

Examples:  $m = 4$ .

	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$   
 $m \times (n+1)$

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$   
 $m$ -dimensional vector

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

# Matrix form to be fitted

- For convenience of writing the expressions introduce zeroth variable  $x_0$ ; which always takes the value 1
- Hypothesis function is  $y = y(x_1, x_2, \dots, x_n) = b_0x_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$ ;

- $\mathbf{b} = \begin{pmatrix} b_0 \\ b_1 \\ \dots \\ b_n \end{pmatrix} \in \mathcal{R}^{n+1}$ ,  $\mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{pmatrix} \in \mathcal{R}^{n+1}$ ;  $\mathbf{b}^T = (b_0, b_1, \dots, b_n)$  is  $1 \times (n+1)$  vector and  $\mathbf{x}$  is  $(n+1) \times 1$  vector

- $\mathbf{b}^T \mathbf{x}$  is  $1 \times 1$  vector  $:= (b_0, b_1, \dots, b_n) \cdot \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{pmatrix}$

- $= b_0x_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$
- $y = \mathbf{b}^T \mathbf{x}$  is the hypothesis function

$\mathbf{x}_0$	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\dots$	$\mathbf{x}_n$	$\mathbf{y}$
1	$x_{11}$	$x_{21}$	$x_{31}$	$\dots$	$x_{n1}$	$y_1$
1	$x_{12}$	$x_{22}$	$x_{32}$	$\dots$	$x_{n2}$	$y_2$
1	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
1	$x_{1i}$	$x_{2i}$	$x_{3i}$	$\dots$	$x_{ni}$	$y_i$
1	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
1	$x_{1m}$	$x_{2m}$	$x_{3m}$	$\dots$	$x_{nm}$	$y_m$

# Least Square method

- $e_i = y_i - \hat{y}_i = y_i - \sum_j b_j x_{ji}$  for  $i = 1, 2, \dots, m$
- $e_i^2 = (y_i - \hat{y}_i)^2 = (y_i - \sum_j b_j x_{ji})^2$  for  $i = 1, 2, \dots, m$
- $SSE = \sum_i e_i^2 = \sum_i (y_i - \sum_j b_j x_{ji})^2$
- Observations in  $n+1$  dimensional space
- For least square fit, need to find  $b_j$  for  $j = 0, 1, \dots, n$  so that SSE is minimum. Denote SSE as  $S$ .
- Partial derivatives of  $S$  with respect to each of the parameters must be zero.

$$\frac{\partial S}{\partial b_j} = 0; \text{ for } j = 0, 1, \dots, n$$

# Normal Equations

- $\frac{\partial S}{\partial b_j} = 0$ ; for  $j = 0, 1, \dots, n$
- $S = \sum_i (y_i - \sum_j b_j x_{ji})^2$ ;
- $\frac{\partial S}{\partial b_j} = \frac{\partial}{\partial b_j} \sum_i (y_i - \sum_j b_j x_{ji})^2$
- $\frac{\partial S}{\partial b_j} = \sum_i \frac{\partial}{\partial b_j} (y_i - \sum_j b_j x_{ji})^2$
- $\frac{\partial S}{\partial b_j} = 2 \sum_i (y_i - \sum_j b_j x_{ji}) (-x_{ji}) = 0$  for  $j = 0, 1, \dots, n$
- $\sum_i (y_i - \sum_j b_j x_{ji}) (x_{ji}) = 0$  for  $j = 0, 1, \dots, n$  are **n+1 simultaneous linear equations in  $b_j$ 's for  $j = 0, 1, \dots, n$  called normal equations**
- For example for  $j = 0$ :
- $\sum_i y_i = b_0 \sum_i x_{0i} + b_1 \sum_i x_{1i} + b_2 \sum_i x_{2i} + \dots + b_n \sum_i x_{ni}$
-

# Normal Equations

- $3.8 = 5 b_0 + 40 b_1 + 45 b_2$
- $33 = 40 b_0 + 338 b_1 + 342 b_2$
- $27.9 = 45 b_0 + 342 b_1 + 549 b_2$
- Solve by elimination manually (H.W.)
- $$\begin{pmatrix} 5 & 40 & 45 \\ 40 & 338 & 342 \\ 45 & 342 & 549 \end{pmatrix} \begin{pmatrix} b_0 \\ b_1 \\ b_2 \end{pmatrix} = \begin{pmatrix} 3.8 \\ 33 \\ 27.9 \end{pmatrix}$$
- $$\begin{pmatrix} b_0 \\ b_1 \\ b_2 \end{pmatrix} = \begin{pmatrix} 5 & 40 & 45 \\ 40 & 338 & 342 \\ 45 & 342 & 549 \end{pmatrix}^{-1} \begin{pmatrix} 3.8 \\ 33 \\ 27.9 \end{pmatrix}$$

# Matrix Method to derive normal Equations

- $\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{21} & \dots & x_{n1} \\ 1 & x_{12} & x_{22} & \dots & x_{n2} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{1m} & x_{2m} & \dots & x_{nm} \end{bmatrix}$  is  $m \times (n+1)$  matrix;  $\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_m \end{pmatrix}$  is  $m \times 1$  vector
- $\mathbf{b} = \begin{pmatrix} b_0 \\ b_1 \\ \dots \\ b_n \end{pmatrix} \in \mathcal{R}^{n+1}$  is  $(n+1) \times 1$  vector;  $\hat{\mathbf{y}} = \begin{pmatrix} \widehat{y_1} \\ \widehat{y_2} \\ \dots \\ \widehat{y_m} \end{pmatrix}$  is  $m \times 1$  vector ;  $X^T$  is  $(n + 1) \times m$
- $X^T \mathbf{y}$  is  $((n + 1) \times m) \times (m \times 1)$  giving  $(n+1) \times 1 = \begin{pmatrix} k_0 \\ k_1 \\ \dots \\ k_n \end{pmatrix}$  (say)
- $\mathbf{b}^T = (b_0, b_1, \dots, b_n)$  is  $1 \times (n+1)$  vector;  $\mathbf{b}^T X^T \mathbf{y} = b_0 k_0 + b_1 k_1 + b_2 k_2 + \dots + b_n k_n$

# Predicted Values

- $b_0x_{01} + b_1x_{11} + b_2x_{21} + \dots + b_nx_{n1} = \widehat{y}_1$
- $b_0x_{02} + b_1x_{12} + b_2x_{22} + \dots + b_nx_{n2} = \widehat{y}_2$
- $\dots\dots\dots$
- $b_0x_{0i} + b_1x_{1i} + b_2x_{2i} + \dots + b_nx_{ni} = \widehat{y}_i$
- $\dots\dots\dots$
- $b_0x_{0m} + b_1x_{1m} + b_2x_{2m} + \dots + b_nx_{nm} = \widehat{y}_m$
- $b_0x_{0i} + b_1x_{1i} + b_2x_{2i} + \dots + b_nx_{ni} = \widehat{y}_i \text{ for } i = 1, 2, \dots, m$
- $\sum_j b_j x_{ji} = \widehat{y}_i \text{ for } i = 1, 2, \dots, m; j = 0, 1, \dots, n$

- $SSE = \sum_i e_i^2$ ;  $\mathbf{e} = \begin{pmatrix} e_1 \\ e_2 \\ \dots \\ e_m \end{pmatrix}$ ;  $\mathbf{e}^T = (e_1, e_2, \dots, e_m)$ ;  $\mathbf{e}^T \mathbf{e} = (e_1, e_2, \dots, e_m) \begin{pmatrix} e_1 \\ e_2 \\ \dots \\ e_m \end{pmatrix}$

- $\mathbf{e}^T \mathbf{e} = e_1^2 + e_2^2 + \dots + e_m^2 = \sum_i e_i^2 = SSE$ ;  $SSE = \mathbf{e}^T \mathbf{e}$

- $e_i = y_i - \hat{y}_i = y_i - \sum_j b_j x_{ji}$  for  $i = 1, 2, \dots, m$

- $X = \begin{bmatrix} 1 & x_{11} & x_{21} & \dots & x_{n1} \\ 1 & x_{12} & x_{22} & \dots & x_{n2} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{1m} & x_{2m} & \dots & x_{nm} \end{bmatrix}$ ;  $\mathbf{b} = \begin{pmatrix} b_0 \\ b_1 \\ \dots \\ b_n \end{pmatrix}$

- $X\mathbf{b} = \hat{y}_i$   $e_i = y_i - \hat{y}_i$  for  $i=1, 2, \dots, m$ ;  $\mathbf{e} = \mathbf{y} - X\mathbf{b}$



- $S = (\mathbf{y} - X\mathbf{b})^T (\mathbf{y} - X\mathbf{b}) = (\mathbf{y}^T - (X\mathbf{b})^T) (\mathbf{y} - X\mathbf{b})$
- $S = (\mathbf{y}^T - \mathbf{b}^T X^T) (\mathbf{y} - X\mathbf{b}) = (\mathbf{y}^T \mathbf{y} - \mathbf{y}^T X\mathbf{b} - \mathbf{b}^T X^T \mathbf{y} + \mathbf{b}^T X^T X\mathbf{b})$
- $\mathbf{y}^T X\mathbf{b}$  is a scalar  $((1 \times m)X(\textcolor{red}{m} \times (\textcolor{teal}{n} + 1))X((\textcolor{teal}{n} + 1) \times 1)) : 1 \times 1$
- $(\mathbf{y}^T X\mathbf{b})^T = \mathbf{b}^T X^T \mathbf{y} = \mathbf{y}^T X\mathbf{b}$  being scalar
- $S = \mathbf{y}^T \mathbf{y} - 2\mathbf{b}^T X^T \mathbf{y} + \mathbf{b}^T X^T X\mathbf{b}$
- $\mathbf{b}^T X^T \mathbf{y} = b_0 k_0 + b_1 k_1 + b_2 k_2 + \dots + b_n k_n$ ;  $k$ 's are constant with respect to  $b$ 's
- $\frac{\partial \mathbf{b}^T X^T \mathbf{y}}{\partial b_j} = k_j ;$

- Contribution of  $\mathbf{b}^T X^T X \mathbf{b}$  is  $b_0^2 d_0 + b_1^2 d_1 + b_2^2 d_2 + \dots + b_n^2 d_n$
- $\frac{\partial \mathbf{b}^T X^T X \mathbf{b}}{\partial b_j} = 2b_j d_j$
- $\frac{\partial S}{\partial b_j} = -2 k_j + 2b_j d_j$
- S is scalar and a function of  $b_0, b_1, \dots, b_n$
- $\nabla f(x, y, z) = (f_x, f_y, f_z)$  called Grad f , gradient of a function
- In order that S is minimum , necessary condition is  $\nabla S = \mathbf{0}$  (zero vector)

- $\nabla S = \frac{\partial S}{\partial \mathbf{b}} = \begin{pmatrix} \frac{\partial S}{\partial b_0} \\ \frac{\partial S}{\partial b_1} \\ \dots \\ \frac{\partial S}{\partial b_n} \end{pmatrix}; \frac{\partial S}{\partial \mathbf{b}} = \begin{pmatrix} -2 k_0 + 2b_0 d_0 \\ -2 k_1 + 2b_1 d_1 \\ \dots \\ -2 k_n + 2b_n d_n \end{pmatrix} = \mathbf{0}$

- $\frac{\partial S}{\partial \mathbf{b}} = 0 - 2X^T \mathbf{y} + 2 X^T X \mathbf{b} = 0 \Rightarrow X^T X \mathbf{b} - X^T \mathbf{y} = 0$

- $X^T X \mathbf{b} = X^T \mathbf{y}$  represent normal e

- $\mathbf{b} = (X^T X)^{-1} X^T \mathbf{y}$

$$\bullet X = \begin{bmatrix} 1 & x_{11} & x_{21} & \dots & x_{n1} \\ 1 & x_{12} & x_{22} & \dots & x_{n2} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{1m} & x_{2m} & \dots & x_{nm} \end{bmatrix}; X^T = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_{11} & x_{12} & x_{13} & \dots & x_{1m} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_{11} & x_{12} & x_{13} & \dots & x_{1m} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & x_{n3} & \dots & x_{nm} \end{bmatrix} \begin{bmatrix} 1 & x_{11} & x_{21} & \dots & x_{n1} \\ 1 & x_{12} & x_{22} & \dots & x_{n2} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{1m} & x_{2m} & \dots & x_{nm} \end{bmatrix}$$

# The Matrix form

$y_1$	1	$x_{1,1}$	$x_{2,1}$	$\dots$	$x_{k,1}$
$y_2$	1	$x_{1,2}$	$x_{2,2}$	$\dots$	$x_{k,2}$
$\dots$	1	$\dots$	$\dots$	$\dots$	$\dots$
$y_n$	1	$x_{1,n}$	$x_{2,n}$	$\dots$	$x_{k,n}$

$$\beta = (X'X)^{-1} X'y$$



# Comparison between analytical and numerical method

## Analytical

Normal Equations

Direct formula/Matrix Formula

No requirement of choosing learning rate

Does not require feature scaling

Becomes slow , when  $n$  is large ; involves calculation of inverse of  $(X^T X)^{-1}$ ;  $n \times n$  matrix,  $O(n^3)$ , upto  $n=10,000$  fine,  $n \geq 10000$ , gradient descent advisable

## Numerical

Iterative: many iterations

Gradient Descent Method

Need to choose learning rate appropriately

Requires feature scaling

Works well when  $n$  is large  $O(n^2)$