Open Elective Course [OE]
Course Code: CSO507
Winter 2023-24

Lecture#

# Deep Learning

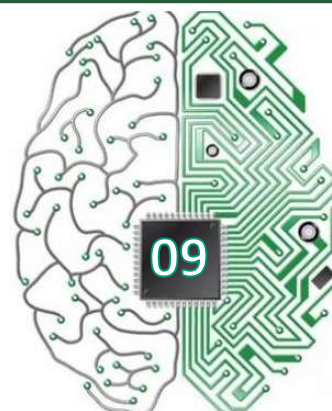## Unit-2: Linear and Logistic Regression (Part-II)

09

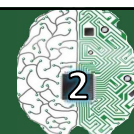Course Instructor:

Dr. Monidipa Das

Assistant Professor

Department of Computer Science and Engineering

Indian Institute of Technology (Indian School of Mines) Dhanbad, Jharkhand 826004, India
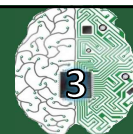
---

# Supervised Learning [revisited]

2

- Given a set of data points $\{x^{(1)}, x^{(2)}, ...., x^{(n)}\}$ associated to a set of outcomes $\{y^{(1)}, y^{(2)}, ...., y^{(n)}\}$ , we want to build a model that learns how to predict $y$ from $x$.

**Type of prediction** — The different types of predictive models are summed up in the table below:

|  | Regression | Classification |
|---|---|---|
| **Outcome** | Continuous | Class |
| **Examples** | Linear regression | Logistic regression, SVM, Naive Bayes |

# Classification: Example

**Email:** Spam / Not Spam?
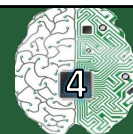**Online Transactions:** Fraudulent (Yes / No)?
**Tumor:** Malignant / Benign ?

$y = \{0, 1\}$     0: "Negative Class" (e.g., benign tumor)   ← **Two-class/Binary Classification**
             1: "Positive Class" (e.g., malignant tumor)

$y = \{0, 1, 2, 3\}$   0: "SMALL"
                1: "MEDIUM"   ← **Multiclass Classification**
                2: "LARGE"
                3: "EXTRA LARGE"

# Classification: Task Description

Given:
- Data $\boldsymbol{X} = \left\{ \boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(n)} \right\}$ where $\boldsymbol{x}^{(i)} \in \mathbb{R}^d$
- Corresponding labels $\boldsymbol{y} = \left\{ y^{(1)}, \ldots, y^{(n)} \right\}$ where $y^{(i)} \in \{0, \ldots, k\}$

$\qquad\qquad\qquad\qquad k = 1$ for Two-class/Binary Classification

**Can the task be performed by Linear Regression?**



Threshold classifier output $h_\theta(x)$ at 0.5:
    If $h_\theta(x) \geq 0.5$, predict "y = 1"
    If $h_\theta(x) < 0.5$, predict "y = 0"

# Classification: Task Description

5

Given:

- Data $X = \left\{ x^{(1)}, \ldots, x^{(n)} \right\}$ where $x^{(i)} \in \mathbb{R}^d$
- Corresponding labels $y = \left\{ y^{(1)}, \ldots, y^{(n)} \right\}$ where $y^{(i)} \in \{0, \ldots, k\}$

$k = 2$ for Two-class/Binary Classification

**Can the task be performed by Linear Regression?**

(Yes) 1

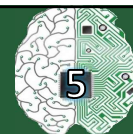Malignant ?     0.5

Threshold classifier output $h_\theta(x)$ at 0.5:

If $h_\theta(x) \geq 0.5$, predict "y = 1"

If $h_\theta(x) < 0.5$, predict "y = 0"

(No) 0

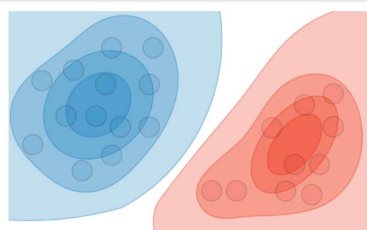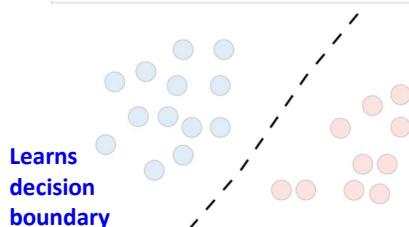$h_\theta(x) = \theta^T x$    Tumor Size

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad
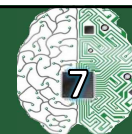
# Logistic Regression

6

- Takes a **probabilistic approach to learning discriminative functions** (i.e., a classifier)

- **Classification based on Probability:** Instead of just predicting the class, give the probability of the instance being in that class. **_Two_** key models:

| Discriminative model | Generative model |
| --- | --- |
| Directly estimate $P(y\|x)$ | Estimate $P(x\|y)$ to then deduce $P(y\|x)$ |

**Learns decision boundary**

**Learns the probability distributions of the data**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Generative vs. Discriminative

**Training Samples**

**Test Sample**
**What is this?**

It's a dog....because dogs have folded ears and they wear collars!

It's a dog.......because it fits well with my generated dog image!

| Discriminative model |
|---|
| Directly estimate $P(y|x)$ |

| Generative model |
|---|
| Estimate $P(x|y)$ to then deduce $P(y|x)$ |

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad
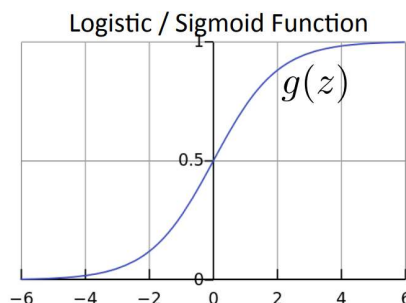
# Logistic Regression: Model Representation

- Takes a **probabilistic approach to learning discriminative functions** (i.e., a classifier)

- $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ should give $p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$
  - Want $0 \le h_{\boldsymbol{\theta}}(\boldsymbol{x}) \le 1$

  Can't just use linear regression with a threshold

- Logistic regression model:

  $$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}\right)$$

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

  $$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}}}$$

Logistic / Sigmoid Function

$g(z)$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Interpretation of Hypothesis Output

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$ = estimated $p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

Example: Cancer diagnosis from tumor size

$$\boldsymbol{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$
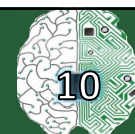
$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant

Note that: $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) + p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1$

Therefore, $p(y = 0 \mid \boldsymbol{x}; \boldsymbol{\theta}) = 1 - p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

# Logistic Regression: Hypothesis

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}\right)$$

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

$g(z)$

$\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}$ should be large <u>negative</u> values for negative instances

$\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}$ should be large <u>positive</u> values for positive instances

- Assume a threshold and...
  - Predict $y = 1$ if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0.5$
  - Predict $y = 0$ if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5$

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

Predict "$y = 1$" if $-3 + x_1 + x_2 \geq 0$
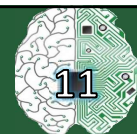
**Decision boundary:** $x_1 + x_2 = 3$

$\mathsf{x}_2$

$\boldsymbol{y = 1}$
$x_1 + x_2 \geq 3$

$\boldsymbol{y = 0}$
$x_1 + x_2 < 3$

**At decision boundary:**
$x_1 + x_2 = 3$
So, $h_{\boldsymbol{\theta}}(\boldsymbol{x})$=0.5

$\mathsf{x}_1$

# Non-Linear Decision Boundary

- Can apply basis function expansion to features, same as with linear regression

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2) \qquad \theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

x_2

$y = 1$   ✖   ✖   ✖   ✖ $y = 1$

✖   ○○○   ✖

y = 0

-1   1   ✖   x_1

✖ ○○○○ ✖

$y = 1$   ✖   -1   ✖ $y = 1$

✖ ✖

**Decision boundary:** $x_1^2 + x_2^2 = 1$

Predict "$y = 1$" if $-1 + x_1^2 + x_2^2 \geq 0$

# Logistic Regression: Cost Function

- Given $\left\{ \left( \boldsymbol{x}^{(1)}, y^{(1)} \right), \left( \boldsymbol{x}^{(2)}, y^{(2)} \right), \ldots, \left( \boldsymbol{x}^{(n)}, y^{(n)} \right) \right\}$

  where $\boldsymbol{x}^{(i)} \in \mathbb{R}^d$, $y^{(i)} \in \{0, 1\}$

- Model: $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left( \boldsymbol{\theta}^\mathsf{T} \boldsymbol{x} \right)$
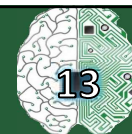
  $$g(z) = \frac{1}{1 + e^{-z}}$$

  $$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \qquad \boldsymbol{x}^\mathsf{T} = \begin{bmatrix} 1 & x_1 & \ldots & x_d \end{bmatrix}$$

- **How to choose parameters?**

# Logistic Regression: Cost Function

13

**Logistic regression objective**:

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

**Why not** $J(\theta) = \left( \frac{1}{1 + e^{-\theta x}} - y \right)^2$ **?**

# Intuition Behind the Objective

14

Aside: Recall the plot of $\log(z)$

$$\text{cost}\left(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y\right) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

If $y$ = 1

- Cost = 0 if prediction is correct
- As $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \to 0, \text{cost} \to \infty$

- Captures intuition that larger mistakes should get larger penalties
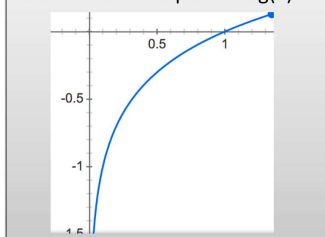  - e.g., predict $h_{\boldsymbol{\theta}}(x) = 0$, but $y$ = 1

If $y$ = 1

cost

0          $h_{\boldsymbol{\theta}}(\boldsymbol{x})$          1

# Intuition Behind the Objective

Aside: Recall the plot of $\log(z)$

$$\text{cost}\,(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$
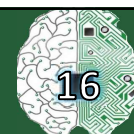
If $y = 0$

- Cost = 0 if prediction is correct
- As $(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \to 0, \text{cost} \to \infty$
- Captures intuition that larger mistakes should get larger penalties

If $y = 1$
If $y = 0$

cost

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$

0          1

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Intuition Behind the Objective

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

- Cost of a single instance:

$$\text{cost}\,(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$
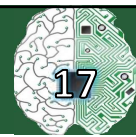
- Can re-write objective function as

$$J(\boldsymbol{\theta}) = \sum_{i=1}^{n} \text{cost}\left( h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}), y^{(i)} \right)$$

Compare to linear regression: $J(\boldsymbol{\theta}) = \dfrac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}\left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$

12

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Regularized Logistic Regression

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n}\left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log\left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right)\right]$$

- We can regularize logistic regression exactly as before:

$$J_{\text{regularized}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda \sum_{j=1}^{d} \theta_j^2$$
$$= J(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

# Gradient Descent for Logistic Regression

$$J_{\text{reg}}(\boldsymbol{\theta}) = -\sum_{i=1}^{n}\left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log\left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right)\right] + \lambda\|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want $\min\limits_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$
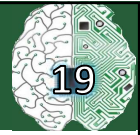- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update for $j = 0 \dots d$

Use the natural logarithm (ln = $\log_e$) to cancel with the exp() in $h_{\boldsymbol{\theta}}(\boldsymbol{x})$

# Gradient Descent for Logistic Regression

$$J_{\mathrm{reg}}(\boldsymbol{\theta}) = -\sum_{i=1}^{n}\left[y^{(i)}\log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right)\log\left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right)\right] + \lambda\|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

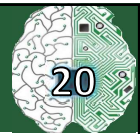Want $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence        (simultaneous update for $j = 0 \ldots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha\sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)$$

$$\theta_j \leftarrow \theta_j - \alpha\left[\sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)x_j^{(i)} + \frac{\lambda}{n}\theta_j\right]$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Gradient Descent for Logistic Regression

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence        (simultaneous update for $j = 0 \ldots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha\sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)$$

$$\theta_j \leftarrow \theta_j - \alpha\left[\sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)x_j^{(i)} + \frac{\lambda}{n}\theta_j\right]$$

This looks IDENTICAL to linear regression!!!
- Ignoring the $1/n$ constant
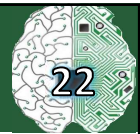- However, the form of the model is very different:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\top}\boldsymbol{x}}}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad
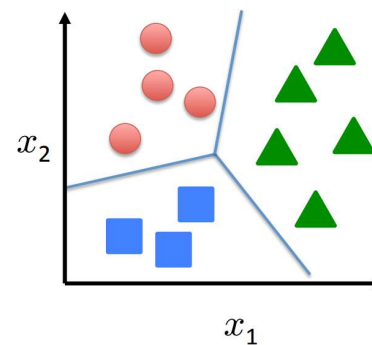
# Multi-Class Classification

---

# Multi-Class Classification

Binary classification:          Multi-class classification:



Disease diagnosis:      healthy / cold / flu / pneumonia
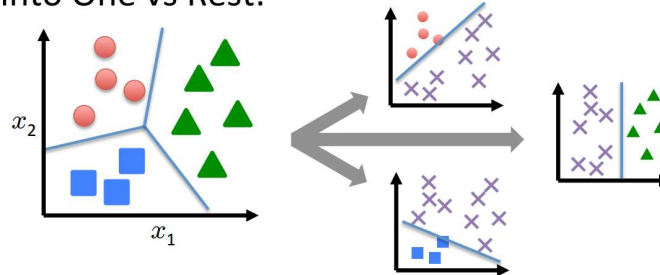
Object classification:  desk / chair / monitor / bookcase

# Multi-Class Logistic Regression
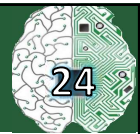
Split into One vs Rest:



- Train a logistic regression classifier for each class $i$ to predict the probability that $y = i$ with

$$h_c(\boldsymbol{x}) = \frac{\exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}{\sum_{c=1}^{C} \exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Multi-Class Logistic Regression

- For 2 classes:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})} = \frac{\exp(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})}{\boxed{1} + \boxed{\exp(\boldsymbol{\theta}^\mathsf{T} \boldsymbol{x})}}$$

weight assigned to $y = 0$       weight assigned to $y = 1$

- For $C$ classes $\{1, \ldots, C\}$:

$$p(y = c \mid \boldsymbol{x}; \boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_C) = \frac{\exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}{\sum_{c=1}^{C} \exp(\boldsymbol{\theta}_c^\mathsf{T} \boldsymbol{x})}$$

  – Called the **softmax** function

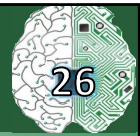Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

## Implementing Multi-Class Logistic Regression

- Use $h_c(\boldsymbol{x}) = \dfrac{\exp(\boldsymbol{\theta}_c^{\mathsf{T}} \boldsymbol{x})}{\sum_{c=1}^{C} \exp(\boldsymbol{\theta}_c^{\mathsf{T}} \boldsymbol{x})}$ as the model for class $c$

- Gradient descent simultaneously updates all parameters for all models
  - Same derivative as before, just with the above $h_c(\boldsymbol{x})$

- Predict class label as the most probable label

$$\max_c h_c(\boldsymbol{x})$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Questions?

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad