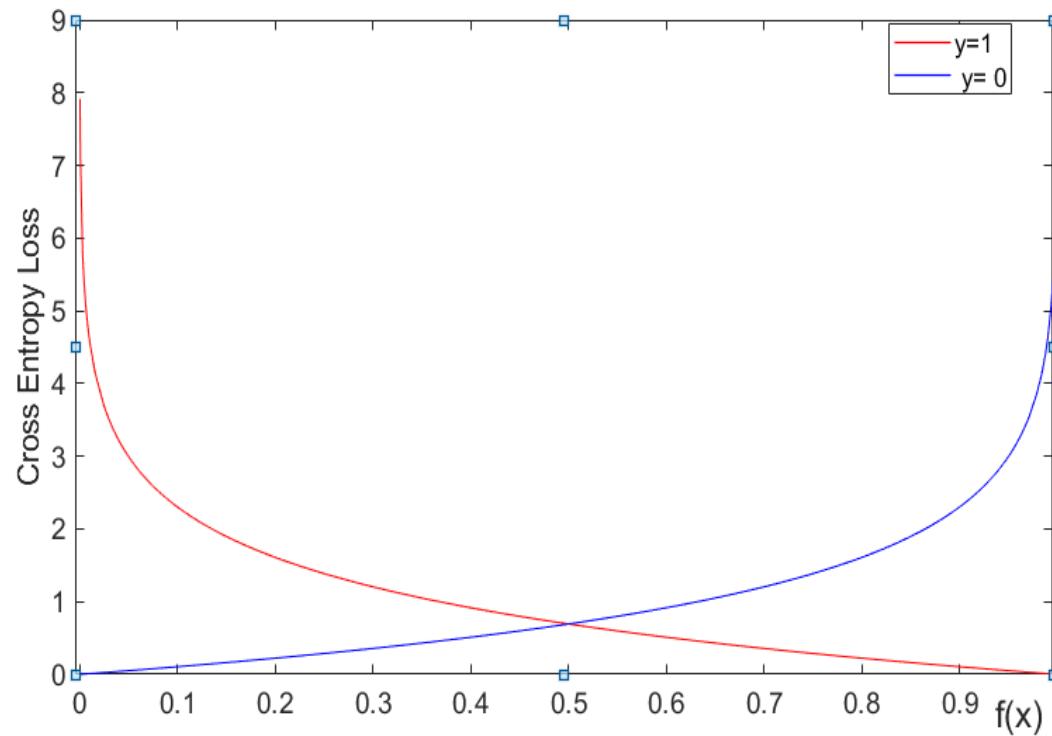
$$L(y, f(x)) = \begin{cases} -\log(f(x)) & \text{if } y = 1 \\ -\log(1 - f(x)) & \text{if } y = 0 \end{cases}$$

$$f(x) = \underline{P(y=1/x)}$$



Logistic Regression and Cross entropy Loss

$$L(y, f(x)) = \begin{cases} -\log(f(x)) & \text{if } y = 1 \\ -\log(1 - f(x)) & \text{if } y = 0 \end{cases}$$



Logistic Regression and Cross entropy Loss

$$L(y, f(x)) = \begin{cases} -\log(f(x)) & \text{if } y = 1 \\ -\log(1 - f(x)) & \text{if } y = 0 \end{cases}$$
$$= -y \log(f(x)) - (1 - y) \log(1 - f(x))$$

$$P(y=1/x) = \frac{1}{1 + e^{-(\beta_1 x + \beta_0)}}$$

Logistic Regression and Cross entropy Loss

$$L(y, f(x)) = \begin{cases} -\log(f(x)) & \text{if } y = 1 \\ -\log(1 - f(x)) & \text{if } y = 0 \end{cases}$$
$$= -y \log(f(x)) - (1 - y) \log(1 - f(x))$$

$$\underset{f}{\text{Min}} \quad \sum_{i=1}^l -y_i \log(f(x_i)) - (1 - y_i) \log(1 - f(x_i))$$

Logistic Regression and Cross entropy Loss

$$f(x) = \frac{1}{1 + e^{-(\beta_1 x + \beta_0)}} = \frac{e^{\beta_1 x + \beta_0}}{1 + e^{\beta_1 x + \beta_0}}$$

$$\underset{f}{\text{Min}} \quad \sum_{i=1}^l -y_i \log(f(x_i)) - (1 - y_i) \log(1 - f(x_i))$$

$$\underset{(\beta_1, \beta_0)}{\text{Min}} \quad \sum_{i=1}^l -y_i \log \left(\frac{1}{1 + e^{-(\beta_1 x_i + \beta_0)}} \right) - (1 - y_i) \log \left(1 - \frac{1}{1 + e^{-(\beta_1 x_i + \beta_0)}} \right)$$


$$\sum_{i=1}^l \left[\ln(1 + e^{B_0 + B_1 x_i}) - y_i(B_0 + B_1 x_i) \right]$$

$$\rightarrow \sum_{i=1}^l \left[-y_i \ln \left(\frac{e^{z_i}}{1+e^{z_i}} \right) - (1-y_i) \ln \left(\frac{1}{1+e^{z_i}} \right) \right] \quad (z_i = B_0 + B_1 x_i)$$

$$= \sum_{i=1}^l \left[-y_i \ln(e^{z_i}) + y_i \cancel{\ln(1+e^{z_i})} + \ln(1+e^{z_i}) - \cancel{y_i \ln(1+e^{z_i})} \right]$$

$$= \sum_{i=1}^L (\ln(1 + e^{B_0 + B_1 x_i}) - y_i(B_0 + B_1 x_i))$$

Logistic Regression and Cross entropy Loss

$$\frac{1}{1+e^{-(\beta_1 x + \beta_0)}} = \frac{e^{\beta_1 x + \beta_0}}{1+e^{\beta_1 x + \beta_0}}$$

$$\begin{aligned}
 & \underset{f}{\text{Min}} \quad \sum_{i=1}^l -y_i \log(f(x_i)) - (1 - y_i) \log(1 - f(x_i)) \\
 & \underset{(\beta_1, \beta_0)}{\text{Min}} \left| \sum_{i=1}^l -y_i \log \left(\frac{1}{1 + e^{-(\beta_1 x_i + \beta_0)}} \right) - (1 - y_i) \log \left(1 - \frac{1}{1 + e^{-(\beta_1 x_i + \beta_0)}} \right) \right| \\
 & = \underset{(\beta_1, \beta_0)}{\text{Min}} - \sum_{i=1}^l \left(y_i (\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right) \\
 & \quad \checkmark
 \end{aligned}$$

$$\begin{aligned}
& -y_i \log \left(\frac{e^{\beta_1 x + \beta_0}}{1 + e^{\beta_1 x + \beta_0}} \right) - (1-y_i) \log \left(\frac{e^{-(\beta_1 x + \beta_0)}}{1 + e^{-(\beta_1 x + \beta_0)}} \right) \\
&= -y_i \log \left((\beta_1 x + \beta_0) - \log \left(1 + e^{(\beta_1 x + \beta_0)} \right) \right) \\
&\quad - (1-y_i) \left(\log \left(e^{-\beta_1 x + \beta_0} \right) - \log \left(1 + e^{-\beta_1 x + \beta_0} \right) \right) \\
&= -y_i (\beta_1 x + \beta_0) + y_i \log \left(1 + e^{(\beta_1 x + \beta_0)} \right) \\
&\quad + (1-y_i) (\beta_1 x + \beta_0) + (1-y_i) \log \left(\frac{1 + e^{\beta_1 x + \beta_0}}{e^{\beta_1 x + \beta_0}} \right)
\end{aligned}$$

$$-\sum (-y_i(\beta_1 x_i + \beta_0) - \log(1 + e^z))$$

$$\begin{bmatrix} \frac{\partial}{\partial \beta_1} \\ \frac{\partial}{\partial \beta_0} \end{bmatrix}$$

$$= \sum_{i=1}^N \left[y_i x_i - \frac{1}{1+e^z} \cdot e^z \cdot x_i \right]$$

$$= \sum_{i=1}^N \left[x_i \left(y_i - \frac{1}{1+e^z} \cdot e^z \right) \right]$$

$$= \sum_{i=1}^N \left[x_i \left(y_i - \frac{e^{\beta_1 x_i + \beta_0}}{1+e^{\beta_1 x_i + \beta_0}} \right) \right]$$

$$= \sum_{i=1}^N \left[(y_i - \sigma(x_i, \beta_1, \beta_0)) x_i \right]$$

$$\begin{aligned}
&= y_i(\beta_1 x + \beta_0) + y_i \log(1 + e^{\beta_1 x + \beta_0}) \\
&\quad + (1-y_i)(\beta_1 x + \beta_0) - (1-y_i)(\beta_1 x + \beta_0) + \\
&\quad (1-y_i) \log(1 + e^{\beta_1 x + \beta_0}) \\
&= -(y_i(\beta_1 x + \beta_0) - \log(1 + e^{\beta_1 x + \beta_0}))
\end{aligned}$$

Logistic Regression and Cross entropy Loss

$$\underset{(\beta_1, \beta_0)}{\text{Min}} - \sum_{i=1}^l \left(y_i(\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right)$$

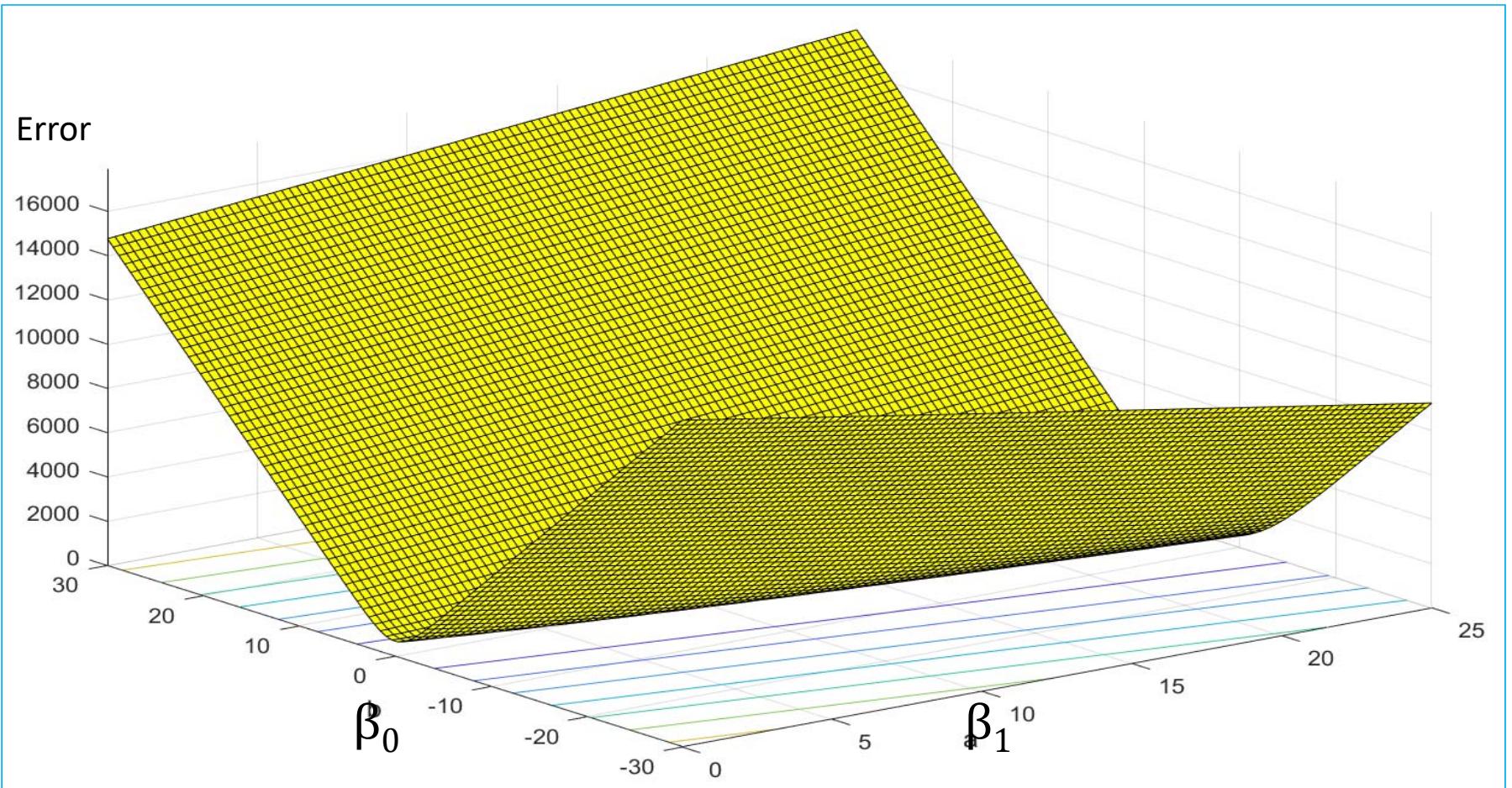
- Only minimization of $- \sum_{i=1}^l \left(y_i(\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right)$ is not enough.
- Our estimated logistic function should not increase sharply as well.

Logistic Regression Optimization Problem

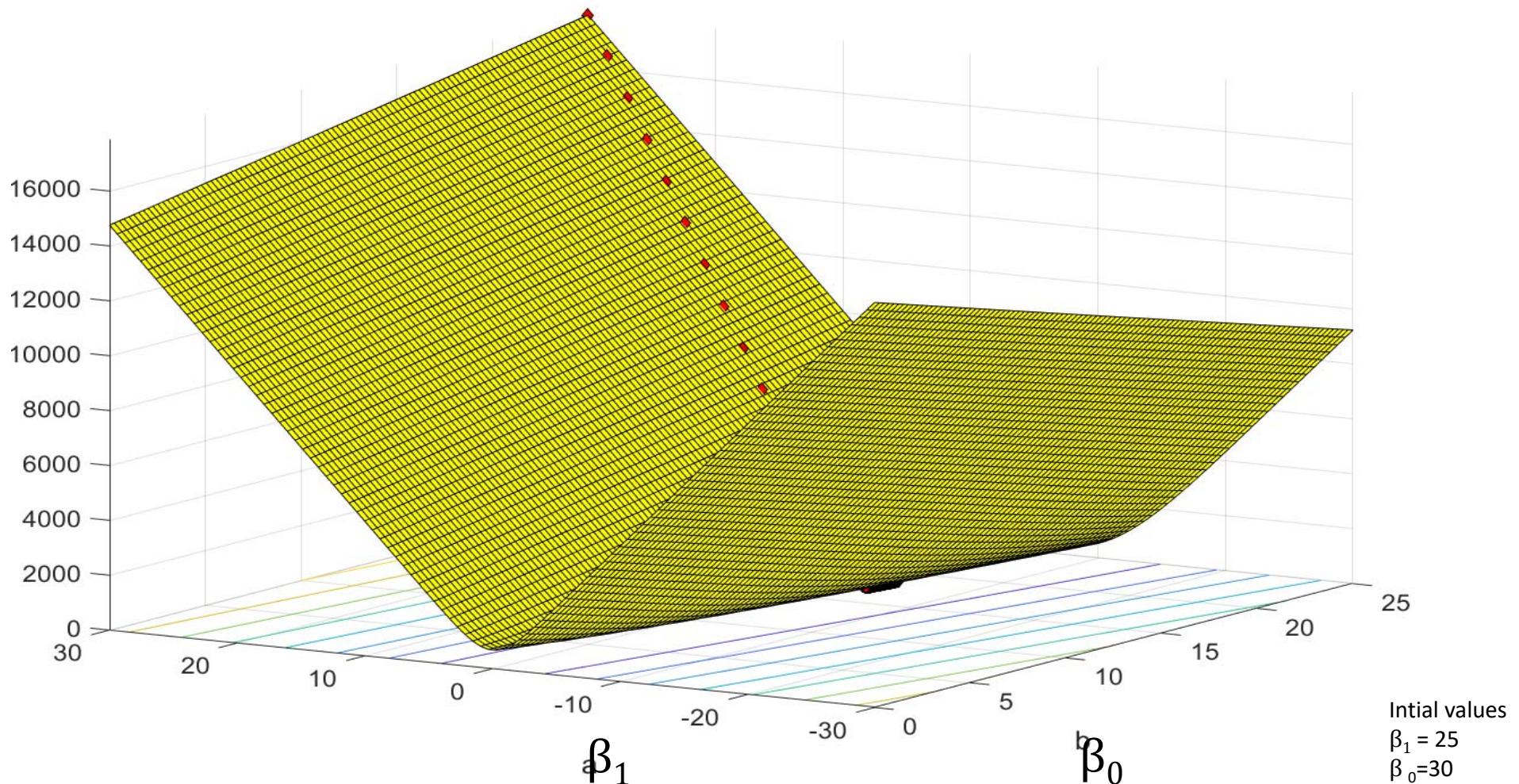
$$\underset{(\beta_1, \beta_0)}{\text{Min}} J(\beta_1, \beta_0) = -\sum_{i=1}^l \left(y_i (\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta_1^2$$

Empirical Loss

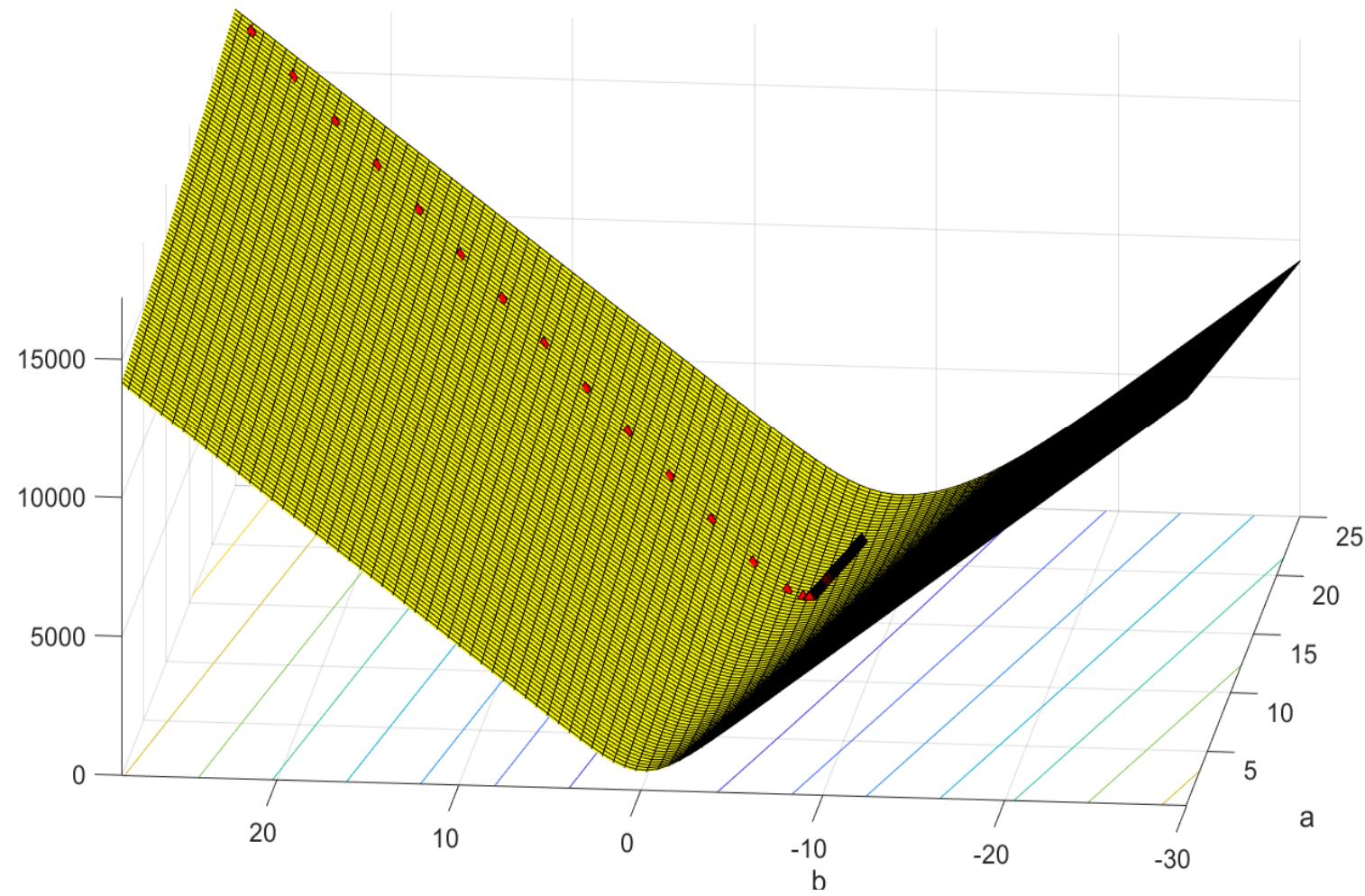
Regularization with
parameter λ



Gradient Descent Method



Gradient Descent Method



Computing Gradients

$$J(\beta_1, \beta_0) = -\sum_{i=1}^l \left(y_i (\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta_1^2$$

Computing Gradients

$$J(\beta_1, \beta_0) = - \sum_{i=1}^l \left(y_i (\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta_1^2$$



$$\nabla_{\beta_1} J(\beta_1, \beta_0) = \lambda \beta_1 - \sum_{i=1}^n (y_i - \left(\frac{1}{1 + e^{-(\beta_1 x_i + \beta_0)}} \right)) x_i$$



Computing Gradients

$$J(\beta_1, \beta_0) = - \sum_{i=1}^l \left(y_i (\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta_1^2$$

$$\begin{aligned} \nabla_{\beta_1} J(\beta_1, \beta_0) &= \lambda \beta_1 - \sum_{i=1}^n (y_i - \underbrace{\left(\frac{1}{1+e^{-(\beta_1 x_i + \beta_0)}} \right)}_{\sigma(x, \beta_1, \beta_0)}) x_i \\ &= \lambda \beta_1 - \sum_{i=1}^n (y_i - \underbrace{\sigma(x, \beta_1, \beta_0)}) x_i \end{aligned}$$

Computing Gradients

$$\beta_1^T x_i + \beta_0$$

$$J(\beta_1, \beta_0) = - \sum_{i=1}^l \left(y_i (\beta_1 x_i + \beta_0) - \log(1 + e^{(\beta_1 x_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta_1^2$$

$$\nabla_{\beta_1} J(\beta_1, \beta_0) = \lambda \beta_1 - \sum_{i=1}^n (y_i - \underbrace{\left(\frac{1}{1+e^{-(\beta_1 x_i + \beta_0)}} \right)}_{\sigma(x, \beta_1, \beta_0)}) x_i$$

$$\frac{1}{1 + e^{(\beta_1^T x + \beta_0)}}$$

$$\nabla_{\beta_0} J(\beta_1, \beta_0) = - \sum_{i=1}^n (y_i - \underbrace{\left(\frac{1}{1+e^{-(\beta_1 x_i + \beta_0)}} \right)}_{\sigma(x, \beta_1, \beta_0)})$$

$$= - \sum_{i=1}^n (y_i - \sigma(x, a^{(k)}, b^{(k)}))$$

Gradient Descent Algorithm for Logistic Regression

Algorithm:- Gradient descent method

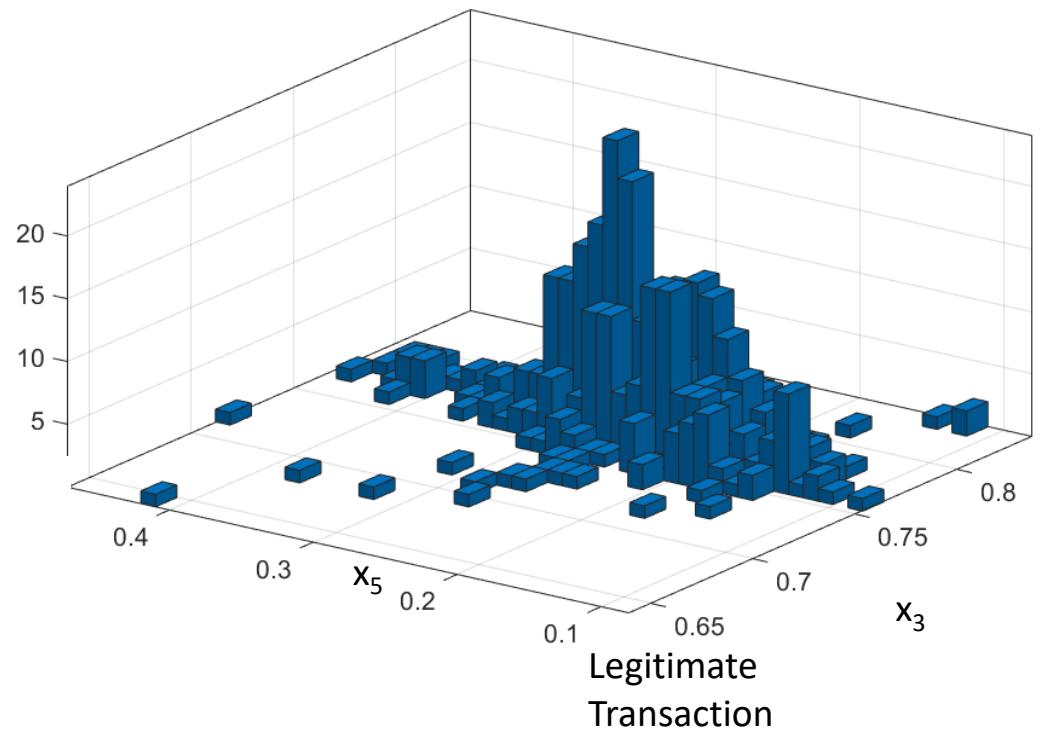
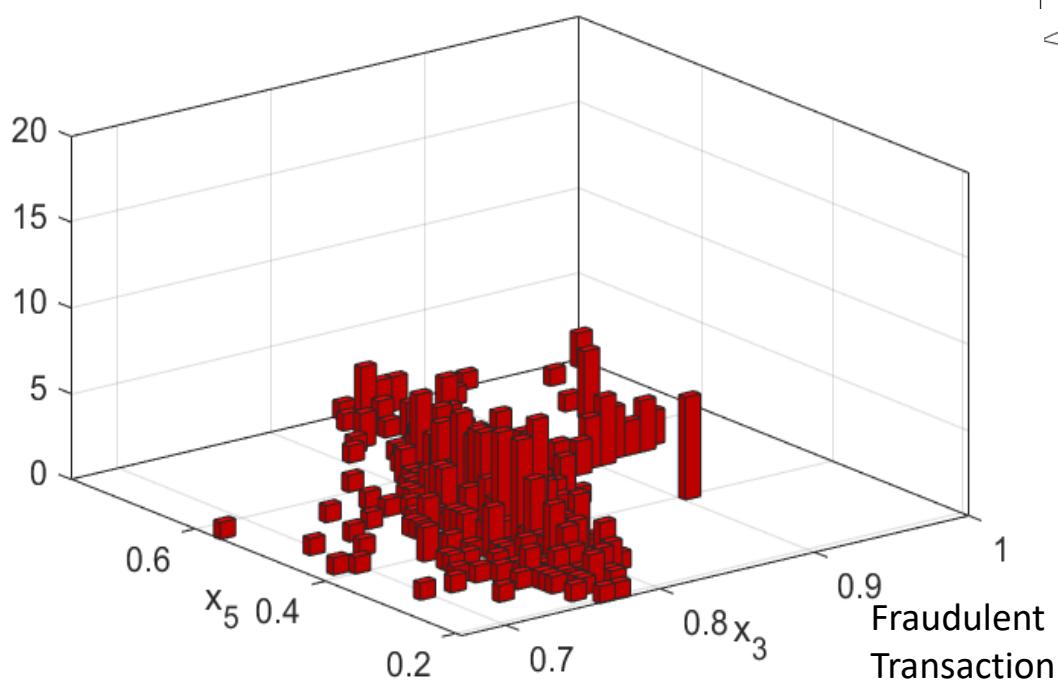
Initialize β_1^0 and β_0^0

Repeat

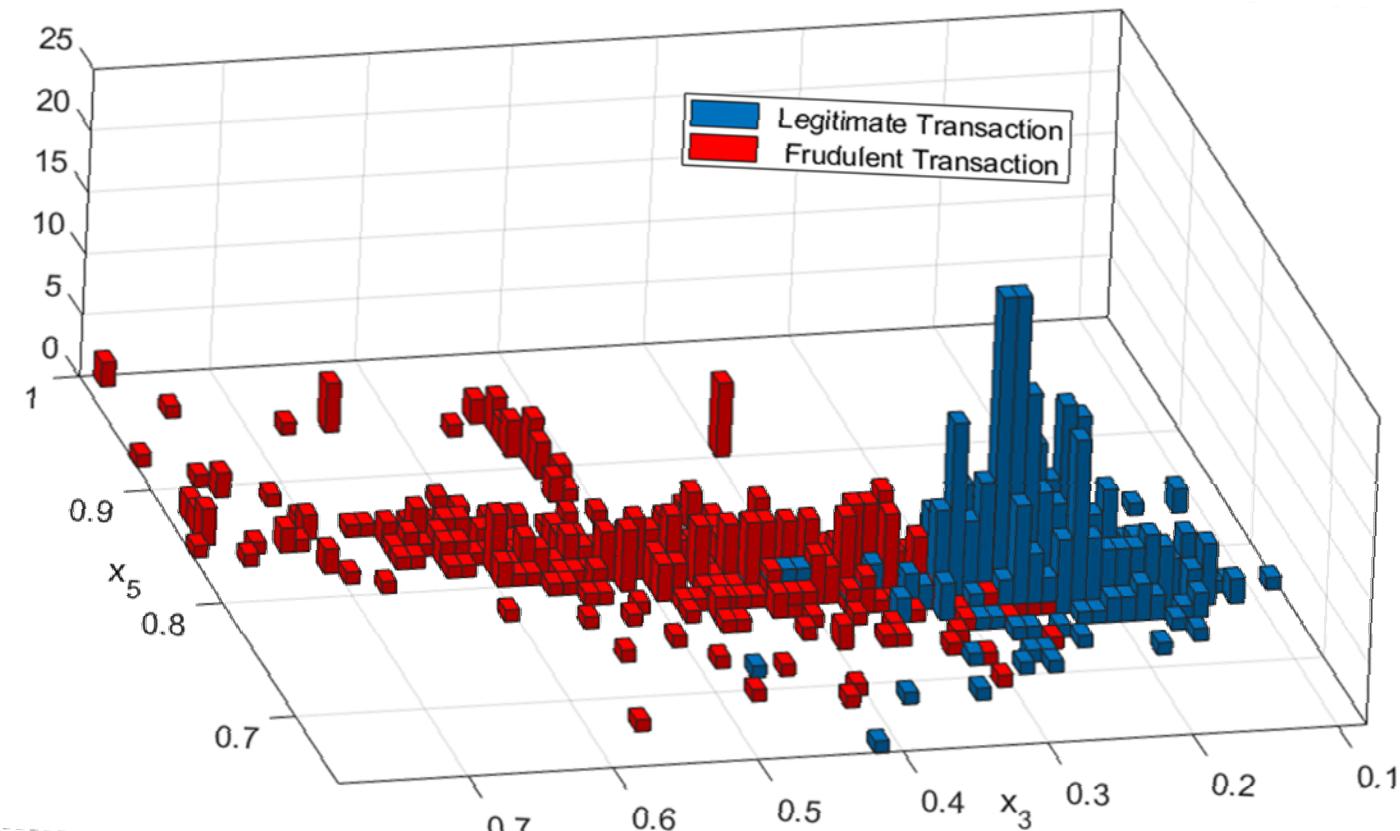
$$\begin{aligned}\beta_1^{(k+1)} &:= \beta_1^{(k)} - \eta \nabla_{\beta_1} J(\beta_1^{(k)}, \beta_0^{(k)}) \\ \beta_0^{(k+1)} &:= \beta_0^{(k)} - \eta \nabla_{\beta_0} J(\beta_1^{(k)}, \beta_0^{(k)})\end{aligned}$$

Until $\left| \left| \begin{bmatrix} \nabla_{\beta_1} J(\beta_1^{(k+1)}, \beta_0^{(k+1)}) \\ \nabla_{\beta_0} J(\beta_1^{(k+1)}, \beta_0^{(k+1)}) \end{bmatrix} \right| \right| \leq \varepsilon$

Working with bivariate data



Working with bivariate data

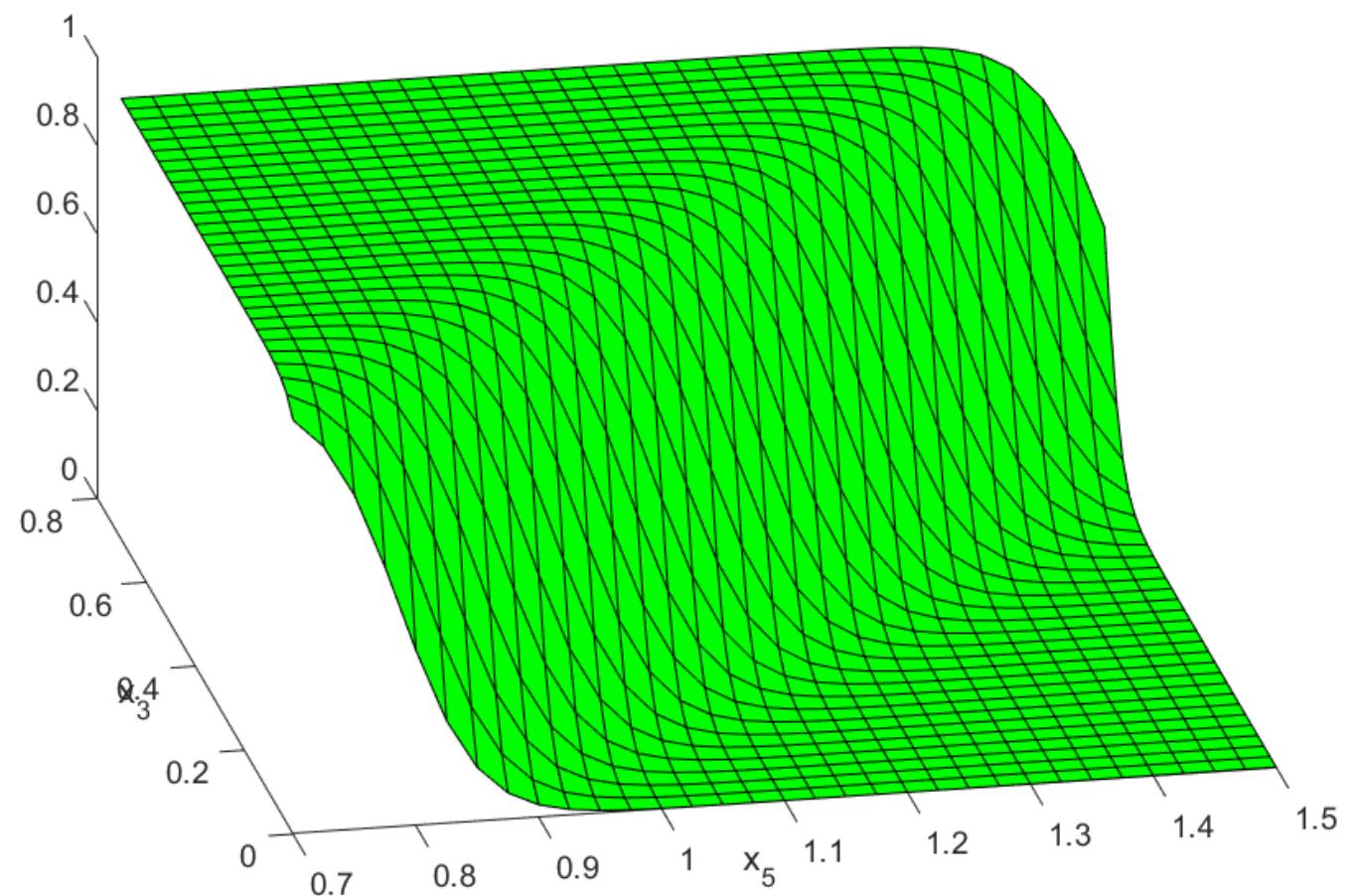


Logistic Regression for bivariate data

$$\frac{1}{1 + e^{-(\beta_1 x_1 + \beta_2 x_2 + \beta_0)}}$$

$$\frac{1}{1 + e^{-(\beta_1 x_1 + \beta_2 x_2 + \beta_0)}} = \frac{e^{(\beta_1 x_1 + \beta_2 x_2 + \beta_0)}}{1 + e^{(\beta_1 x_1 + \beta_2 x_2 + \beta_0)}}$$

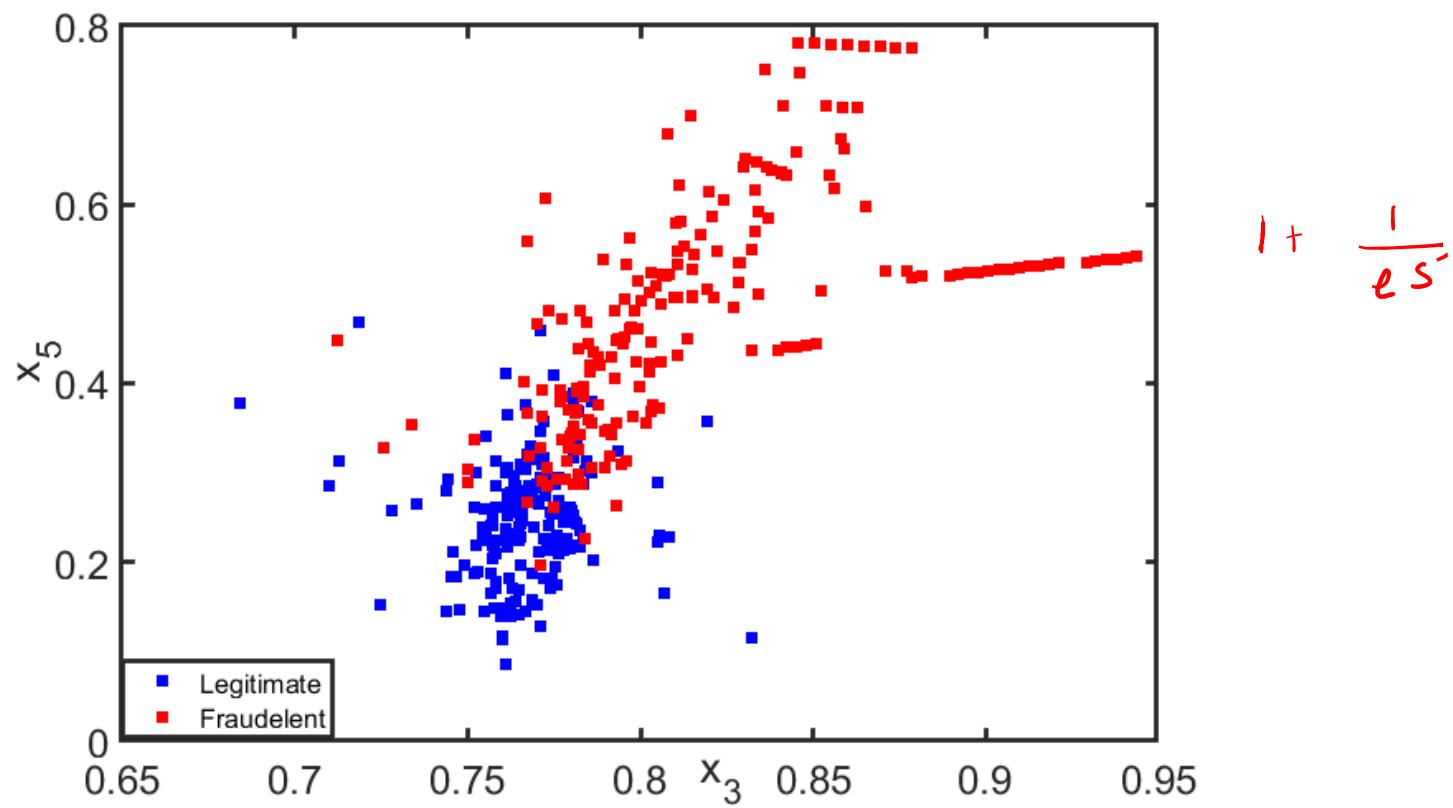
<https://www.desmos.com/calculator/coknirwubg>

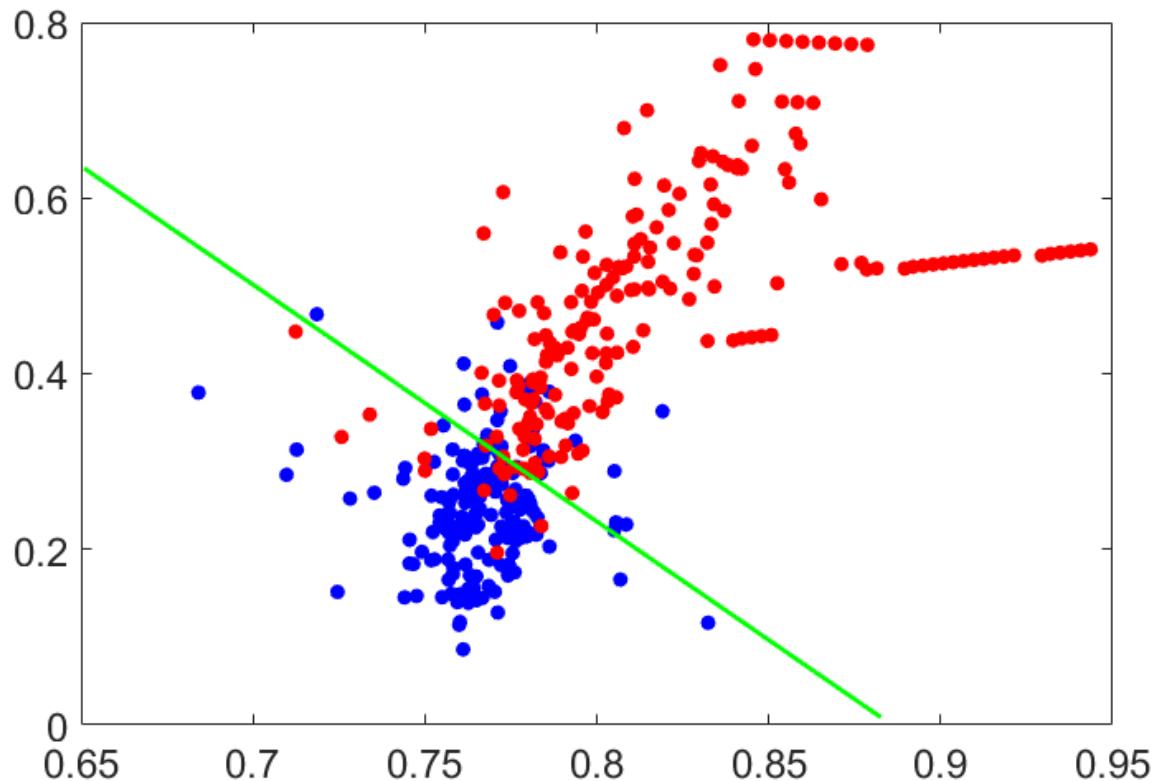


One more Observation

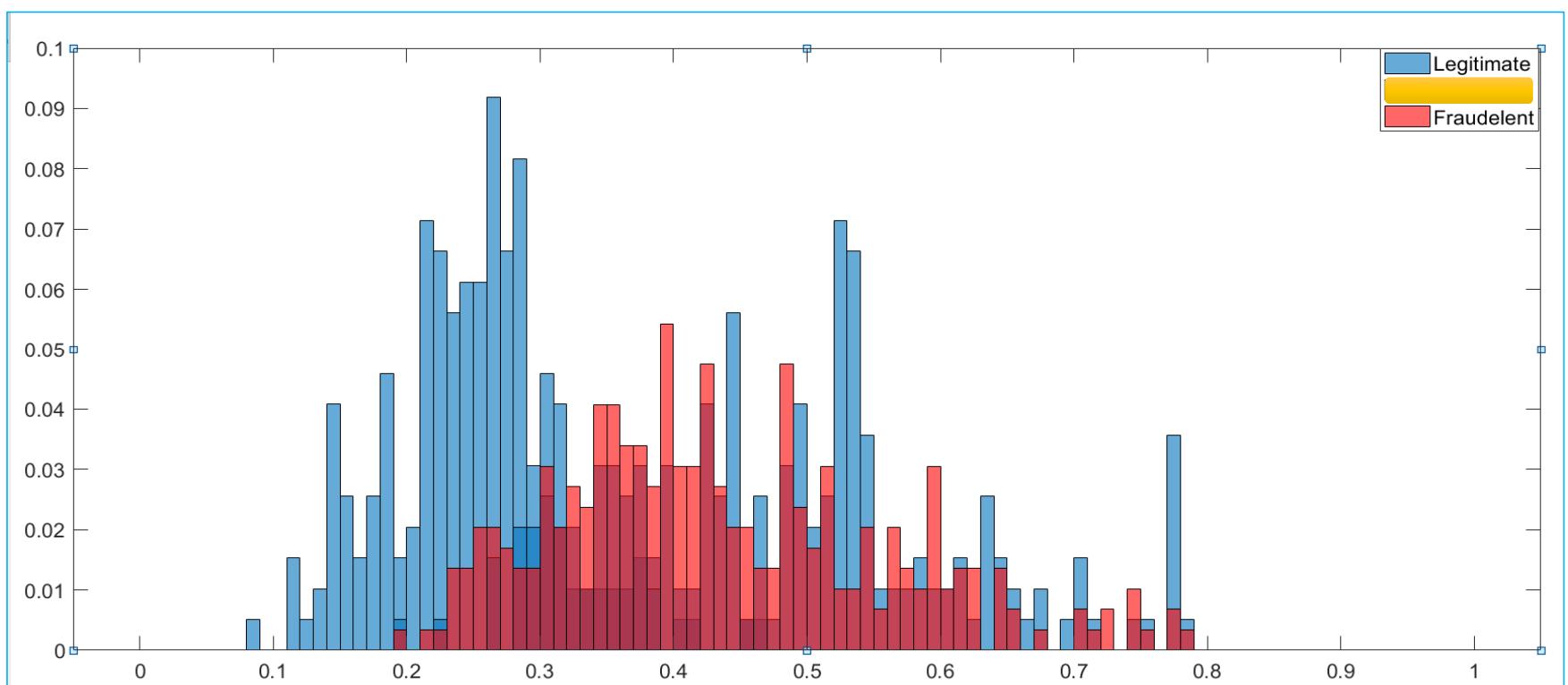
$$\frac{1}{1+e^{-(\beta_1 x + \beta_0)}} = 0.5$$

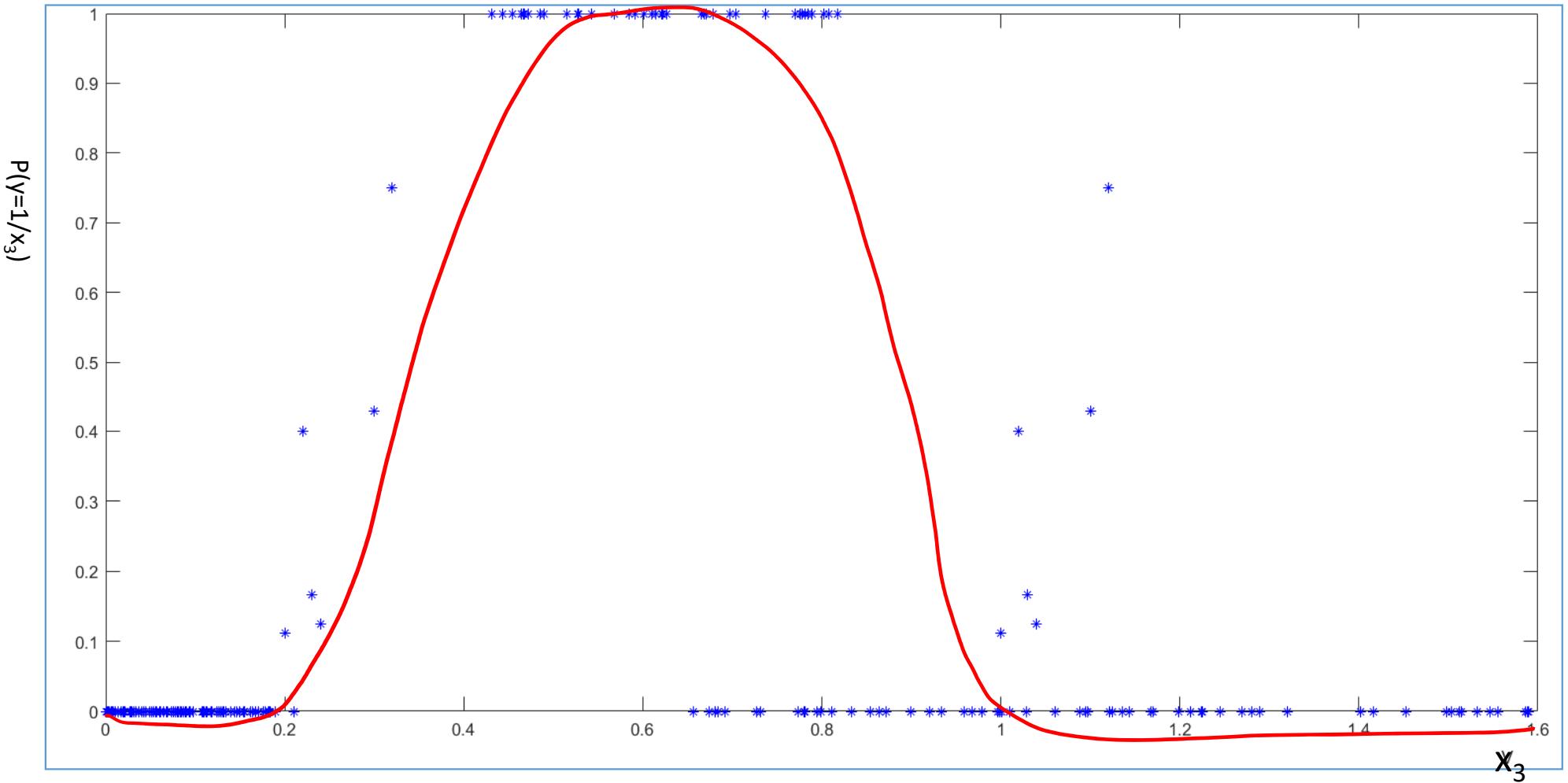
$$\beta_1 x + \beta_0 = 0$$





Need for non-linear logistic regression





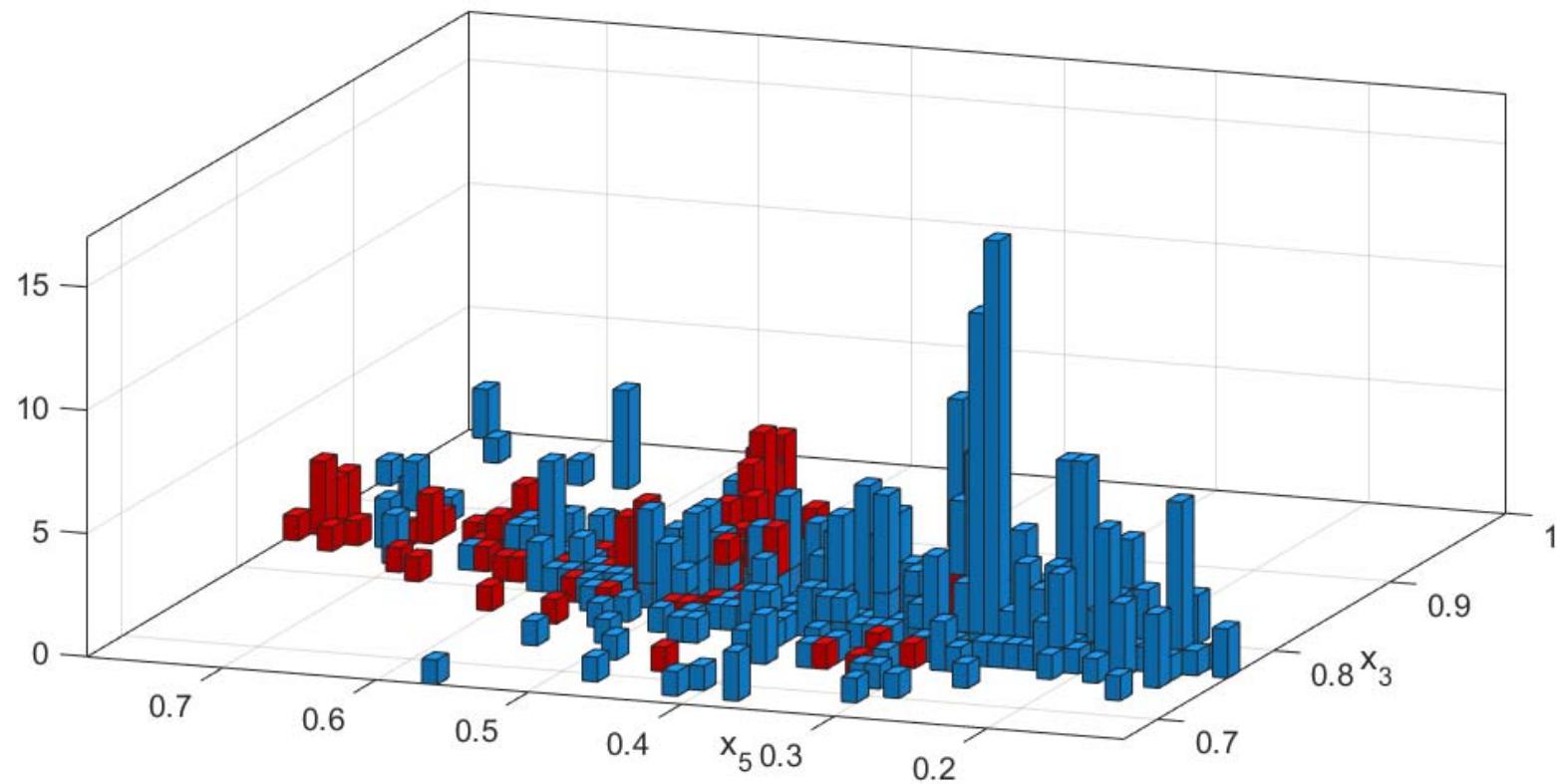
Non-linear Logistic Regression

$$\frac{1}{1 + e^{-(\beta_1 x + \beta_0)}}$$

$$f(x) = \frac{1}{1 + e^{-(\beta_2 x^2 + \beta_1 x + \beta_0)}} = \frac{e^{(\beta_2 x^2 + \beta_1 x + \beta_0)}}{1 + e^{(\beta_2 x^2 + \beta_1 x + \beta_0)}}$$

<https://www.desmos.com/calculator/coknirwubg>

Non-linear logistic Regression for bivariate data



$$\frac{e^{\beta_1 x_1 + \beta_2 x_2 + \beta_0}}{1 + e^{\beta_1 x_1 + \beta_2 x_2 + \beta_0}}$$

$$\frac{e^{\beta^T \Phi(x) + \beta_0}}{1 + e^{\beta^T \Phi(x) + \beta_0}}$$

$$\frac{e^{\beta^T x + \beta_0}}{1 + e^{\beta^T x + \beta_0}}$$

$$f(x) = \frac{1}{1 + e^{-(\beta_1 x_1^2 + \beta_2 x_2^2 + \beta_3 x_1 x_2 + \beta_4 x_1 + \beta_5 x_2 + \beta_6)}}$$

$$= \frac{e^{(\beta_1 x_1^2 + \beta_2 x_2^2 + \beta_3 x_1 x_2 + \beta_4 x_1 + \beta_5 x_2 + \beta_6)}}{1 + e^{(\beta_1 x_1^2 + \beta_2 x_2^2 + \beta_3 x_1 x_2 + \beta_4 x_1 + \beta_5 x_2 + \beta_6)}}$$

$$\Phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1 \\ x_2 \\ b \end{bmatrix}$$

<https://www.desmos.com/calculator/coknirwubg>

x_1, x_2, \dots, x_n, y

$x_i \in R^n$

$$\sigma(\beta_1, \beta_0, x) = \frac{1}{1 + e^{-(\beta_1^T x + \beta_0)}}$$

$$\sigma(\beta_1, \beta_0, x) = \frac{1}{1 + e^{-(\beta_1^T \Phi(x) + \beta_0)}}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

$$\beta_1 = \begin{bmatrix} \beta_{11} \\ \vdots \\ \vdots \\ \beta_{1m} \end{bmatrix}$$

$$\beta_0 \in R$$

Consider the NAND dataset. Which one of the following may be the solution of logistic regression model ?

$$(a) \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_0 \end{bmatrix} = \begin{bmatrix} -80.26 \\ -110.85 \\ 21 \end{bmatrix}$$

$$(b) \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_0 \end{bmatrix} = \begin{bmatrix} 80.26 \\ 110.85 \\ -152.61 \end{bmatrix}$$

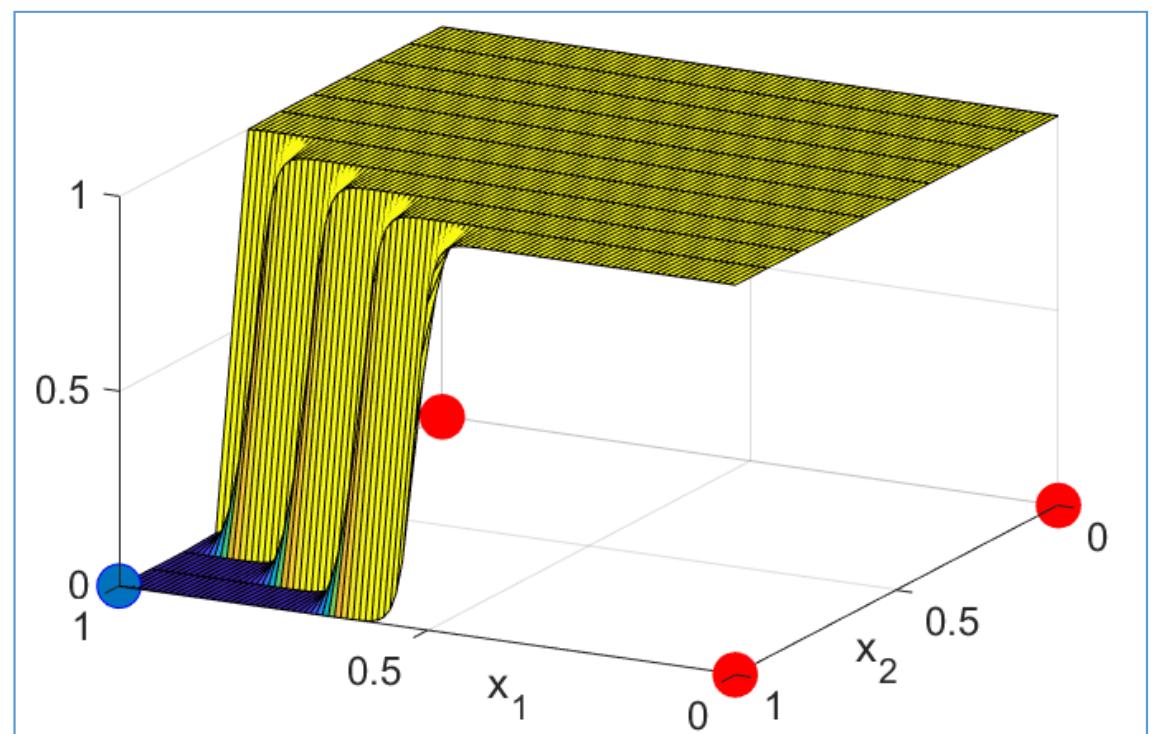
$$(c) \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_0 \end{bmatrix} = \begin{bmatrix} -80.26 \\ 0 \\ 152.61 \end{bmatrix}$$

$$(d) \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_0 \end{bmatrix} = \begin{bmatrix} -80.26 \\ -110.85 \\ 152.61 \end{bmatrix}$$

$$\beta_1 x_1 + \beta_2 x_2 + \beta_0 > 0 \quad \curvearrowleft 1$$

x1	x2	Y
0	0	1
0	1	1
1	0	1
1	1	0

NAND Dataset



Logistic Regression with n features

$$\frac{1}{1 + e^{-(\beta^T X_i + \beta_0)}}$$

$\beta \in \mathbb{R}^n$
 $\beta_0 \in \mathbb{R}$
 $X_i \in \mathbb{R}^n$

$$\underset{f}{\operatorname{Min}} \sum_{i=1}^l -y_i \log(f(X_i)) - (1 - y_i) \log(1 - f(X_i))$$

$$\underset{(\beta, \beta_0)}{\operatorname{Min}} \sum_{i=1}^l -y_i \log \left(\frac{1}{1 + e^{-(\beta^T X_i + \beta_0)}} \right) - (1 - y_i) \log \left(1 - \frac{1}{1 + e^{-(\beta^T X_i + \beta_0)}} \right)$$

$$= \underset{(\beta_1, \beta_0)}{\operatorname{Min}} - \sum_{i=1}^l \left(y_i (\beta^T X_i + \beta_0) - \log(1 + e^{(\beta^T X_i + \beta_0)}) \right)$$

Computing Gradients

$$J(\beta, \beta_0) = - \sum_{i=1}^l \left(y_i (\beta^T X_i + \beta_0) - \log(1 + e^{(\beta^T X_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta^T \beta$$

$$\nabla_{\beta} J(\beta, \beta_0) = \lambda \beta - \sum_{i=1}^n (y_i - \left(\frac{1}{1 + e^{-(\beta^T X_i + \beta_0)}} \right)) X_i$$

Computing Gradients

$$J(\beta, \beta_0) = - \sum_{i=1}^l \left(y_i (\beta^T X_i + \beta_0) - \log(1 + e^{(\beta^T X_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta^T \beta$$

$$\begin{aligned}\nabla_{\beta} J(\beta, \beta_0) &= \lambda \beta - \sum_{i=1}^n (y_i - \left(\frac{1}{1+e^{-(\beta^T X_i + \beta_0)}} \right)) X_i \\ &= \lambda \beta - \sum_{i=1}^n (y_i - \sigma(\cancel{x}, \beta, \beta_0)) X_i\end{aligned}$$

Computing Gradients

$$J(\beta, \beta_0) = - \sum_{i=1}^l \left(y_i (\beta^T X_i + \beta_0) - \log(1 + e^{(\beta^T X_i + \beta_0)}) \right) + \frac{\lambda}{2} \beta^T \beta$$

$$\begin{aligned} \nabla_{\beta} J(\beta, \beta_0) &= \lambda \beta - \sum_{i=1}^n (y_i - \left(\frac{1}{1 + e^{-(\beta^T X_i + \beta_0)}} \right)) X_i \\ &= \lambda \beta - \sum_{i=1}^n (y_i - \sigma(X, \beta, \beta_0)) X_i \end{aligned}$$

$$\begin{aligned} \nabla_{\beta_0} J(\beta, \beta_0) &= - \sum_{i=1}^n (y_i - \left(\frac{1}{1 + e^{-(\beta^T X_i + \beta_0)}} \right)) \\ &= - \sum_{i=1}^n (y_i - \sigma(X, \beta, \beta_0)) \end{aligned}$$

Gradient Descent Algorithm for Logistic Regression

Algorithm:- Gradient descent method

Initialize β^0 and β_0^0

Repeat

$$\beta^{(k+1)} := \beta^{(k)} - \eta \nabla_{\beta} J(\beta^{(k)}, \beta_0^{(k)})$$

$$\beta_0^{(k+1)} := \beta_0^{(k)} - \eta \nabla_{\beta_0} J(\beta^{(k)}, \beta_0^{(k)})$$

Until $\left| \left| \begin{bmatrix} \nabla_{\beta_1} J(\beta^{(k+1)}, \beta_0^{(k+1)}) \\ \nabla_{\beta_0} J(\beta^{(k+1)}, \beta_0^{(k+1)}) \end{bmatrix} \right| \right| \leq \varepsilon$

Gaussian Basis Function

$$\frac{1}{1 + e^{-f(x)}}$$

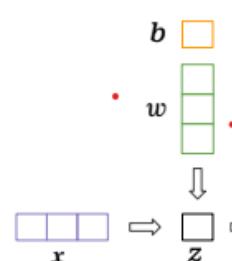
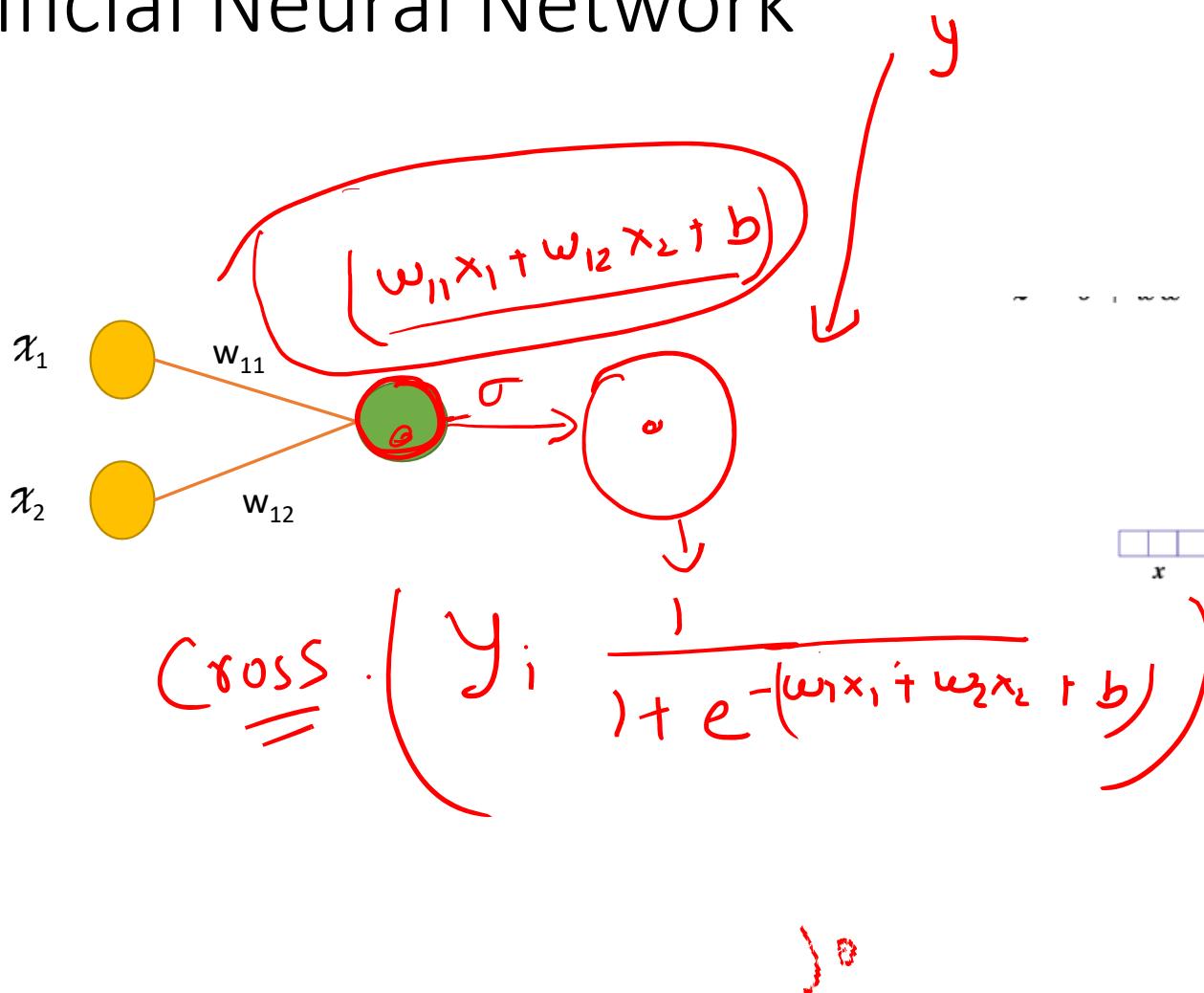
$$f(X) = \beta_M \phi_M(X) + \beta_{M-1} \phi_{M-1}(X) + \dots + \beta_2 \phi_2(X) + \beta_1 \phi_1(X) + \beta_0 \\ = W^T \phi(X)$$

where ,

$$W = \begin{bmatrix} \beta_M \\ \beta_{M-1} \\ \vdots \\ \beta_1 \\ \beta_0 \end{bmatrix}$$

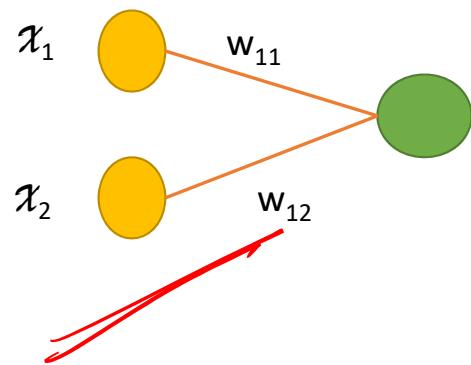
$$\text{and } \phi(X) = \begin{bmatrix} \phi_M(X) \\ \phi_{M-1}(X) \\ \vdots \\ \phi_1(X) \\ \phi_0(X) \end{bmatrix} = \begin{bmatrix} \sigma(\bar{\omega}_M^T X + \bar{b}_M) \\ \sigma(\bar{\omega}_{M-1}^T X + \bar{b}_{M-1}) \\ \vdots \\ \sigma(\bar{\omega}_1^T X + \bar{b}_1) \\ 1 \end{bmatrix}$$

Artificial Neural Network

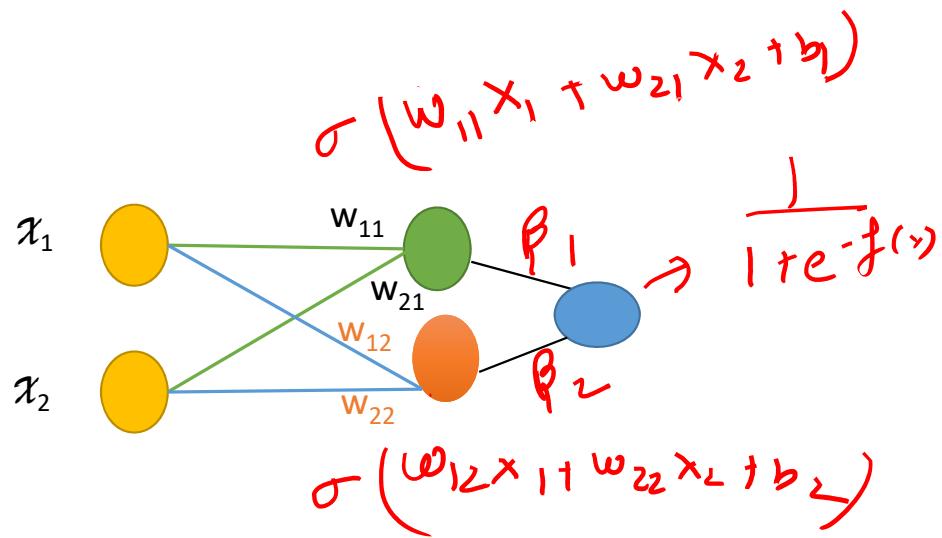


Although artificial neurons are inspired by real neurons, really all we're doing is the dot product of two vectors, followed by element-wise application of the activation function.

Artificial Neural Network



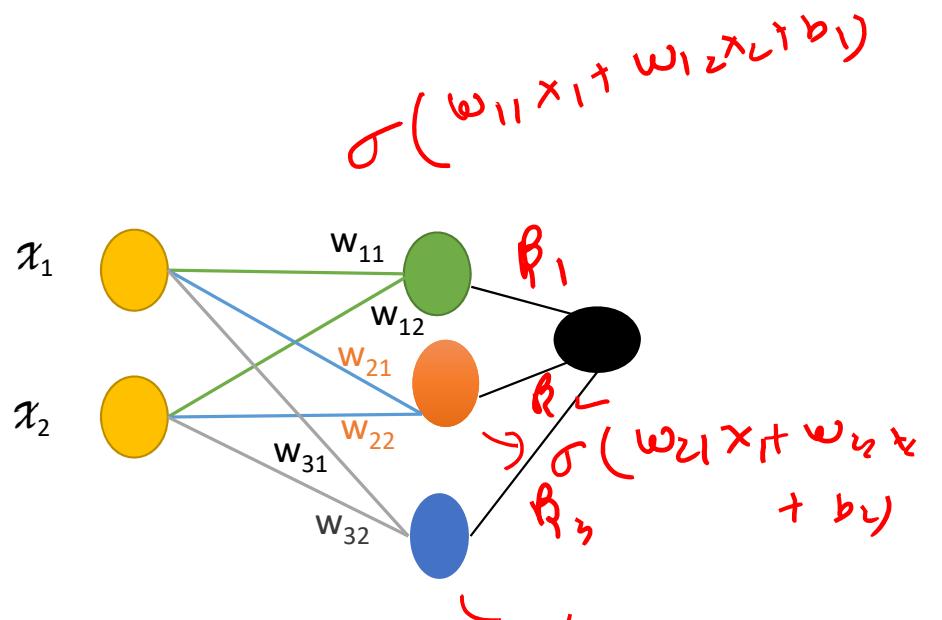
$$f(x) = \beta_1 \sigma(w_{11}x_1 + w_{21}x_2 + b_1) + \beta_2 \sigma(w_{12}x_1 + w_{22}x_2 + b_2)$$

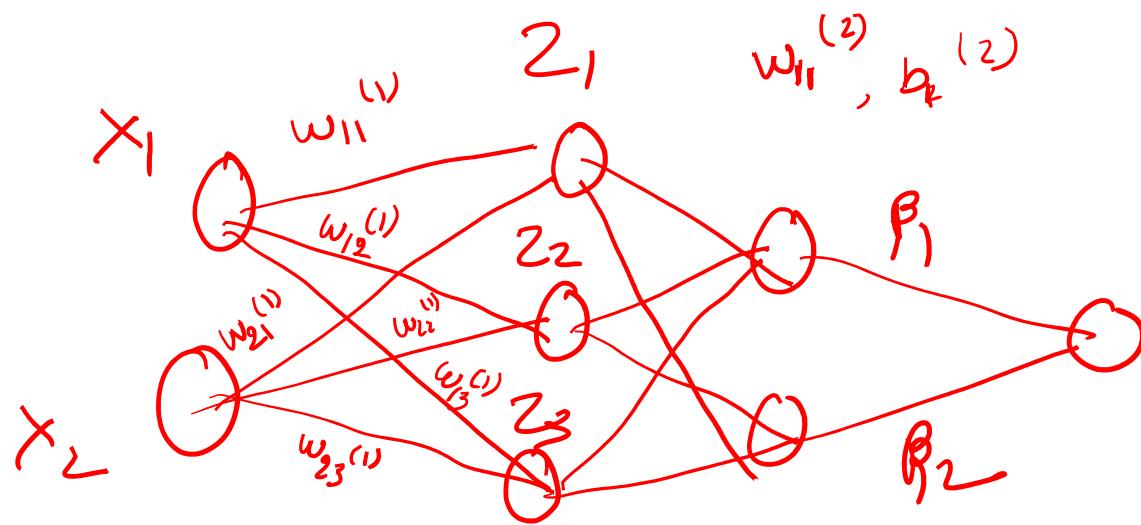


Artificial Neural Network

$$y_i = \frac{1}{1 + e^{-f(x)}}$$

$$f(x) = \beta_1 \sigma(w_{11}x_1 + w_{12}x_2 + b_1) + \beta_2 (\underbrace{w_{21}x_1 + w_{22}x_2 + b_2}_{\sigma(w_{31}x_1 + w_{32}x_2 + b_3)}) + \beta_3 (w_{31}x_1 + w_{32}x_2 + b_3)$$





$$w_1^{(1)} \quad b_1^{(1)}$$

$$z_1 = \sigma(w_{11}^{(1)}x_1 + w_{21}^{(1)}x_2 + b_1)$$

$$z_2 = \sigma(w_{12}^{(1)}x_1 + w_{22}^{(1)}x_2 + b_2)$$

$$z_3 = \sigma(w_{13}^{(1)}x_1 + w_{23}^{(1)}x_2 + b_3)$$

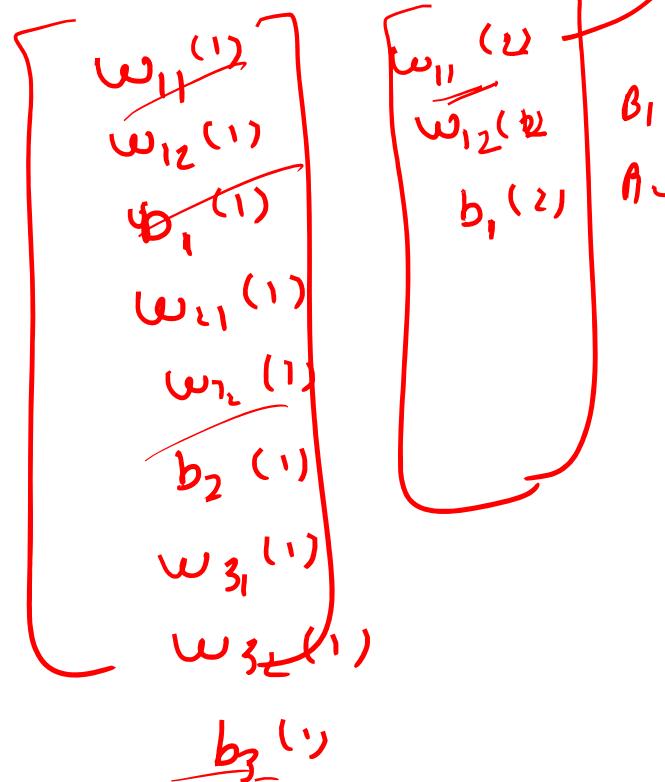
$$\rightarrow X_1 = G(\underbrace{z_1 w_{11}^{(2)} + z_2 w_{21}^{(2)} + z_3 w_{31}^{(2)} + c_1}_{(2)})$$

$$Y_2 = G(z_1 w_{12}^{(2)} + z_2 w_{22}^{(2)} + z_3 w_{32}^{(2)} + c_2)$$

\$ \rightarrow \$ Final function : $\frac{1}{1 + e^{-(\beta_1 Y_1 + \beta_2 Y_2)}}$

$$L = -y_i \log \left(\frac{1}{1 + e^{-(\beta_1 y_1 + \beta_2 y_2)}} \right) - (1-y_i) \log \left(1 - \frac{1}{1 + e^{-(\beta_1 y_1 + \beta_2 y_2)}} \right)$$

$$\frac{\partial L}{\partial y_1} \cdot \frac{\partial y_1}{\partial z_1} \quad \frac{\partial L}{\partial w_{11}^{(1)}}$$

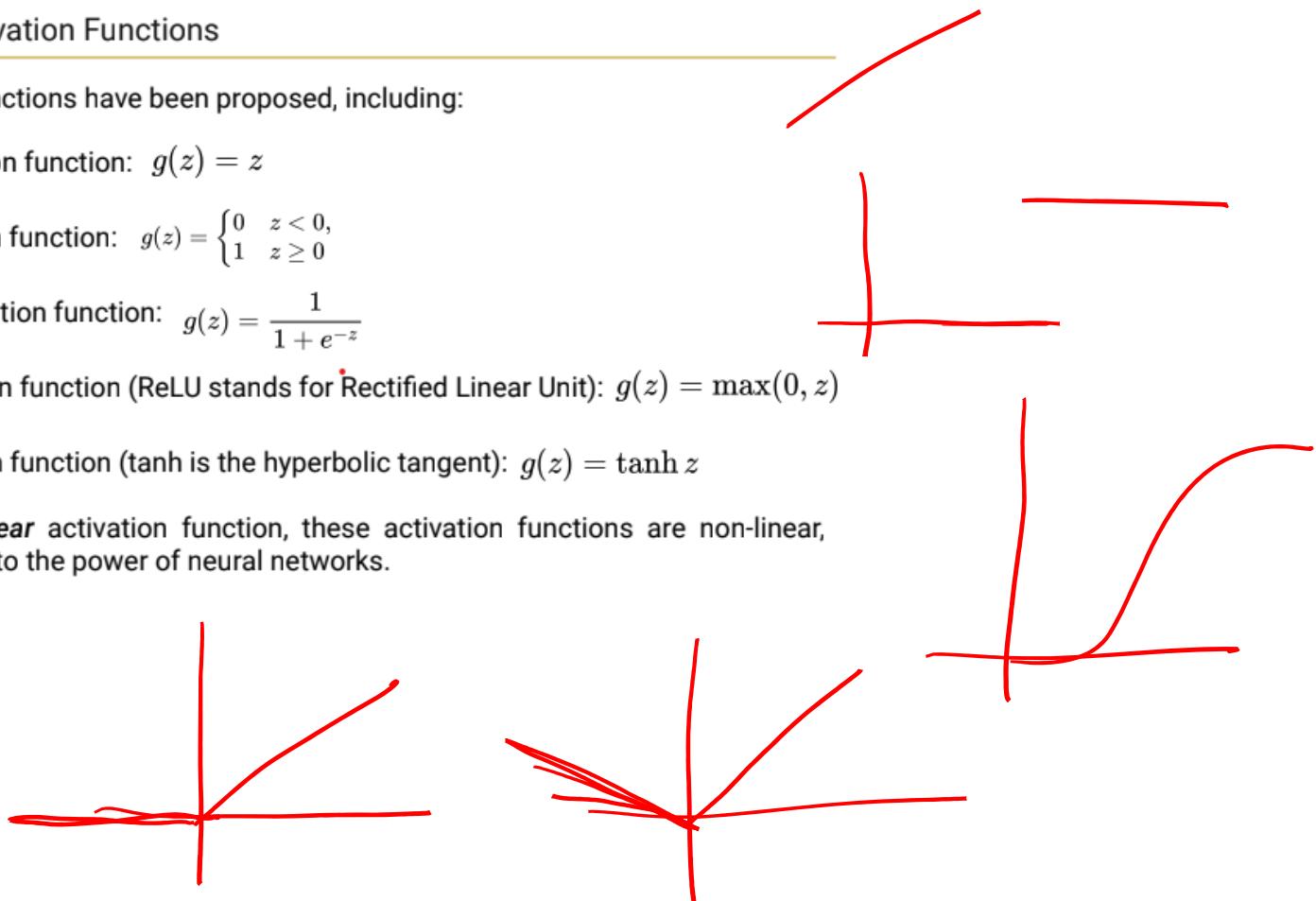


Introduction: Activation Functions

Many activation functions have been proposed, including:

- *linear* activation function: $g(z) = z$
- *step* activation function: $g(z) = \begin{cases} 0 & z < 0, \\ 1 & z \geq 0 \end{cases}$
- *sigmoid* activation function: $g(z) = \frac{1}{1 + e^{-z}}$
- *ReLU* activation function (ReLU stands for Rectified Linear Unit): $g(z) = \max(0, z)$
- *tanh* activation function (*tanh* is the hyperbolic tangent): $g(z) = \tanh z$

Apart from the *linear* activation function, these activation functions are non-linear, which is important to the power of neural networks.



$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_L, y_L)\}$$

$$\max \underline{P}(y_1, y_2, \dots, y_L | x_1, x_2, \dots, x_L)$$

$$= \max_{\text{over } i=1} \prod_{i=1}^L P(y_i | x_i)$$

$$= \max \sum_{i=1}^L \log P(y_i | x_i)$$

$$= \max \sum_{i=1}^L \log \left(\hat{P}_{x_i}^{y_i} (1 - \hat{P}_{x_i})^{1-y_i} \right) = \hat{P}_{x_i}^{y_i} (1 - \hat{P}_{x_i})^{1-y_i}$$

$$= \max \sum_{i=1}^L \left(y_i \log \hat{P}_{x_i} + (1-y_i) \log (1 - \hat{P}_{x_i}) \right)$$

$$y_i | x_i \sim \beta(\rho_{x_i})$$

$$y_i = 1 | x_i = \hat{P}_{x_i}$$

$$y_i = 0 | x_i = 1 - \hat{P}_{x_i}$$

$$\min_f \sum_{i=1}^n -y_i \log f(x_i) - (1-y_i) \log (1-f(x_i))$$

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\} \quad y_i | x_i \sim B(P_{x_i})$$

$$y_i = 0, 1, i=1, 2, \dots, N$$

$$x_i \in R^n$$

$$\max \prod_{i=1}^N P(y_1, y_2, \dots, y_N) | (x_1, x_2, \dots, x_N)$$

$$\max \prod_{i=1}^N P(y_i | x_i)$$

$$\max \log \prod_{i=1}^N P(y_i | x_i) = \max \sum_{i=1}^N \log P(y_i | x_i)$$

$$= \max \sum_{i=1}^N \left(y_i \log P_{x_i} + (1-y_i) \log (1-P_{x_i}) \right)$$

$$P_{x_i}^{y_i} (1-P_{x_i})^{1-y_i}$$

$$P_{x_i}^{y_i} P(x_i | x_i)$$

$$= P_{x_i} \text{ if } y_i \\ = 1 - P_{x_i} \text{ if } y_i$$

$$\max \sum_{i=1}^N \frac{\log P(y_i | x_i)}{\log(P_{x_i}^{y_i} (1-P_{x_i})^{1-y_i})}$$

$$\min \sum_{i=1}^n -y_i \log(p_{x_i}) - (1-y_i) \log(1-p_{x_i})$$

$$P(y_i=1/x_i) = \frac{1}{1+e^{-(\beta_1 x_i + \beta_0)}} = f(x_i)$$

$$p_{x_i} = \boxed{P(y_i/x_i)}$$

$$\min \sum_{i=1}^n -y_i \log(f(x_i)) - (1-y_i) \log(1-f(x_i))$$

Linear Logistic Regression with k variables

$$f(x) = \frac{1}{1+e^{-(\beta^T x + \beta_0)}} = \frac{e^{(\beta^T x + \beta_0)}}{1+e^{(\beta^T x + \beta_0)}}$$

where $\beta = \begin{bmatrix} \beta_k \\ \vdots \\ \beta_1 \end{bmatrix}$

Some notes on Evaluation Metric

		1	0
		Predicted	Predicted
Predicted		Fraudulent	Legitimate
Fraudulent	Transaction	True Positive	False Positive
Legitimate	Transaction	False Negative	True Negative

Some notes on Evaluation Metric

		Predicted Fraudulent	Predicted Legitimate
Predicted	Fraudulent Transaction	True Positive	False Positive
	Legitimate Transaction	False Negative	True Negative

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+FN+TN} \times 100$$

$$\text{Error} = 1 - \text{Accuracy}$$

Some notes on Evaluation Metric

		1	0
		Predicted Fraudulent	Predicted Legitimate
Predicted	Fraudulent Transaction	70	30
	Legitimate Transaction	100	1800

Accuracy = ??

Error = ??

$$\frac{130 + 1870}{2000} \times 100$$

Some notes on Evaluation Metric

		Predicted Fraudulent	Predicted Legitimate
Predicted	Fraudulent Transaction	True Positive	False Positive
	Legitimate Transaction	False Negative	True Negative

What proportion of positive identifications was actually correct?

$$\text{Precision} = \frac{TP}{TP+FP}$$

Some notes on Evaluation Metric

	Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	70	30
Legitimate Transaction	100	1800

Precision = ??

$$\frac{70}{100} \times 100$$

Some notes on Evaluation Metric

		Predicted Fraudulent	Predicted Legitimate
Predicted Fraudulent Transaction	Predicted Fraudulent	True Positive	False Positive
	Predicted Legitimate	False Negative	True Negative

What proportion of actual positives was identified correctly?

$$\text{Recall} = \frac{TP}{TP+FN}$$

Also called sensitivity some time.

Some notes on Evaluation Metric

	Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	70	30
Legitimate Transaction	100	1800

Recall = ??

Some notes on Evaluation Metric

		1	0
		Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	Fraudulent	True Positive	False Positive
	Legitimate	False Negative	True Negative

$$\text{Specificity} = \frac{TN}{TN+FP}$$

Some notes on Evaluation Metric

	Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	70	30
Legitimate Transaction	100	1800

Specificity = ??

Some notes on Evaluation Metric

	Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	0	10
Legitimate Transaction	190	1800

Accuracy = ??

Some notes on Evaluation Metric

	Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	True Positive	False Positive
Legitimate Transaction	False Negative	True Negative

10,000
99.00 100
↓ ↓
0 1

$$F1 \text{ score} = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

$$Gmean = \sqrt{\text{Specificity} \times \text{Recall}}$$

Some notes on Evaluation Metric

	Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	True Positive	False Positive
Legitimate Transaction	False Negative	True Negative

$$F1 \text{ score} = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

$$Gmean = \sqrt{Specificity \times Recall}$$

Some notes on Evaluation Metric

	Predicted Fraudulent	Predicted Legitimate
Fraudulent Transaction	0	10
Legitimate Transaction	190	1800

F1 Score ??

G-mean ??