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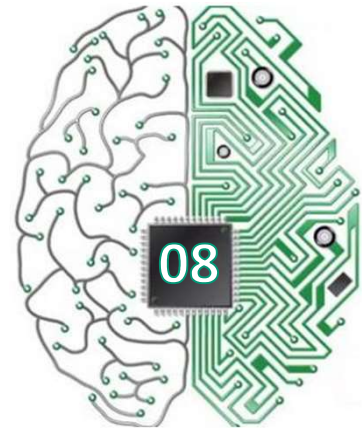**Course Code: CSO507**
**Winter 2023-24**

**Lecture#**

# Deep Learning

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

08

**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**
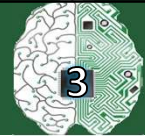
---

# Supervised Learning

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 $x$.

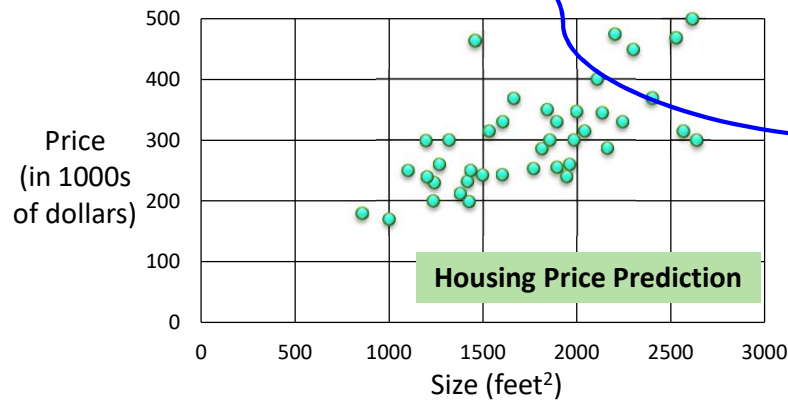**Type of prediction** — The different types of predictive models are summed up in the table below:

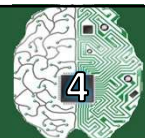|  | Regression | Classification |
|---|---|---|
| **Outcome** | Continuous | Class |
| **Examples** | Linear regression | Logistic regression, SVM, Naive Bayes |

# Regression: Task Description

③

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 $y^{(i)} \in \mathbb{R}$



**Housing Price Prediction**

Price (in 1000s of dollars) vs Size (feet²)

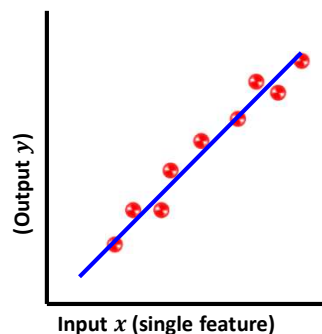| Living area (feet²) | Price (1000$s) |
|---|---|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| ⋮ | ⋮ |

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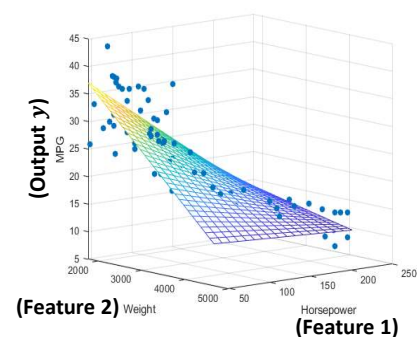
# Linear Regression Model

④

▪**Linear regression is like fitting a line or (hyper)plane to a set of points**

- **Univariate linear regression:** A single independent variable is used to predict

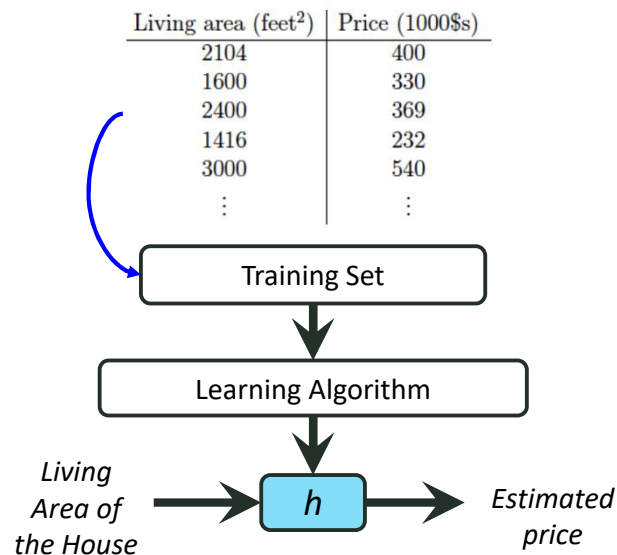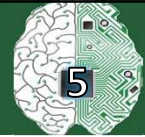- **Multivariate linear regression:** Two or more independent variables are used to predict



**(Output $y$)** vs **Input $x$ (single feature)**

**(Output $y$)** MPG vs **(Feature 2)** Weight, **(Feature 1)** Horsepower

However, we can even fit a curve using a linear model after suitably transforming the inputs
$$y \approx w^\top \phi(x)$$

The transformation $\phi(.)$ can be predefined or learned (e.g., using kernel methods or a deep neural network based feature extractor).
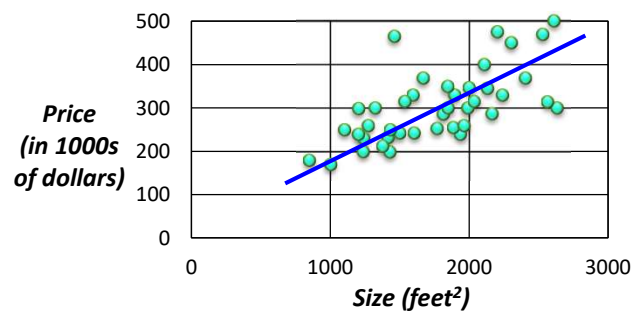
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

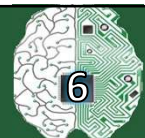# Model Representation:
# Univariate Linear Regression

5

| Living area (feet$^2$) | Price (1000$s) |
|---|---|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| ⋮ | ⋮ |

Training Set

Learning Algorithm

*Living Area of the House* → *h* → *Estimated price*

**How do we represent *h* ?**

**Hypothesis:** $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$'s: **Parameters**
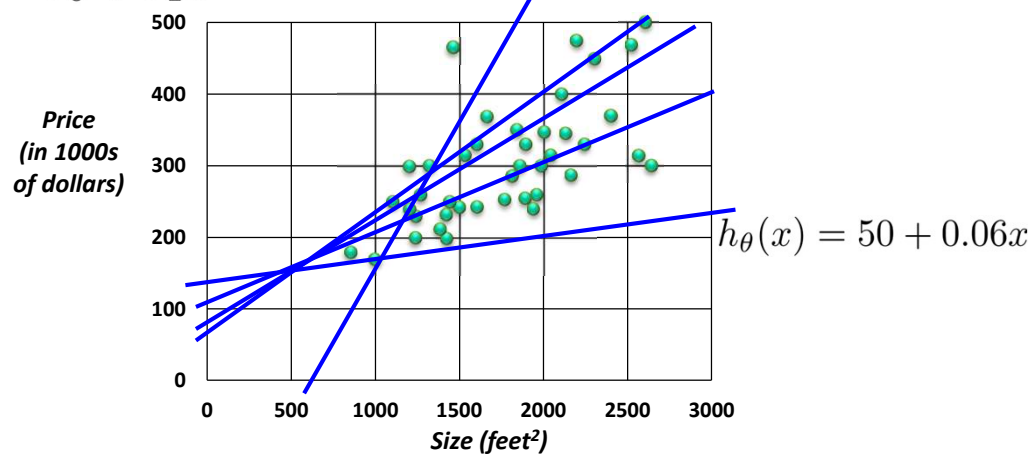
*Price (in 1000s of dollars)*

*Size (feet$^2$)*

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

---

# How to choose $\theta_i$'s ?

6

$$h_\theta(x) = \theta_0 + \theta_1 x$$
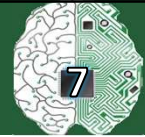
*Price (in 1000s of dollars)*
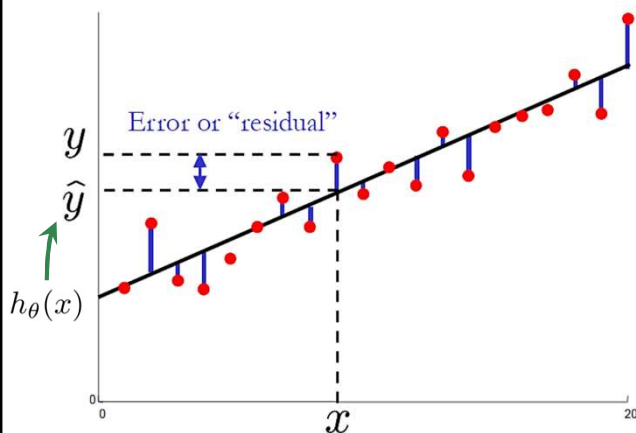
$$h_\theta(x) = 50 + 0.06x$$

*Size (feet$^2$)*

**Idea:** Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training examples $(x, y)$

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

# Cost Function

**How good is the prediction given by the straight line?**



$y$

$\widehat{y}$

$h_\theta(x)$

Error or "residual"

$x$

**Hypothesis:** $h_\theta(x) = \theta_0 + \theta_1 x$

**Parameters:** $\theta_0, \theta_1$

**Cost Function:**

$\widehat{y}^{(i)}$

$$J(\theta_0, \theta_1) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

**Goal:** $\underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

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

# Cost Function: Contour Plot

*Price (in 1000s of dollars)*

*Size (feet²)*

$J(\theta_0, \theta_1)$

$\theta_1$

$\theta_0$

**Goal:** $\underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

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

# Optimization using Gradient Descent

**Gradient descent algorithm**

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$
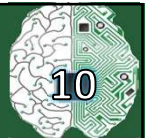
(for $j = 1$ and $j = 0$)

}

Univariate Linear Regression Model

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$J(\theta_0, \theta_1) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

# GD for Univariate Linear Regression Model
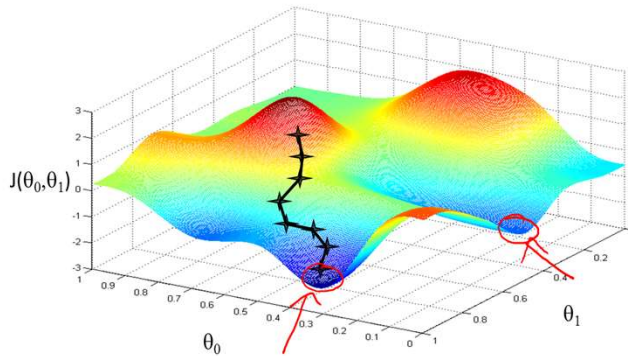
**Gradient descent algorithm**

repeat until convergence {

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

$$\theta_1 := \theta_1 - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$$

}

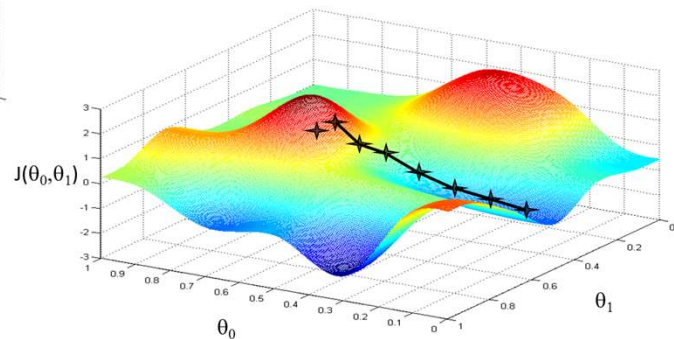update $\theta_0$ and $\theta_1$ simultaneously

# Non-Convex Cost Function
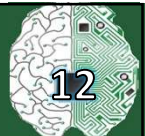
- **Initialization of $\theta$ matters**



Courtesy: Andrew Ng

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

# Multiple Features (variables)

| Size (feet²) | Number of bedrooms | Number of floors | Age of home (years) | Price ($1000) |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

$x_3^{(2)}$

$y^{(2)}$

**Notations:**

$n$ : Number of training samples

$d$ : Number of features. E.g. $d = 4$ in the above example

$x^{(i)}$: training example $i$.

$y^{(i)}$: Label/target for training example $i$.

$x_j^{(i)}$: value of feature $j$ in training example $i$.

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

## Multivariate Linear Regression: **Hypothesis**

- Hypothesis for univariate linear regression:

$$h_\theta(x) = \theta_0 + \theta_1 x$$

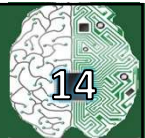- Hypothesis for **multivariate linear regression**:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_d x_d$$

*For convenience of notation, define $x_0 = 1$ for all sample i* $\qquad h_\theta(x^{(i)}) = \sum_{k=0}^{d} \theta_k x_k^{(i)}$

## Multivariate Linear Regression: **Cost Function**

Hypothesis: $\quad h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_d x_d$
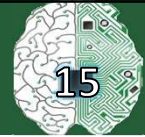
Parameters: $\quad \theta_0, \theta_1, \ldots, \theta_d$

Cost function:

$$J(\theta_0, \theta_1, \ldots, \theta_d) = \frac{1}{2n} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)})^2$$

## Multivariate Linear Regression: **Gradient Descent**

Gradient descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_d) \qquad \equiv J(\boldsymbol{\theta})$$

}  (simultaneously update for every $j = 0, \ldots, d$ )

Repeat {

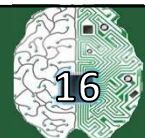$$\theta_j := \theta_j - \alpha \frac{1}{n} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

simultaneously update $\theta_j$ for $j = 0, \ldots, d$

}

**For Linear Regression**

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}$$

$$\theta_0 := \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{n} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{n} \sum_{i=1}^{n} (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...
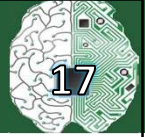
## Linear Regression: More Complex Models

- The inputs X for linear regression can be:
  - Original quantitative inputs: $x$
  - Transformation of quantitative inputs
    - **example:** $\log(x)$, $\exp(x)$, $\sqrt{x}$, $x^2$ etc.
  - Polynomial transformation
    - **example:** $y = \beta 0 + \beta 1 \cdot x + \beta 2 \cdot x^2 + \beta 3 \cdot x^3$
  - Basis expansions
  - Dummy coding of categorical inputs:
    - **example:** R: [1 0 0], G: [0 1 0], B: [0 0 1]
  - Interactions between variables
    - **example:** x3 = x1 · x2

- **These allows use of linear regression techniques to fit non-linear datasets.**

# Linear Basis Function Models

**17**

- Generally,

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{j=0}^{d} \theta_j \underbrace{\phi_j(\boldsymbol{x})}_{\text{basis function}}$$

- Typically, $\phi_0(\boldsymbol{x}) = 1$ so that $\theta_0$ acts as a bias
- In the simplest case, we use linear basis functions :

$$\phi_j(\boldsymbol{x}) = x_j$$
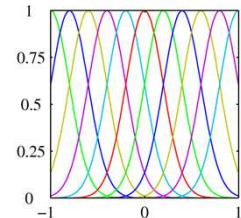
- Polynomial basis functions:

$$\phi_j(x) = x^j$$

- Gaussian basis functions:

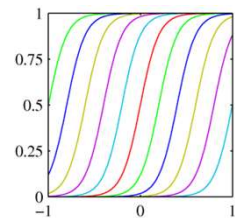$$\phi_j(x) = \exp\left\{-\frac{(x-\mu_j)^2}{2s^2}\right\}$$

- Sigmoidal basis functions:

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$$

where

$$\sigma(a) = \frac{1}{1+\exp(-a)}$$

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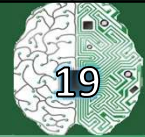
---

# Linear Basis Function Models

**18**

- Basic Linear Model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{j=0}^{d} \theta_j x_j$$

- Generalized Linear Model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{j=0}^{d} \theta_j \phi_j(\boldsymbol{x})$$
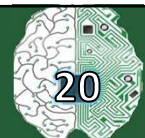
- Once we have replaced the data by the outputs of the basis functions, fitting the generalized model is exactly the same problem as fitting the basic model
  - Unless we use the kernel trick
  - Therefore, there is no point in cluttering the math with basis functions

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

# Model Representation through Vectorization

# Vectorization

20

- Benefits of vectorization
  - More compact equations
  - Faster code (using optimized matrix libraries)

- Consider our model::

$$h(\boldsymbol{x}) = \sum_{j=0}^{d} \theta_j x_j$$

- Let

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

- Can write the model in vectorized form as $h(\boldsymbol{x}) = \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}$

# Vectorization

- Consider our model for $n$ instances:

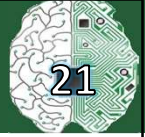$$h\left(\boldsymbol{x}^{(i)}\right) = \sum_{j=0}^{d} \theta_j x_j^{(i)}$$
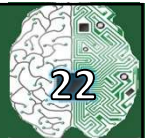
- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \qquad \boldsymbol{X} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \cdots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix}$$

$$\mathbb{R}^{(d+1)\times 1} \qquad\qquad \mathbb{R}^{n\times(d+1)}$$

- Can write the model in vectorized form as $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = \boldsymbol{X}\boldsymbol{\theta}$

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

# Vectorization

- For the linear regression cost function:

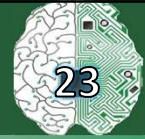$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)} \right)^2$$

$$= \frac{1}{2n} \sum_{i=1}^{n} \left( \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{2n} \left( \boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y} \right)^{\mathsf{T}} \left( \boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y} \right)$$

$$\mathbb{R}^{n\times(d+1)}$$
$$\mathbb{R}^{(d+1)\times 1}$$
$$\mathbb{R}^{1\times n} \qquad \mathbb{R}^{n\times 1}$$

Let:

$$\boldsymbol{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$
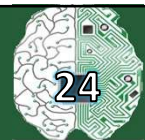
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# GD vs. Closed Form Solution for Linear Regression

---

# Closed Form Solution

- Instead of using GD, solve for optimal $\boldsymbol{\theta}$ analytically
  - Notice that the solution is when $\dfrac{\partial}{\partial \boldsymbol{\theta}} J(\boldsymbol{\theta}) = 0$

- Derivation:

$$\mathcal{J}(\boldsymbol{\theta}) = \frac{1}{2n} (\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y})^{\mathsf{T}} (\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y})$$

1 x 1

$$\propto \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\boldsymbol{\theta} - \boxed{\boldsymbol{y}^{\mathsf{T}} \boldsymbol{X}\boldsymbol{\theta}} - \boxed{\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y}} + \boldsymbol{y}^{\mathsf{T}} \boldsymbol{y}$$

$$\propto \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\boldsymbol{\theta} - 2\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y} + \boldsymbol{y}^{\mathsf{T}} \boldsymbol{y}$$

Take derivative and set equal to 0, then solve for $\boldsymbol{\theta}$:

$$\frac{\partial}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\boldsymbol{\theta} - 2\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y} + \boldsymbol{y}^{\mathsf{T}} \boldsymbol{y}) = 0$$

$$(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X})\boldsymbol{\theta} - \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y} = 0$$

$$(\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X})\boldsymbol{\theta} = \boldsymbol{X}^{\mathsf{T}} \boldsymbol{y}$$
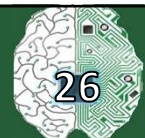
# Closed Form Solution vs. GD

- Can obtain $\theta$ by simply plugging $X$ and $y$ into

$$\theta = (X^{\mathsf{T}} X)^{-1} X^{\mathsf{T}} y$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \cdots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \cdots & x_d^{(n)} \end{bmatrix} \qquad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

| Gradient Descent | Closed Form Solution |
|---|---|
| • Requires multiple iterations<br>• Need to choose $\alpha$<br>• Works well when $n$ is large<br>• Can support incremental learning | • Non-iterative<br>• No need for $\alpha$<br>• Slow if $n$ is large<br>   – Computing $(X^{\mathsf{T}} X)^{-1}$ is roughly O($n^3$) |

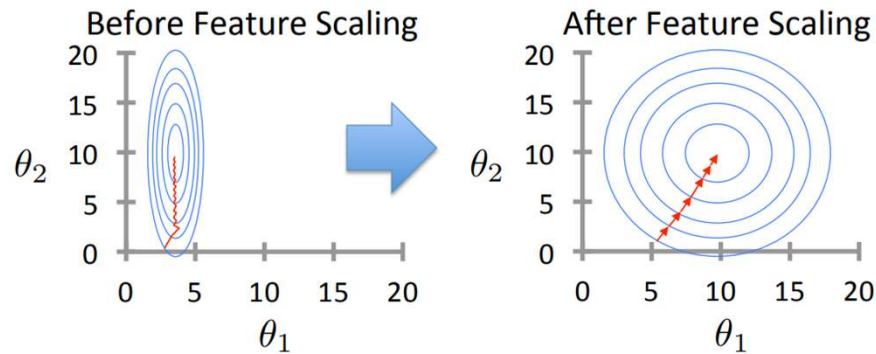Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Improving Learning

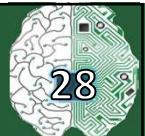Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

# Feature Scaling

- **Idea:** Ensure that feature have similar scales



- Makes gradient descent converge much faster

# Feature Standardization

- Rescales features to have zero mean and unit variance
  - Let $\mu_j$ be the mean of feature $j$: $\quad \mu_j = \dfrac{1}{n}\sum_{i=1}^{n} x_j^{(i)}$
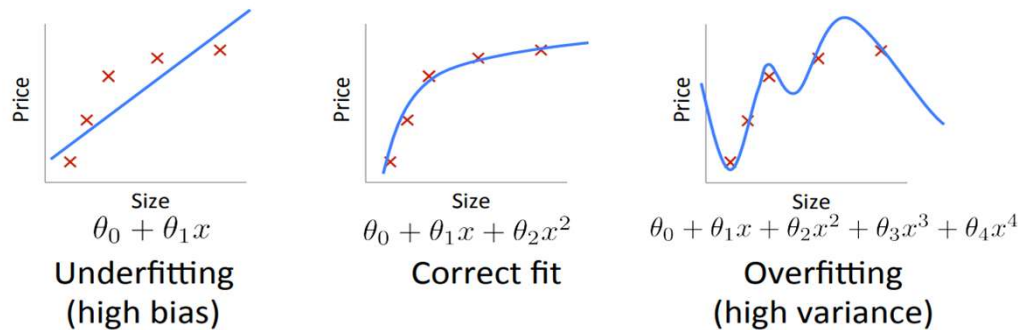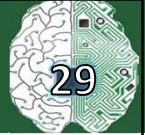
  - Replace each value with:

  $$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j} \qquad \begin{array}{l}\text{for } j = 1...d \\ (\text{not } x_0!)\end{array}$$

    - $s_j$ is the standard deviation of feature $j$
    - Could also use the range of feature $j$ ($\max_j - \min_j$) for $s_j$

- Must apply the same transformation to instances for both training and prediction
- Outliers can cause problems

# Quality of Fit

$$\theta_0 + \theta_1 x$$

Underfitting
(high bias)

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

Correct fit

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

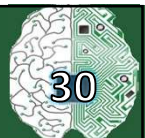Overfitting
(high variance)

**Overfitting:**

- The learned hypothesis may fit the training set very well ( $J(\boldsymbol{\theta}) \approx 0$ )

- ...but fails to generalize to new examples
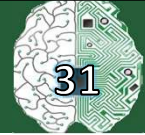
# Regularization

- A method for automatically controlling the complexity of the learned hypothesis

- **Idea**: penalize for large values of $\theta_j$
  - Can incorporate into the cost function
  - Works well when we have a lot of features, each that contributes a bit to predicting the label

- Can also address overfitting by eliminating features (either manually or via model selection)

# Regularization

- Linear regression objective function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{d} \theta_j^2$$

$\underbrace{\qquad\qquad}_{\text{model fit to data}}$ $\underbrace{\qquad}_{\text{regularization}}$

- $\lambda$ is the regularization parameter ($\lambda \geq 0$)
- No regularization on $\theta_0$!

# Understanding Regularization

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{d} \theta_j^2$$

- Note that $\displaystyle\sum_{j=1}^{d} \theta_j^2 = \|\boldsymbol{\theta}_{1:d}\|_2^2$
  - This is the magnitude of the feature coefficient vector!

**Let $d = 4$**
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

- We can also think of this as:

$$\sum_{j=1}^{d} (\theta_j - 0)^2 = \|\boldsymbol{\theta}_{1:d} - \vec{\mathbf{0}}\|_2^2$$

- $L_2$ regularization pulls coefficients toward 0

**What happens if we set $\lambda$ to be huge (e.g., $10^{10}$)?** $\theta_1, \theta_2, \theta_3, \theta_4$ **becomes 0**

*Price (in 1000s of dollars)*

500
400
300
200
100
0

0    1000    2000    3000
*Size (feet²)*

# Regularized Linear Regression

- Cost Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \left[ \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{d} \theta_j^2 \right]$$

- Fit by solving $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- We can rewrite the gradient step as:
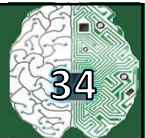
$$\theta_j \leftarrow \theta_j \left( 1 - \alpha \frac{\lambda}{n} \right) - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)}$$

- Gradient update:

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

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

regularization

# Various Ways of Regularization

- $$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{d} \theta_j^2$$

  $\lambda ||\theta||_2^2$

  $\ell_2$ regularization

  **Ridge Regression**

- $$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2 + \lambda \sum_{j=1}^{d} |\theta_j|$$

  $\lambda ||\theta||_1$

  $\ell_1$ regularization

  **Lasso Regression**

| LASSO | Ridge |
|---|---|
| • Shrinks coefficients to 0<br>• Good for variable selection | Makes coefficients smaller |

Least Absolute Shrinkage and Selection Operator (LASSO)
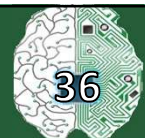
## Non-regularization Approaches for Overfitting

- *Early-stopping* (stopping training just when we have a decent validation set accuracy)

- *Injecting noise* in the inputs

- *Dropout* (in each iteration, don't update some of the weights) [e.g. used in deep network-based models]

# Questions?