

CS60010: Deep Learning

Spring 2021

Sudeshna Sarkar and Abir Das

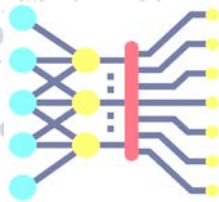
Module 2 Part 2

Linear Models for Classification

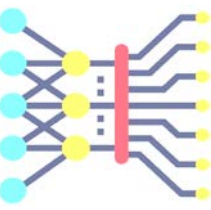
Sudeshna Sarkar

18 Jan 2021

Announcements



- **Class Test 1 on 19th Jan 2021 (tomorrow)**
- Assignment 1 due 22nd Jan 12 pm



ML Background and Linear Models

Logistic Regression



Classification

A binary classifier is a mapping from $R^d \rightarrow \{-1, +1\}$

$$x \rightarrow h \rightarrow y$$

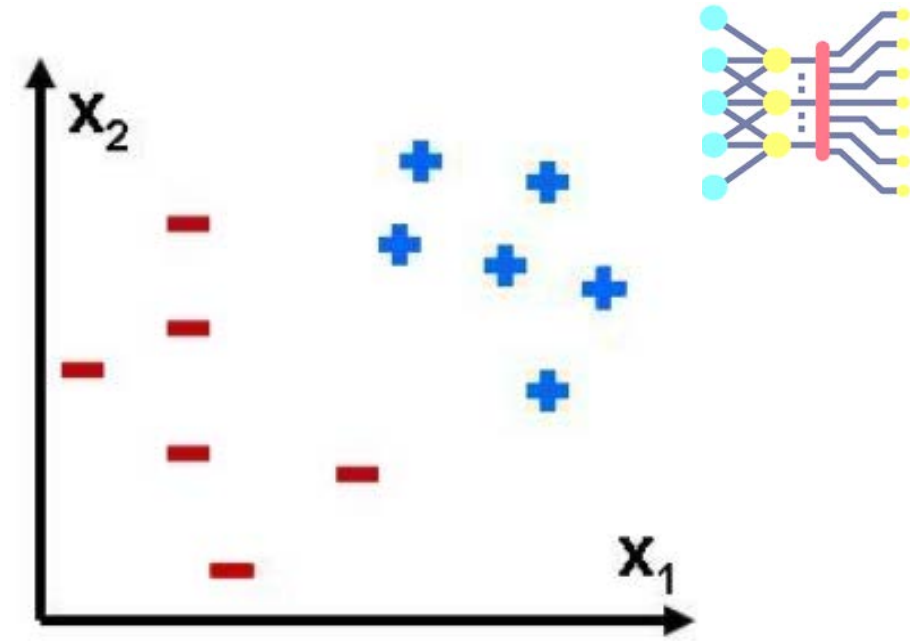
Training data set

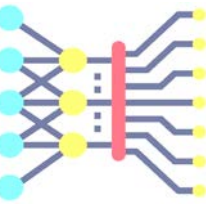
$$\mathcal{D}_m = \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$$

- Assume that each $x^{(i)}$ is a $d \times 1$ column vector
- Given a training set \mathcal{D}_m and a classifier h , we can define the training error of h to be

$$\varepsilon_m(h) = \frac{1}{m} \sum_{i=1}^m \begin{cases} 1 & h(x^{(i)}) \neq y^{(i)} \\ 0 & \text{otherwise} \end{cases}$$

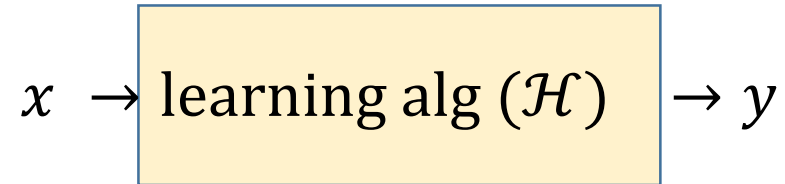
- For now, we will try to find a classifier with small training error (and hope it generalizes well to new data, and has a small test error)



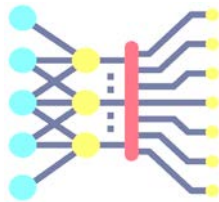


Learning algorithm

- A hypothesis class \mathcal{H} is a set (finite or infinite) of possible classifiers, each of which represents a mapping from $R^d \rightarrow \{-1, +1\}$
- A learning algorithm is a procedure that takes a data set \mathcal{D}_n as input and returns an element $h \in \mathcal{H}$



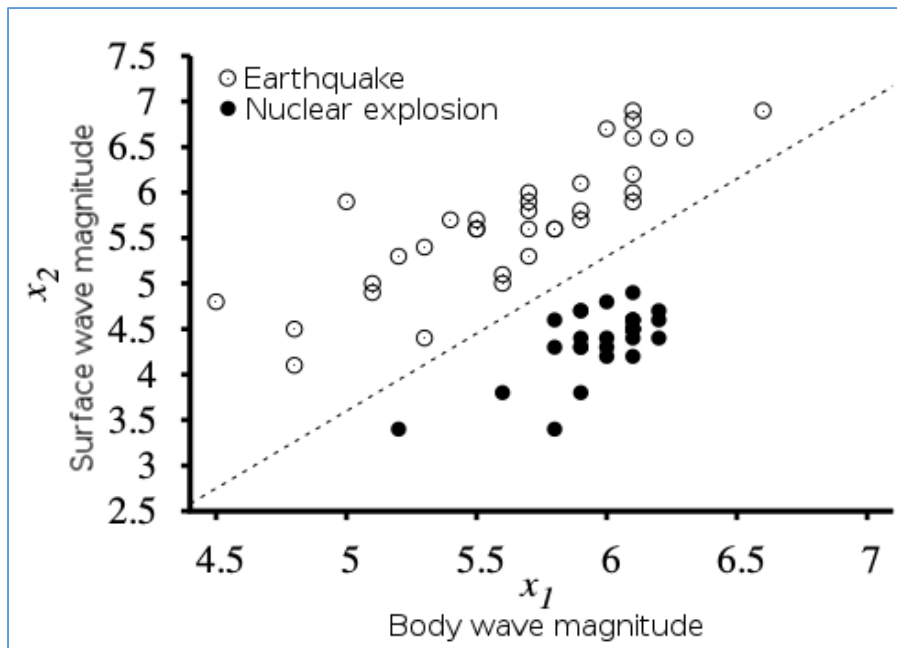
- Choice of \mathcal{H} so as to get low test error



Hypothesis class : Linear classifiers

- A linear classifier in d dimensions is defined by
 - A vector of parameters $\theta \in R^d$ and scalar $\theta_0 \in \mathcal{R}$
Assume a $d \times 1$ column vector

$$h(x; \theta, \theta_0) = \text{sign}(\theta^T x + \theta_0) = \begin{cases} +1 & \text{if } \theta^T x + \theta_0 > 0 \\ -1 & \text{otherwise} \end{cases}$$



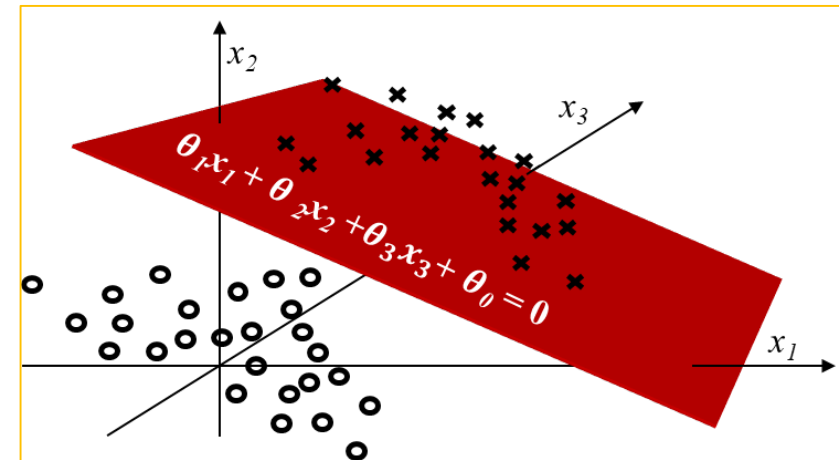
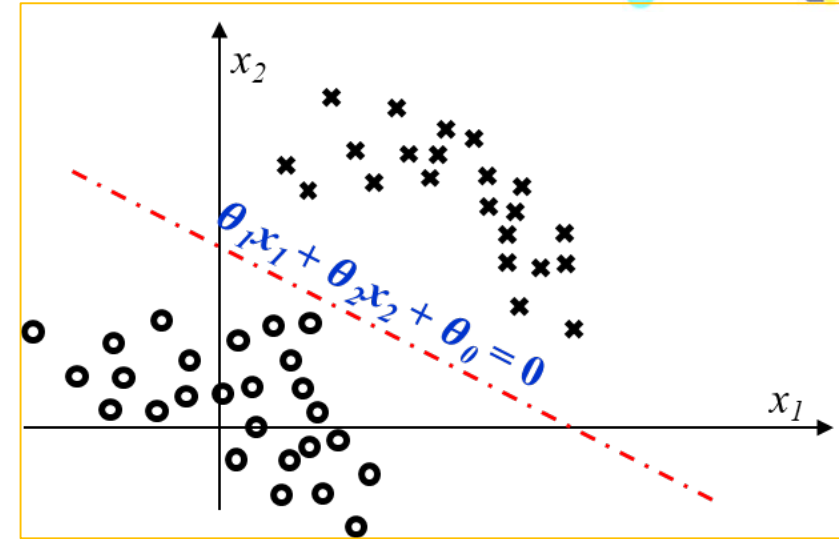
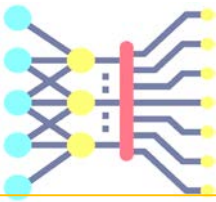
θ, θ_0 specifies a hyperplane (decision boundary) that divides the instance space into two half-spaces.

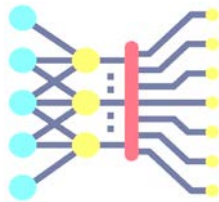
Linear classifiers

$$h(x; \theta, \theta_0) = \text{sign}(\theta^T x + \theta_0)$$

$$= \begin{cases} +1 & \text{if } \theta^T x + \theta_0 > 0 \\ -1 & \text{otherwise} \end{cases}$$

- θ, θ_0 specifies a hyperplane that divides the instance space into two half-spaces.
 - The one that is on the same side as the normal vector is the positive half-space, and we classify all points in that space as positive.
 - The half-space on the other side is negative and all points in it are classified as negative.





Linear Classifier with Hard Threshold

- The linear separator in the associated fig is given by

$$x_2 = 1.7x_1 - 4.9$$

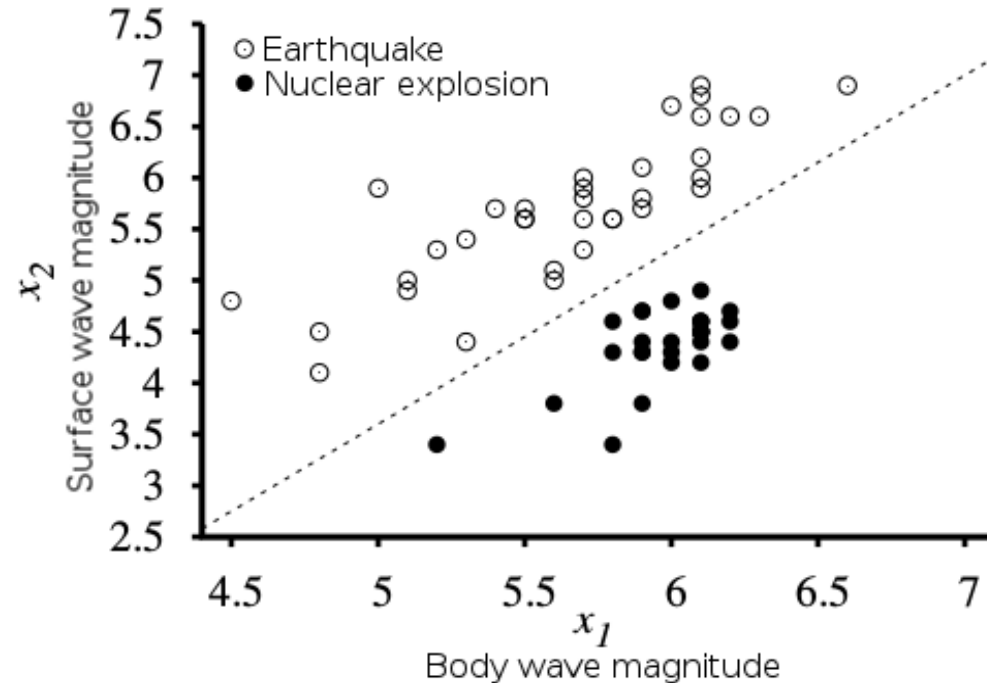
$$\rightarrow -4.9 + 1.7x_1 - x_2 = 0$$

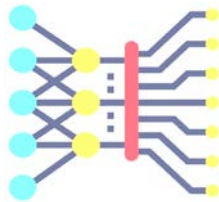
$$\rightarrow [-4.9 \quad 1.7 \quad -1] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = 0$$

$$\boldsymbol{\theta}^T \mathbf{x} = 0$$

Classification Rule:

$$y(x) = \begin{cases} +1 & \text{if } \boldsymbol{\theta}^T \mathbf{x} > 0 \\ -1 & \text{otherwise} \end{cases}$$





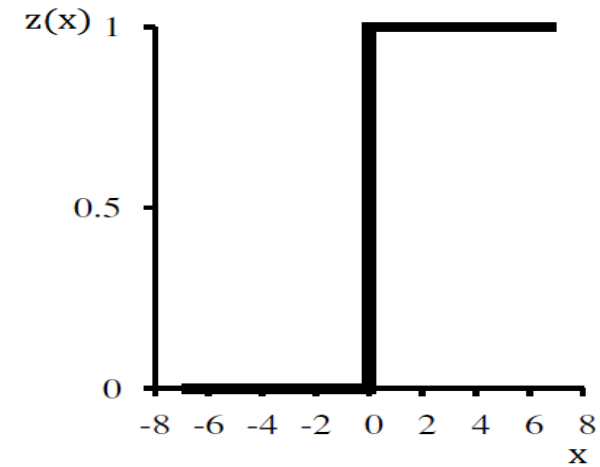
Linear Classifier with Hard Threshold

Classification Rule:

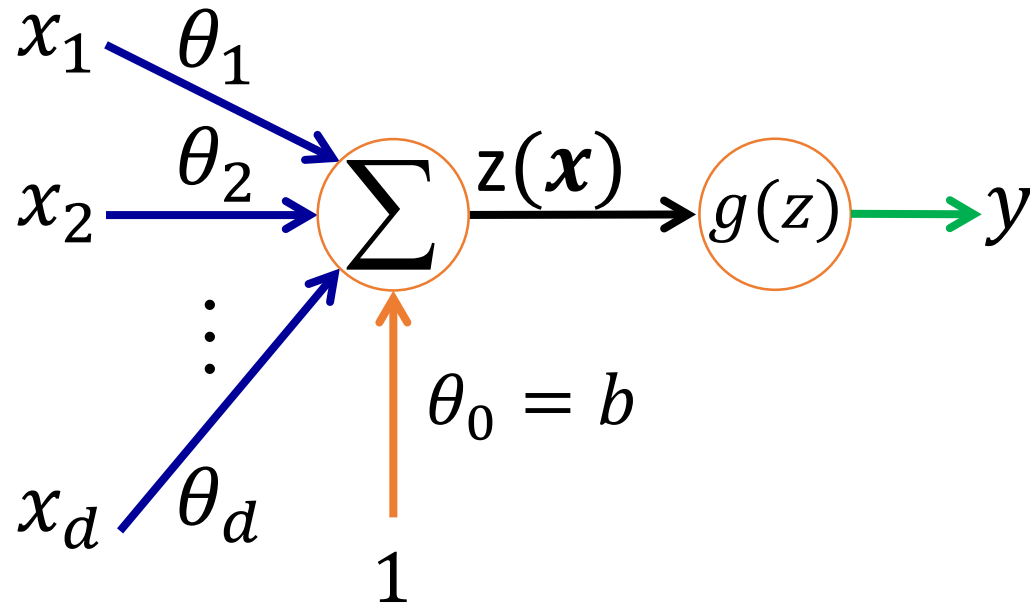
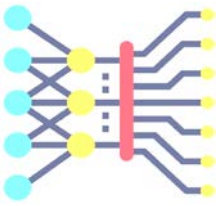
$$y(x) = \begin{cases} +1 & \text{if } \theta^T x > 0 \\ -1 & \text{otherwise} \end{cases}$$

We can think y as the result of passing the linear function $\theta^T x$ through a threshold function.

- Find the θ which minimizes classification error on the training set.
- We cannot use gradient descent at all points for the above threshold function



Perceptron



$$\boldsymbol{\theta} = [\theta_1 \ \theta_2 \ \dots \ \theta_d]^T \text{ and } \mathbf{x} = [x_1 \ x_2 \ \dots \ x_d]^T$$

$$z(\mathbf{x}) = \theta_0 + \sum_{i=1}^d \theta_i x_i = [\boldsymbol{\theta}^T \ \mathbf{b}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$
$$y = g(z(\mathbf{x}))$$

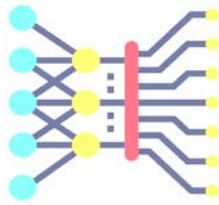
Terminologies

\mathbf{x} : input, $\boldsymbol{\theta}$: weights, \mathbf{b} : bias

z : pre-activation (input activation)

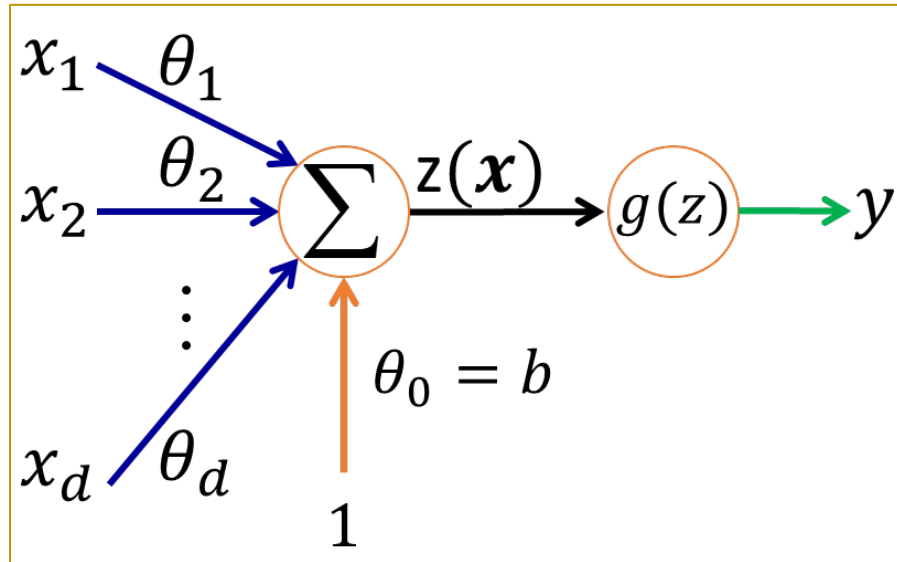
g : activation function

y : activation (output activation)



Perceptron

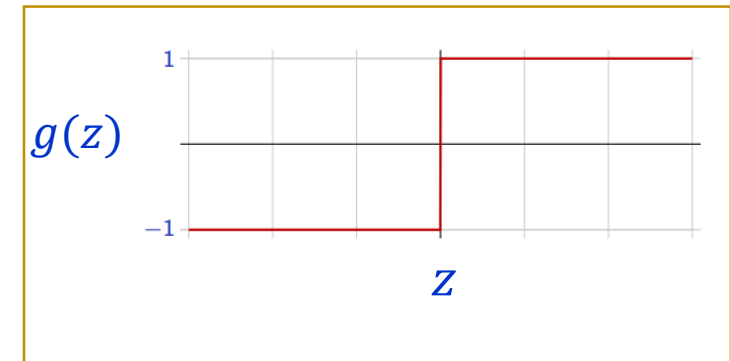
$\mathbf{x} \in \mathcal{R}^d$ and $y \in \{0, 1\}$ for Binary Classification



$$g(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases} \quad (\text{Rosenblatt, 1957})$$

Or, the response may be taken as $y \in \{-1, 1\}$

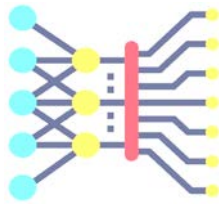
$$g(z) = \begin{cases} 1, & z \geq 0 \\ -1, & z < 0 \end{cases}$$



The perceptron classification rule, thus, translates to

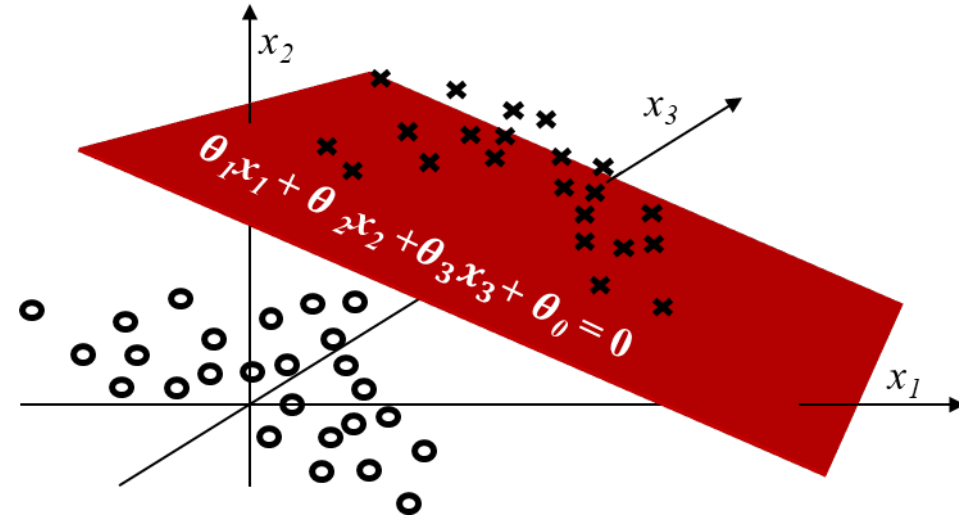
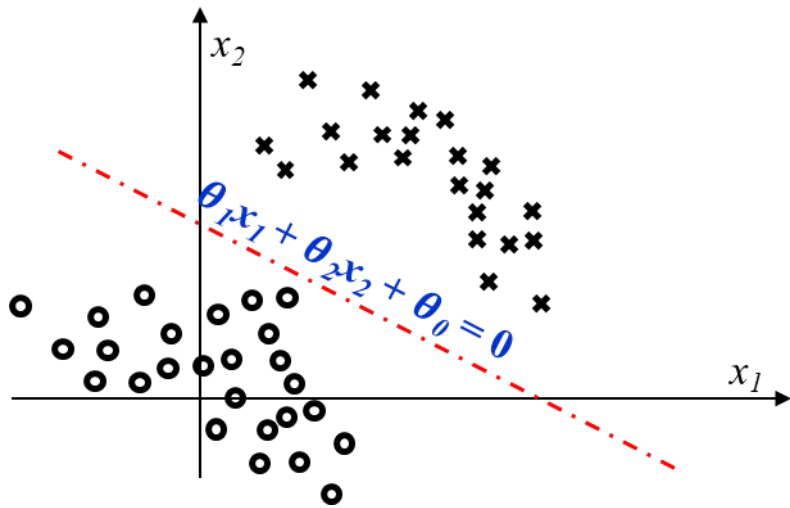
$$y = \begin{cases} 1, & \boldsymbol{\theta}^T \mathbf{x} + b \geq 0 \\ -1, & \boldsymbol{\theta}^T \mathbf{x} + b < 0 \end{cases}$$

$\boldsymbol{\theta}^T \mathbf{x} + b = 0$ represents a hyperplane.



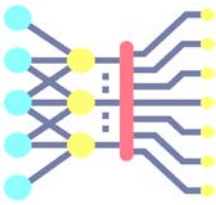
Perceptron (Geometrically)

$$h_{\theta}(\mathbf{x}) = \text{sign}(\theta^T \mathbf{x})$$



- $\theta^T \mathbf{x}$ is the (signed) distance of point \mathbf{x} to hyperplane
 - Example: http://mathinsight.org/distance_point_plane

Perceptron Learning Algorithm



Training Set: $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

$x^{(i)} \in R^d; y^{(i)} \in \{-1, +1\}$

1. $t \leftarrow 1; \theta^{(t)} = \mathbf{0}$

2. // Loop until all examples are correctly classified

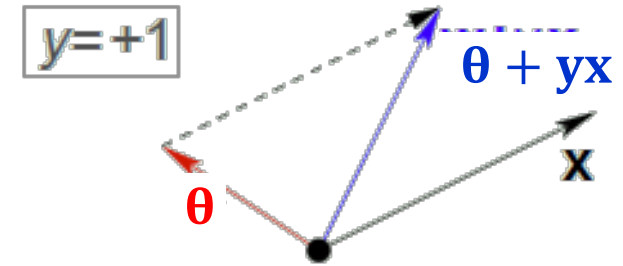
While exists i such that $x^{(m)}$ is not correctly classified

Pick a j s.t. $y^{(j)} \cdot \langle \theta^{(t)}, x^{(j)} \rangle \leq 0$ then

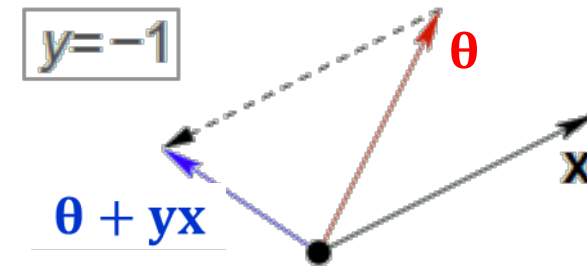
$$\theta^{(t+1)} = \theta^{(t)} + y^{(j)} x^{(j)}$$

$$t \leftarrow t + 1$$

3. Output $\theta^{(t)}$

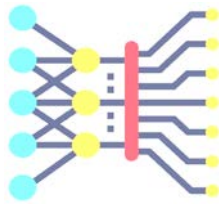


$\theta + x$ pushes vector θ towards x



$\theta - x$ pushes vector θ away from x

Perceptron Learning Algorithm

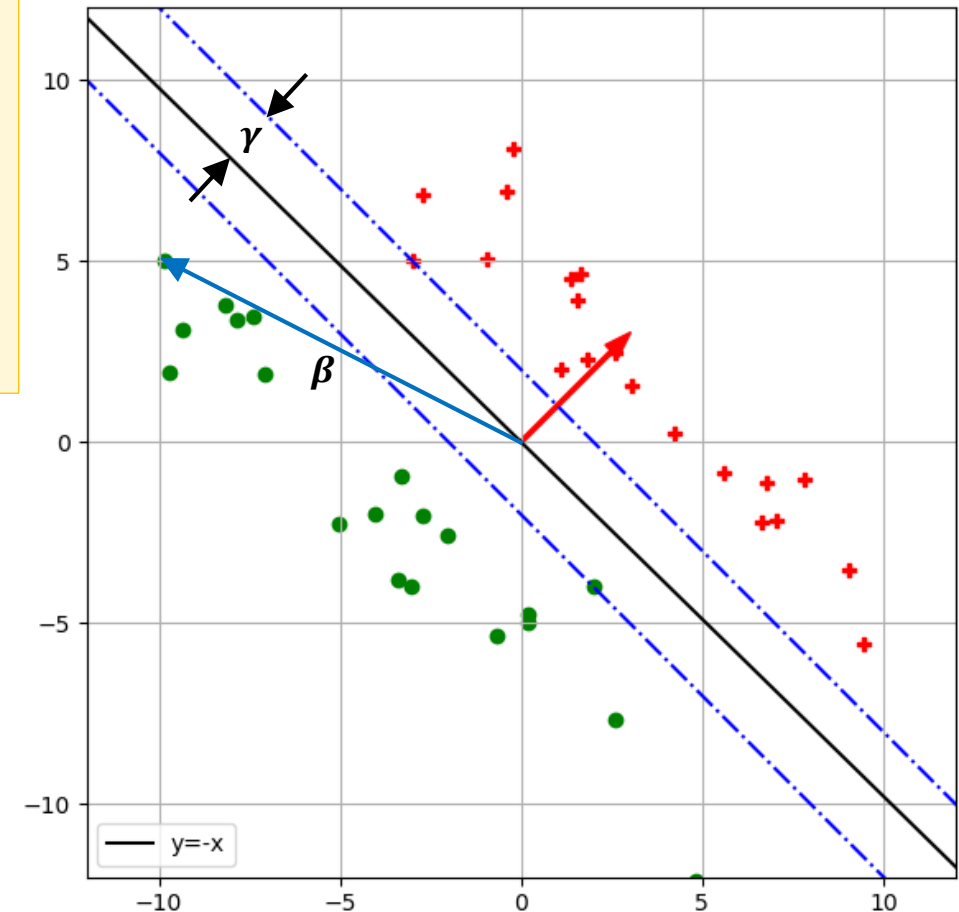


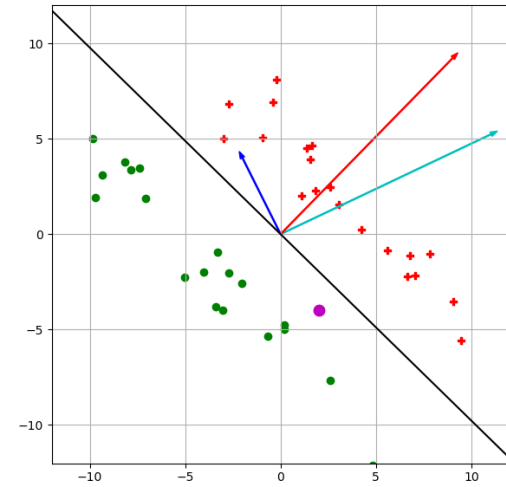
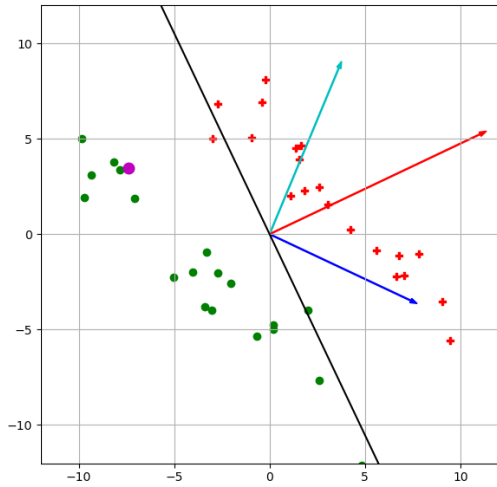
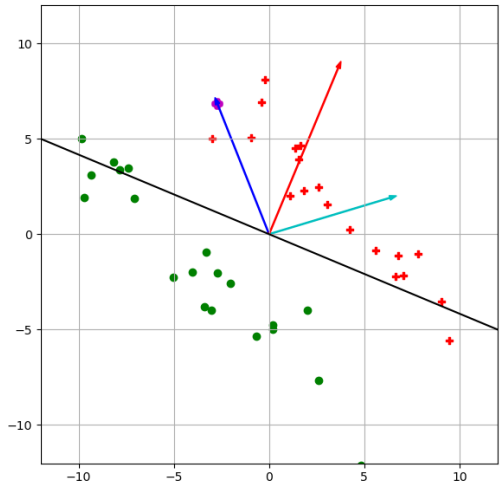
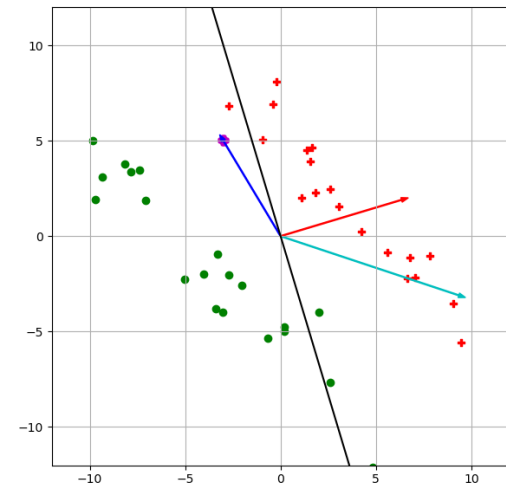
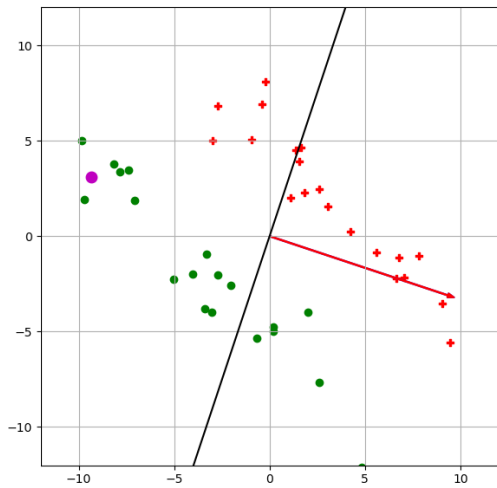
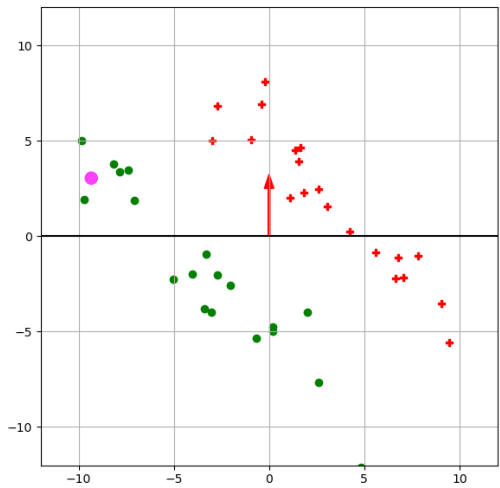
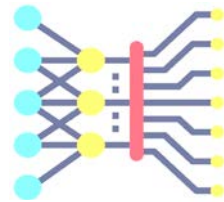
1. $t \leftarrow 1; \theta^{(t)} = \mathbf{0}$
2. While there exists i such that $x^{(i)}$ is not correctly classified
Pick a j s.t. $y^{(j)} \cdot \langle \theta^{(t)}, x^{(j)} \rangle \leq 0$
 $\theta^{(t+1)} = \theta^{(t)} + y^{(j)} x^{(j)}$
 $t \leftarrow t + 1$
3. Output $\theta^{(t)}$

If such a separating hyperplane exists, then the data is known to be linearly separable.

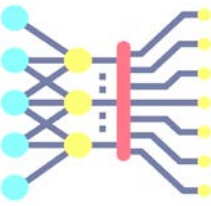
Convergence Theorem

For a finite and linearly separable set of data, the perceptron learning algorithm will find a linear separator in at most $\frac{\beta^2}{\gamma^2}$ iterations where the maximum length of any data point is β and γ is the maximum margin of the linear separators.

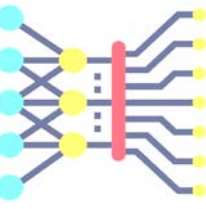




Perceptron Rule



- Perceptron Learning Rule can find a linear separator given the data is *linearly separable*.
- For data that are not linearly separable, the Perceptron algorithm fails.



Linear Classifiers by Gradient Descent

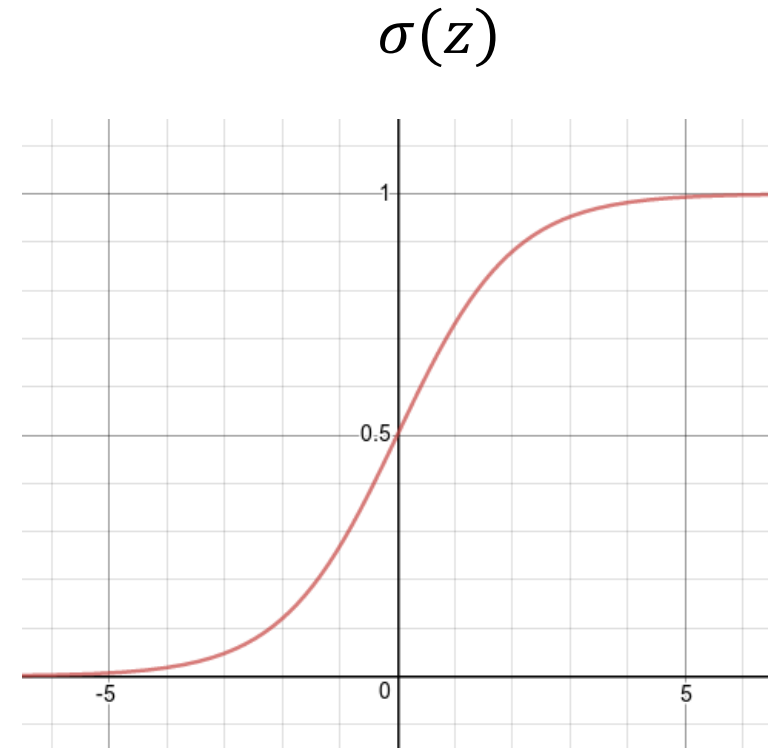
For a gradient based optimization approach, we need to approximate hard threshold function with something smooth.

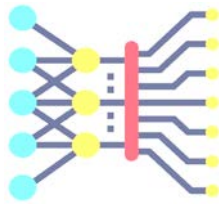
- Logistic regression classifier

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

$$y = \sigma(h_{\theta}(\mathbf{x})) = \sigma(\boldsymbol{\theta}^T \mathbf{x})$$

$$z = \boldsymbol{\theta}^T \mathbf{x}$$





Likelihood Function for Logistic Regression

- The probability that an example belongs to class 1 is

$$P(y^{(i)} = 1|x^{(i)}; \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)})$$

Thus $P(y^{(i)} = 0|x^{(i)}; \boldsymbol{\theta}) = 1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)})$

Thus

$$\begin{aligned} &P(y^{(i)} | x^{(i)}; \boldsymbol{\theta}) \\ &= \left(\sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) \right)^{y^{(i)}} \left(1 - \sigma(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) \right)^{1-y^{(i)}} \end{aligned}$$

$$z^{(i)} = \boldsymbol{\theta}^T x^{(i)}$$

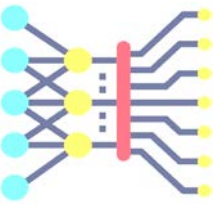
$$P(y^{(i)} = 1|x^{(i)})\sigma(z^{(i)})$$

Thus $P(y^{(i)} = 0|x^{(i)}) = 1 - \sigma(z^{(i)})$

$$P(y^{(i)} | x^{(i)})$$

$$= \left(\sigma(z^{(i)}) \right)^{y^{(i)}} \left(1 - \sigma(z^{(i)}) \right)^{1-y^{(i)}}$$

Maximum Likelihood Estimation of Logistic Regression



The probability that an example belongs to class 1 is $P(y^{(i)} = 1|x^{(i)}; \theta) = \sigma(\theta^T \mathbf{x}^{(i)})$

Thus $P(y^{(i)} = 0|x^{(i)}; \theta) = 1 - \sigma(\theta^T \mathbf{x}^{(i)})$

Thus $P(y^{(i)}|x^{(i)}; \theta) = \left(\sigma(\theta^T \mathbf{x}^{(i)})\right)^{y^{(i)}} \left(1 - \sigma(\theta^T \mathbf{x}^{(i)})\right)^{1-y^{(i)}}$

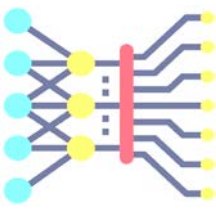
The joint probability of all the labels

$$\prod_{i=1}^m \left(\sigma(\theta^T \mathbf{x}^{(i)})\right)^{y^{(i)}} \left(1 - \sigma(\theta^T \mathbf{x}^{(i)})\right)^{1-y^{(i)}}$$

So the log likelihood for logistic regression is given by

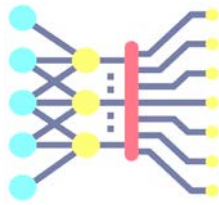
$$l(\theta) = \sum_{i=1}^m y^{(i)} \log \left(\sigma(\theta^T \mathbf{x}^{(i)})\right) + (1 - y^{(i)}) \log \left(1 - \sigma(\theta^T \mathbf{x}^{(i)})\right)$$

Maximum Likelihood Estimation of Logistic Regression



Derivative of log likelihood w.r.t. one component of θ

$$\begin{aligned}\frac{\partial l(\theta)}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \sum_{i=1}^m y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - \sigma(\theta^T \mathbf{x}^{(i)})) && \text{derivative of sum of terms} \\ &= \sum_{i=1}^m \left[\frac{y^{(i)}}{\sigma(\theta^T \mathbf{x}^{(i)})} - \frac{1 - y^{(i)}}{1 - \sigma(\theta^T \mathbf{x}^{(i)})} \right] \frac{\partial}{\partial \theta_j} \sigma(\theta^T \mathbf{x}^{(i)}) && \text{derivative of log f(x)} \\ &= \sum_{i=1}^m \left[\frac{y^{(i)}}{\sigma(\theta^T \mathbf{x}^{(i)})} - \frac{1 - y^{(i)}}{1 - \sigma(\theta^T \mathbf{x}^{(i)})} \right] \sigma(\theta^T \mathbf{x}^{(i)}) (1 - \sigma(\theta^T \mathbf{x}^{(i)})) \mathbf{x}_j^{(i)} && \text{chain rule + derivative of } \sigma \\ &= \sum_{i=1}^m \left[\frac{y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})}{\sigma(\theta^T \mathbf{x}^{(i)}) (1 - \sigma(\theta^T \mathbf{x}^{(i)}))} \right] \sigma(\theta^T \mathbf{x}^{(i)}) (1 - \sigma(\theta^T \mathbf{x}^{(i)})) \mathbf{x}_j^{(i)} \\ &= \sum_{i=1}^m \left[y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)}) \right] \mathbf{x}_j^{(i)} && (12)\end{aligned}$$

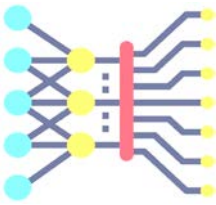


Calculating derivatives

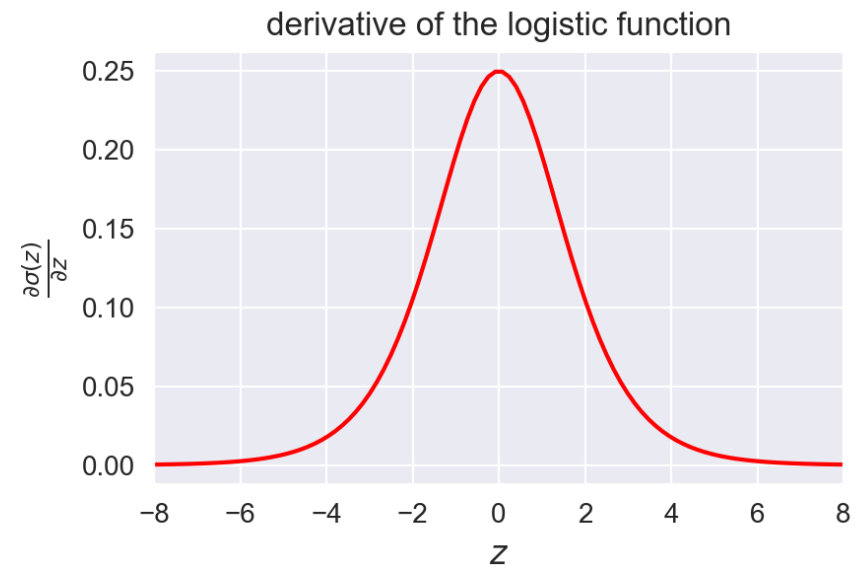
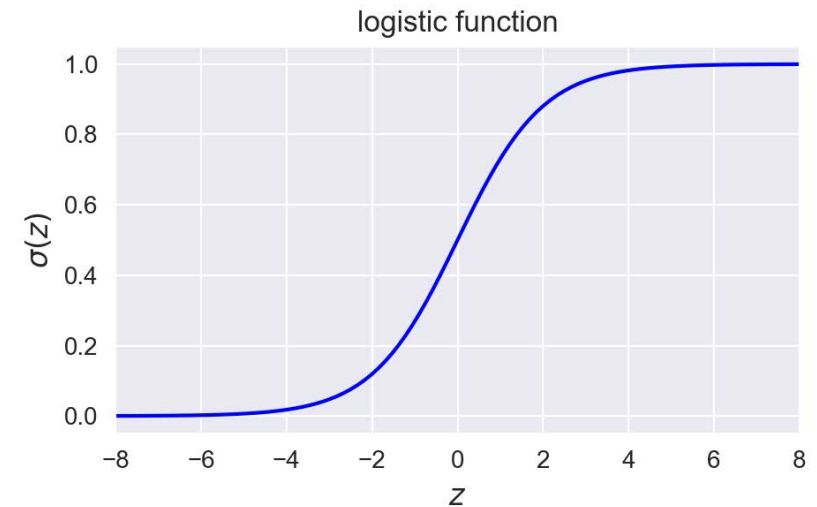
- Since the likelihood function is a sum over all of the data, and in calculus the derivative of a sum is the sum of derivatives, we can focus on computing the derivative of one example. The gradient of theta is simply the sum of this term for each training data point.
- The derivative of gradient for one data point (x, y):
- $p = \sigma(\theta^T x)$
- $z = \theta^T x$

$$\begin{aligned}\frac{\partial l(\theta)}{\partial \theta_j} &= \frac{\partial l(\theta)}{\partial p} \cdot \frac{\partial p}{\partial \theta_j} \\ &= \frac{\partial l(\theta)}{\partial p} \cdot \frac{\partial p}{\partial z} \cdot \frac{\partial z}{\partial \theta_j}\end{aligned}$$

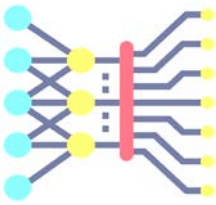
Derivative of logistic function



$$\begin{aligned}\frac{\partial}{\partial z} \sigma(z) &= \frac{\partial}{\partial z} \frac{1}{1 + e^{-z}} \\ &= \frac{-1}{(1 + e^{-z})^2} \cdot e^{-z} \cdot -1 \\ &= \frac{1}{1 + e^{-z}} \frac{e^{-z}}{1 + e^{-z}} \\ &= \sigma(z)(1 - \sigma(z)) \\ &= y(1 - y)\end{aligned}$$



Calculating derivatives



$$\frac{\partial l(\theta)}{\partial \theta_j} = \frac{\partial l(\theta)}{\partial p} \cdot \frac{\partial p}{\partial z} \cdot \frac{\partial z}{\partial \theta_j}$$

$$l(\theta) = y \log p + (1 - y) \log(1 - p)$$

$$\frac{\partial l(\theta)}{\partial p} = \frac{y}{p} - \frac{1-y}{1-p}$$

$$\frac{\partial p}{\partial z} = \sigma(z)(1 - \sigma(z)) \quad p = \sigma(z)$$

$$\frac{\partial z}{\partial \theta_j} = x_j \quad z = \theta^T x$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = \frac{\partial l(\theta)}{\partial p} \cdot \frac{\partial p}{\partial z} \cdot \frac{\partial z}{\partial \theta_j}$$

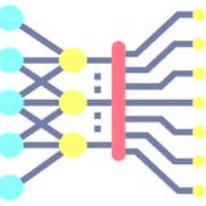
$$= \left[\frac{y}{p} - \frac{1-y}{1-p} \right] \cdot \sigma(z)(1 - \sigma(z)) \cdot x_j$$

$$= \left[\frac{y}{p} - \frac{1-y}{1-p} \right] \cdot p(1 - p) \cdot x_j$$

$$= [y(1 - p) - p(1 - y)] \cdot x_j$$

$$= [y - p] \cdot x_j$$

$$= [y - \sigma(\theta^T x)] \cdot x_j \quad p = \sigma(\theta^T x)$$

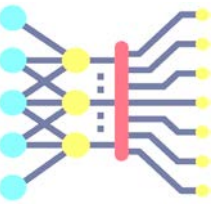


Gradient of Log Likelihood

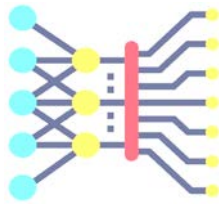
$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m [y^{(i)} - \sigma(\theta^T x^{(i)})] x_j^{(i)}$$

- We need to choose the values of theta that maximize the log-likelihood.
- Unfortunately, if we try just setting the derivative equal to zero, there's no closed form for the maximum.
- However, we can find the best values of theta by using an optimization algorithm.

Gradient Ascent Optimization

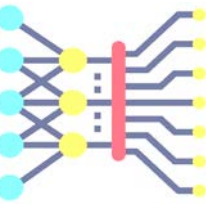


$$\begin{aligned}\theta_j^{\text{new}} &= \theta_j^{\text{old}} + \eta \cdot \frac{\partial l(\theta^{\text{old}})}{\partial \theta_j} \\ &= \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^m [y^{(i)} - \sigma(\theta^T x^{(i)})] x_j^{(i)}\end{aligned}$$



Cross-Entropy Loss Function

- We need a loss function $L(\hat{y}, y)$ that expresses, for an observation x , how close the classifier output \hat{y} is to the correct output y (which is 0 or 1).
- A loss function that prefers the correct class labels of the training examples to be more likely. This is called conditional maximum likelihood estimation: we choose the parameters that maximize the log probability of the true y labels in the training data given the observations x .
- The resulting loss function is the negative log likelihood loss, generally called the cross-entropy loss
- Minimizing the negative of this function (minimizing the negative log likelihood) corresponds to maximizing the likelihood. This error function $L(\hat{y}, y)$ is typically known as the [cross-entropy error function](#) (also known as log-loss):



Cross entropy loss

$$p(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$$

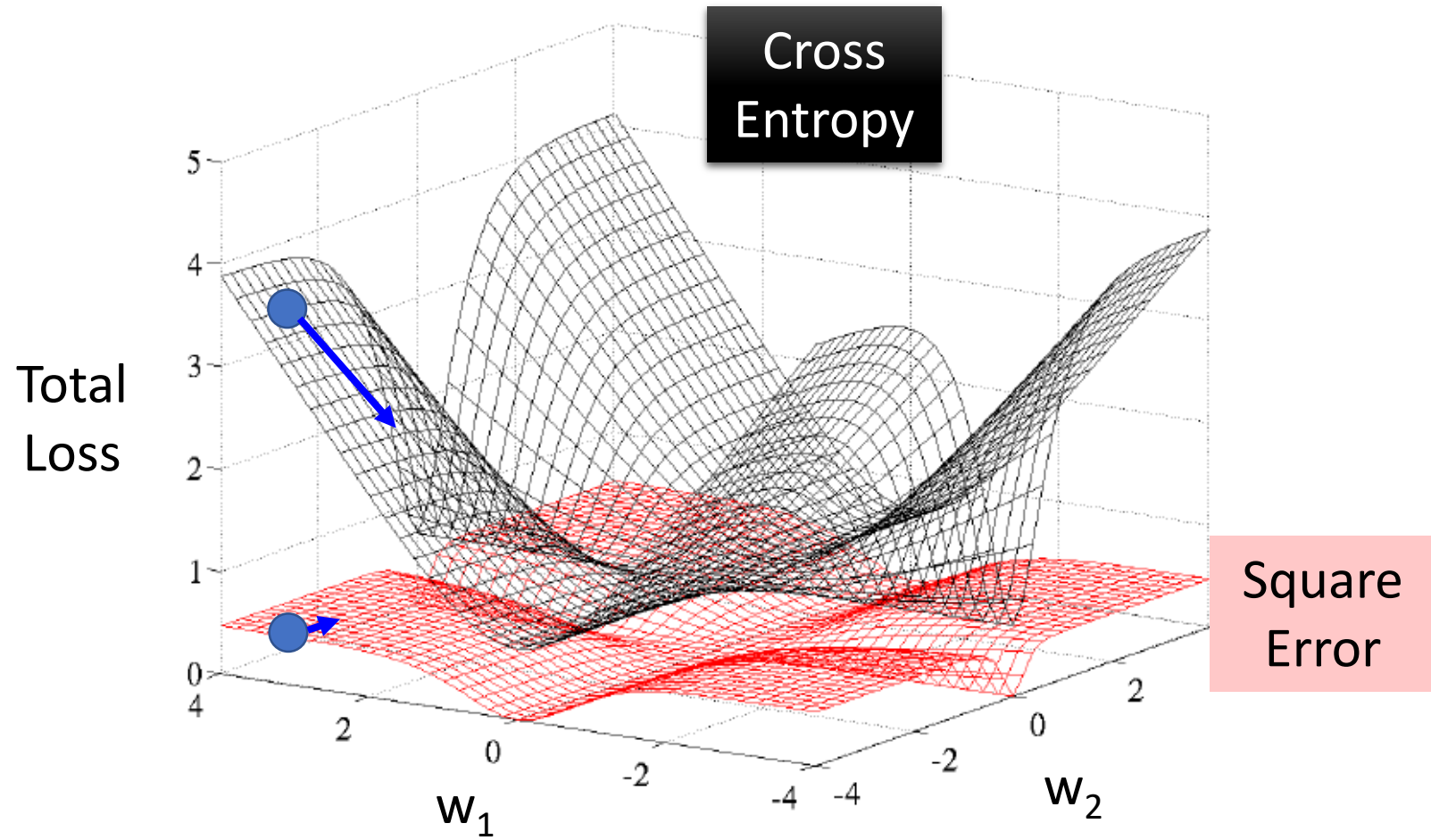
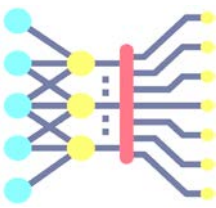
$$\log(p(y|x)) = y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \quad \text{log likelihood}$$

cross-entropy loss:

$$L_{CE}(\hat{y}, y) = -\log(p(y|x)) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})]$$

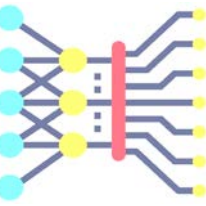
$$L_{CE}(\hat{y}, y) = -[y \log \sigma(\theta^T x) + (1 - y) \log(1 - \sigma(\theta^T x))]$$

Cross Entropy v.s. Square Error

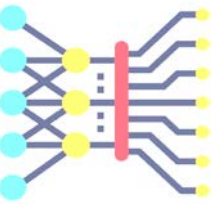


<http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>

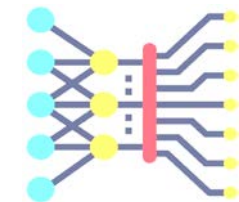
Gradient Descent



Multiclass Classification



Multi-class Classification



$$C_1: \theta_1, b_1 \quad z_1 = \theta_1 \cdot x + b_1$$

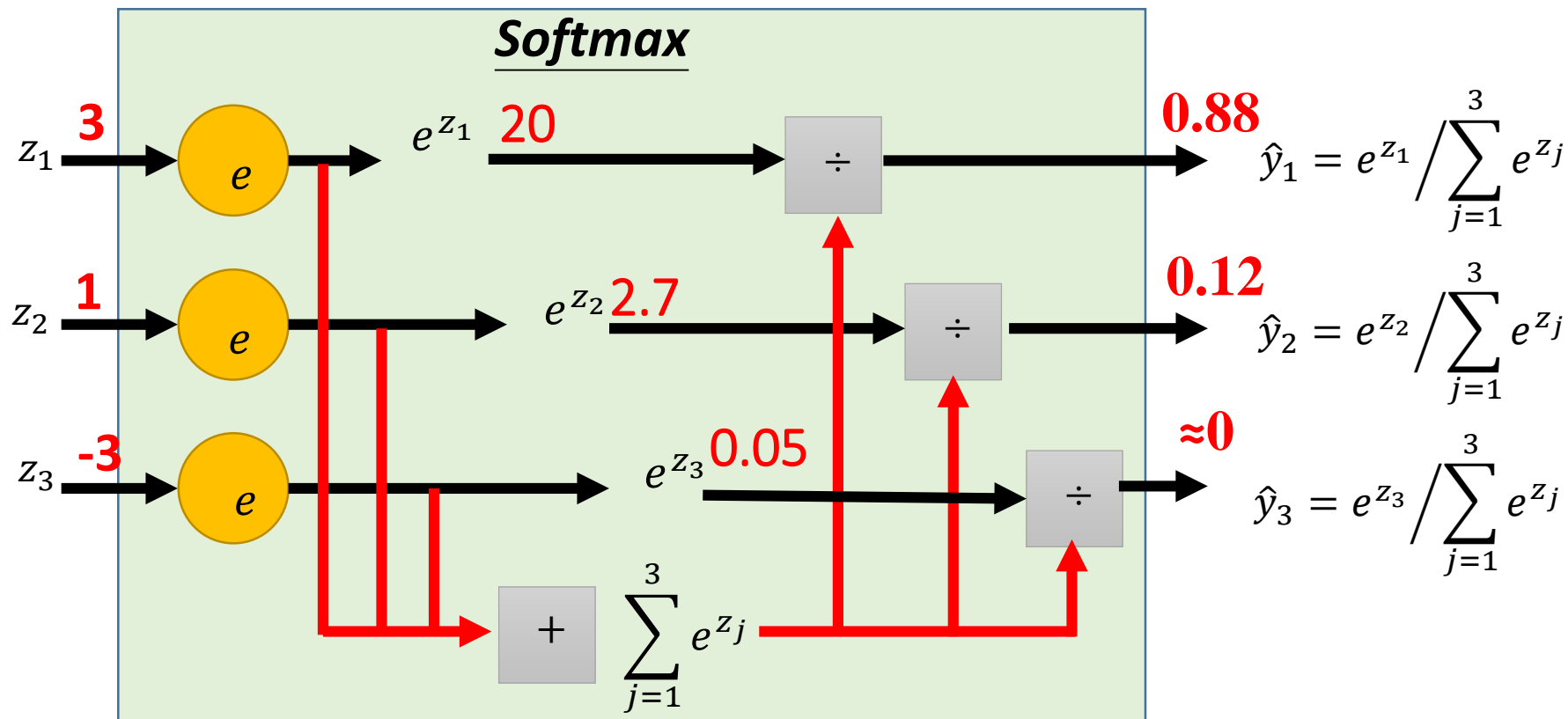
$$C_2: \theta_2, b_2 \quad z_2 = \theta_2 \cdot x + b_2$$

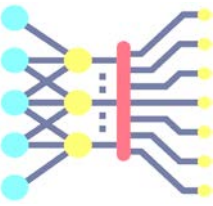
$$C_3: \theta_3, b_3 \quad z_3 = \theta_3 \cdot x + b_3$$

Probability:

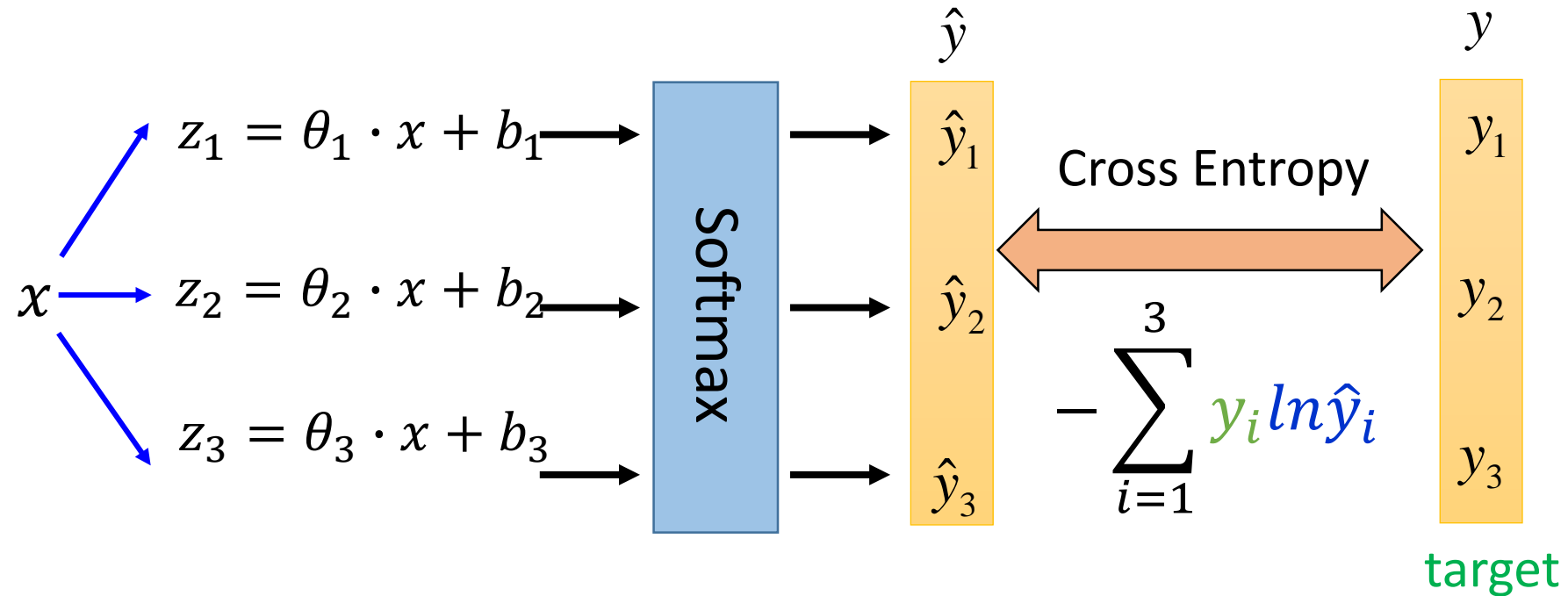
- $1 > \hat{y}_i > 0$
- $\sum_i \hat{y}_i = 1$

$$y_i = P(C_i|x)$$





Multi-class Classification



If $x \in$ class 1

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$-\ln \hat{y}_1$$

If $x \in$ class 2

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$-\ln \hat{y}_2$$

If $x \in$ class 3

$$y = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$-\ln \hat{y}_3$$