

11/06/2020

## Simple Regression with mathematics [University of Washington]

- \* dataset which is going to use in this is on "Predicting house prices"
  - so starting with the Question:
    - How much is my house worth?
    - look at recent sales in my neighborhood. How much did they sell for?
- Regression fundamentals
  - 1) Data
  - 2) Model
  - 3) Task

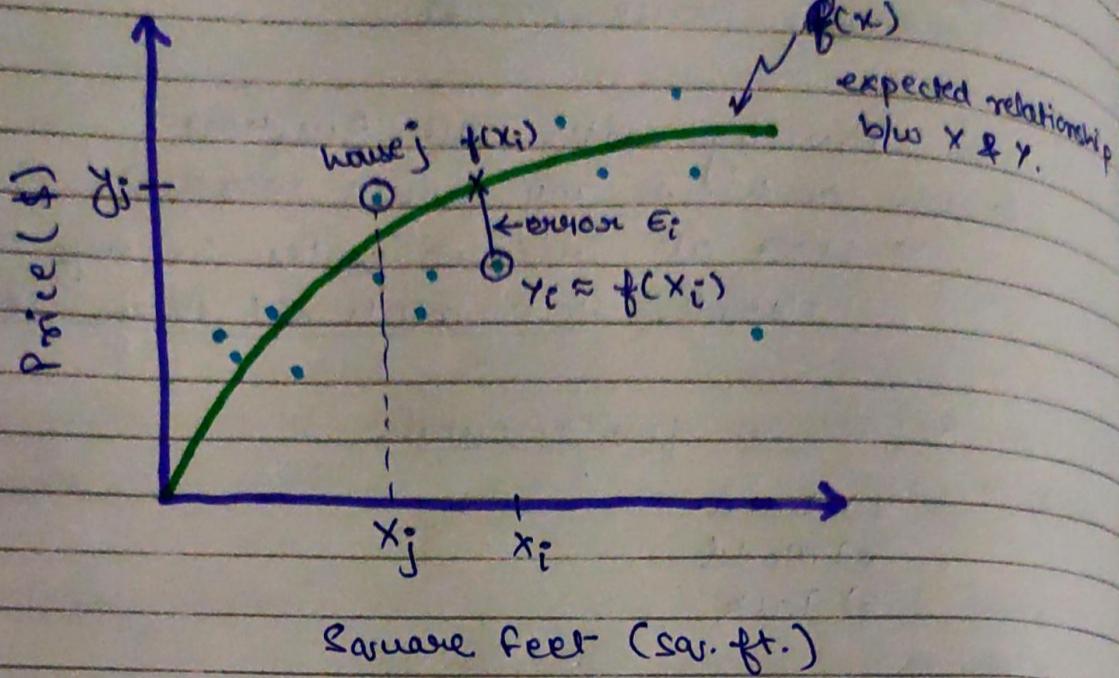
1) Data	Input	Output
1	( $x_1 = \text{sq.ft.}$ , $y_1 = \$$ )	
2	( $x_2 = \text{sq.ft.}$ , $y_2 = \$$ )	
3	( $x_3 = \text{sq.ft.}$ , $y_3 = \$$ )	
4	( $x_4 = \text{sq.ft.}$ , $y_4 = \$$ )	
⋮	⋮	

Input vs Output:

- $y$  is the quantity of interest
- assume  $y$  can be predicted from  $x$ .

## 2) Model :

How we assume the model works?



Regression Model :

$$Y_i = f(x_i) + \epsilon_i$$

$E[\epsilon_i] = 0$  ← equally likely that  
error is +ve or -ve.

↑  
expected value

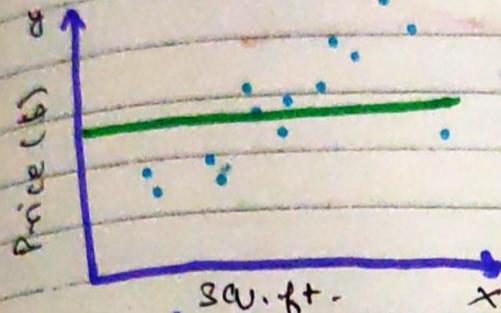
↓  
 $Y_i$  is equally likely  
to be above or  
below  $f(x_i)$

• "Essentially, all models are wrong,  
but some are useful."

- George Box, 1987.

### TASK 1:

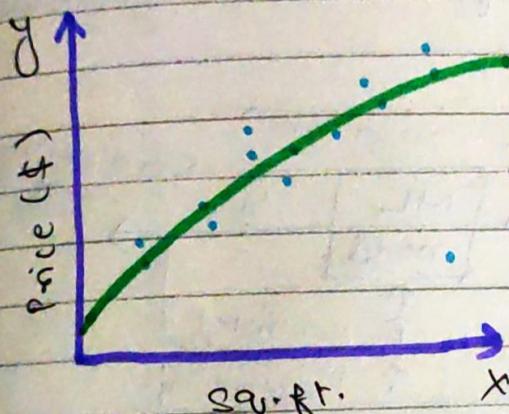
We have to choose which model  $f(x)$ ?  
suppose here are 4 examples:



(a)

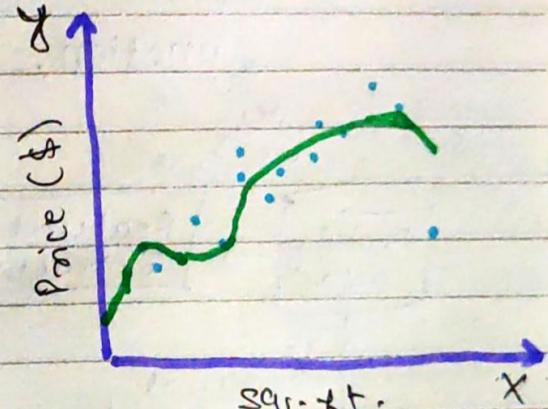


(b)



✓

(c)

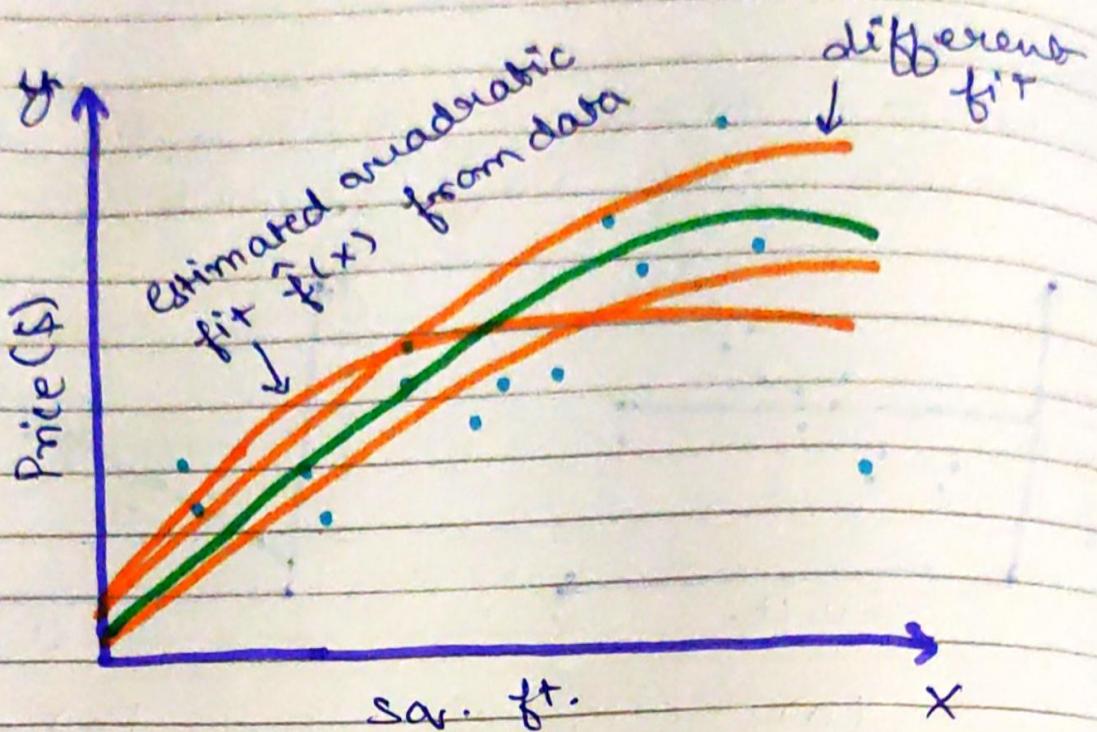


(d)

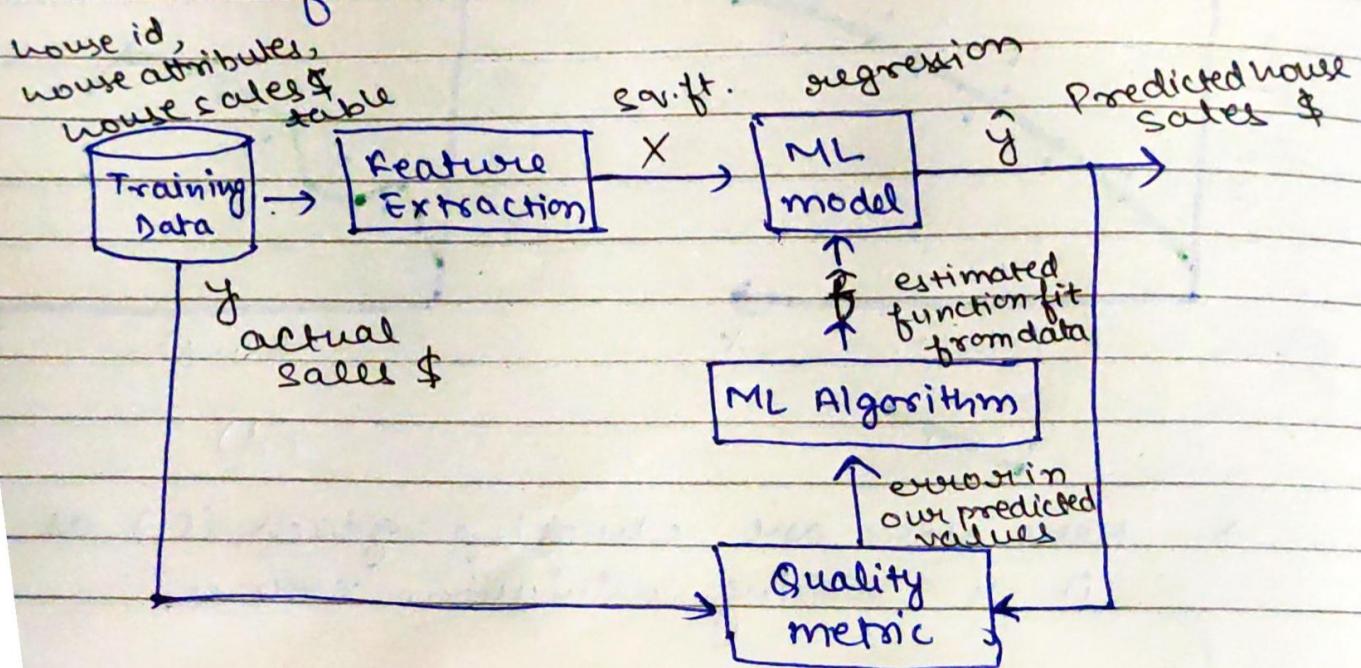
→ Here we are choosing graph (c) as  
it is having minimum error.

## Task 2 :

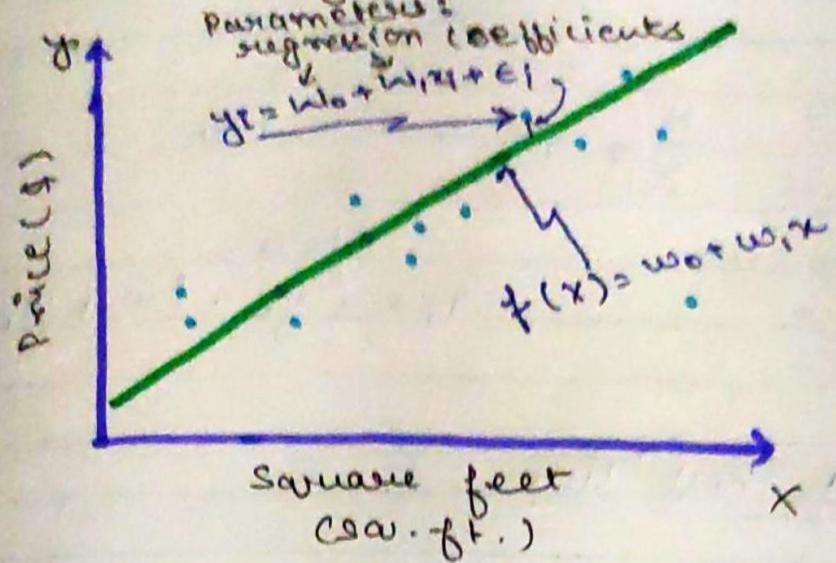
for a given model  $f(x)$ ,  
estimate function  $\hat{f}(x)$  from data



→ Assume mode  $f(x)$  is a quadratic function.



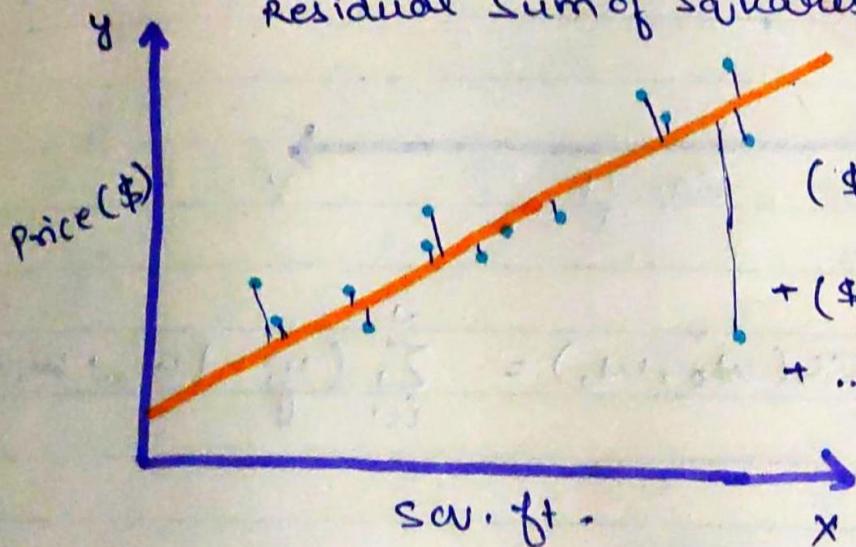
## Simple Linear Regression Model



⇒ Fitting a line to data

"cost" of using a given line

residual sum of squares (RSS)



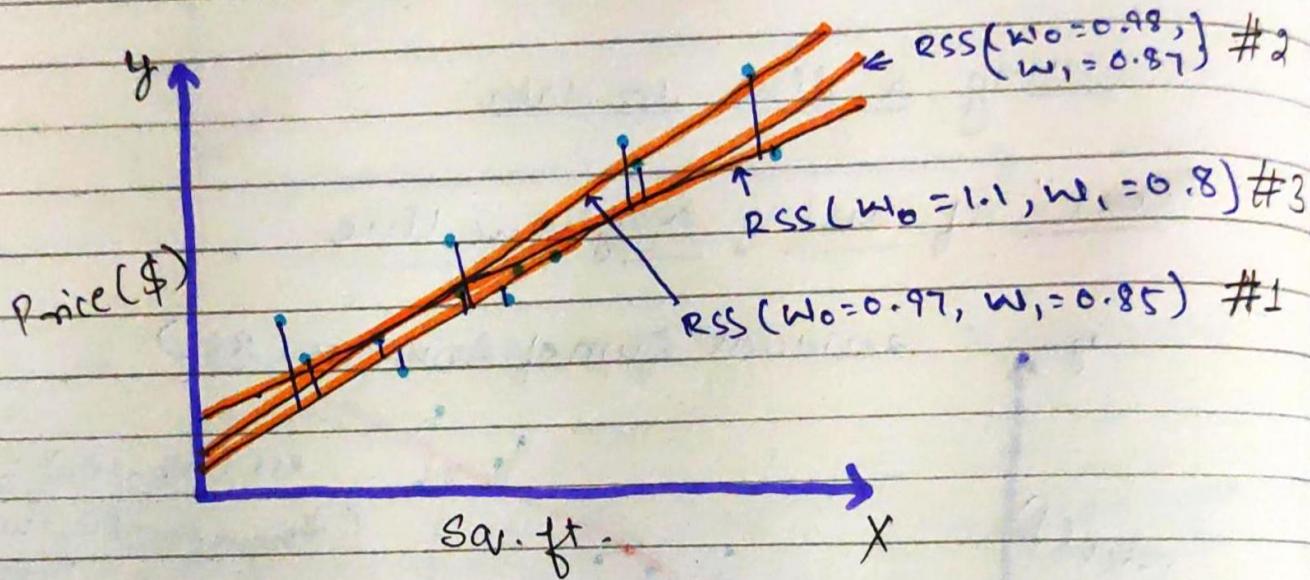
$$\begin{aligned}
 \text{RSS } (w_0, w_1) = & (\$_{\text{house}_1} - [w_0 + w_1 \text{sq. ft.}_{\text{house}_1}])^2 \\
 & + (\$_{\text{house}_2} - [w_0 + w_1 \text{sq. ft.}_{\text{house}_2}])^2 \\
 & + \dots \text{ [including all training houses]}
 \end{aligned}$$

$$\rightarrow \text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

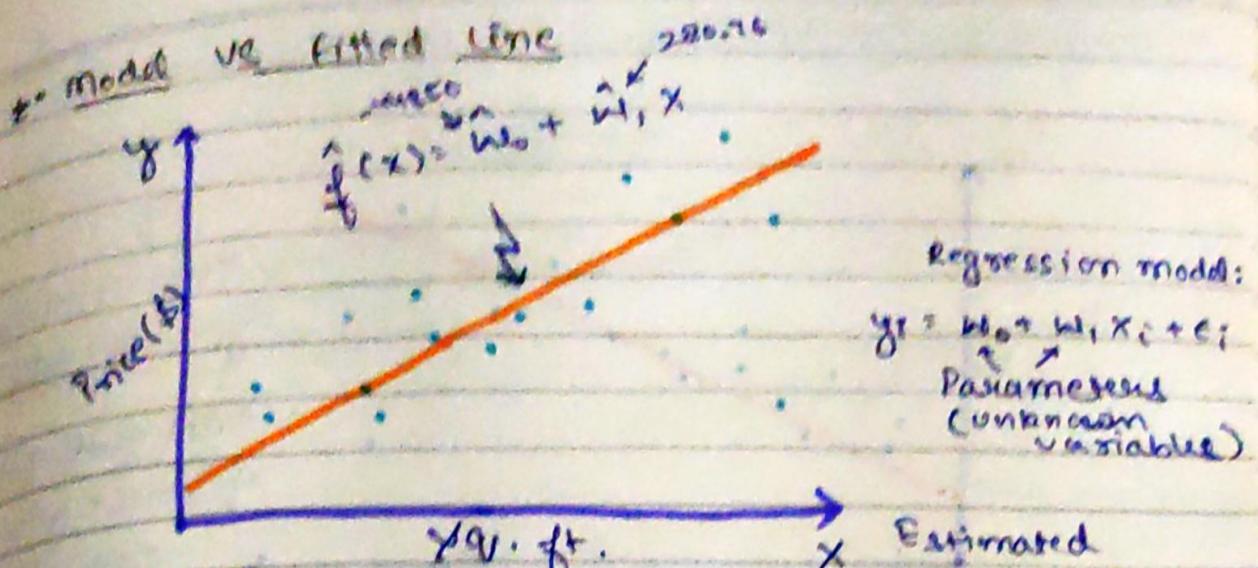
$\sum_{i=1}^N a_i = a_1 + a_2 + \dots + a_N$

Here,  
 $a_i = (y_i - [w_0 + w_1 x_i])^2$

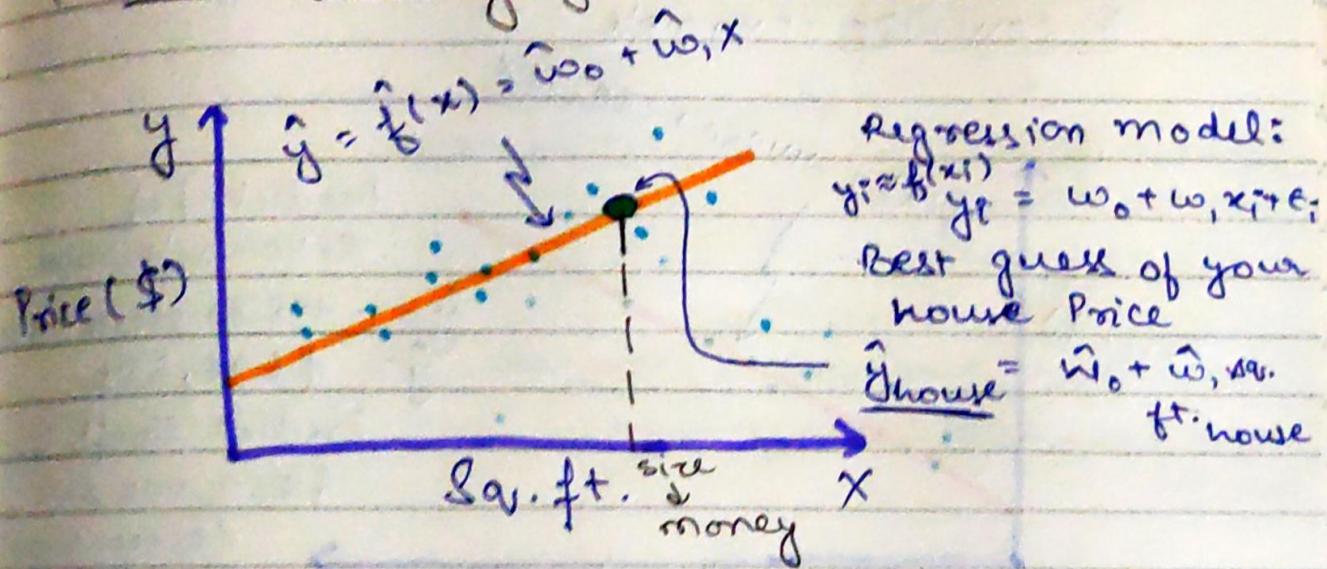
Find "Best" line



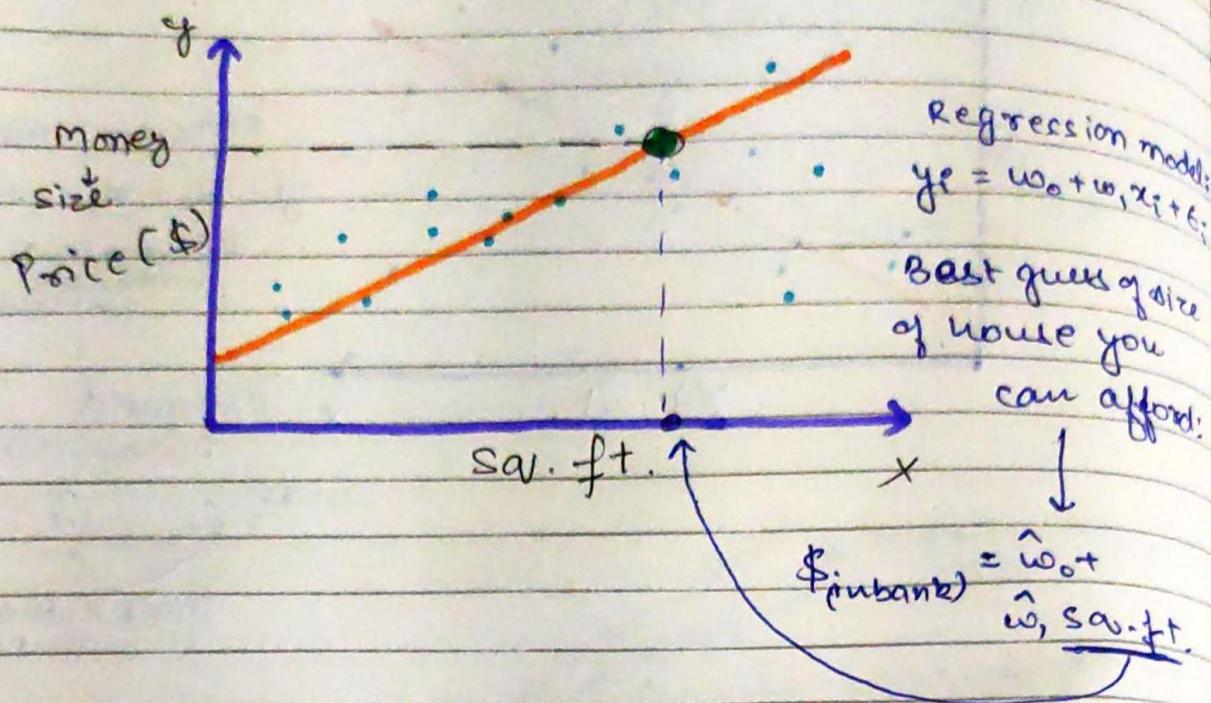
$$\text{RSS}(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$



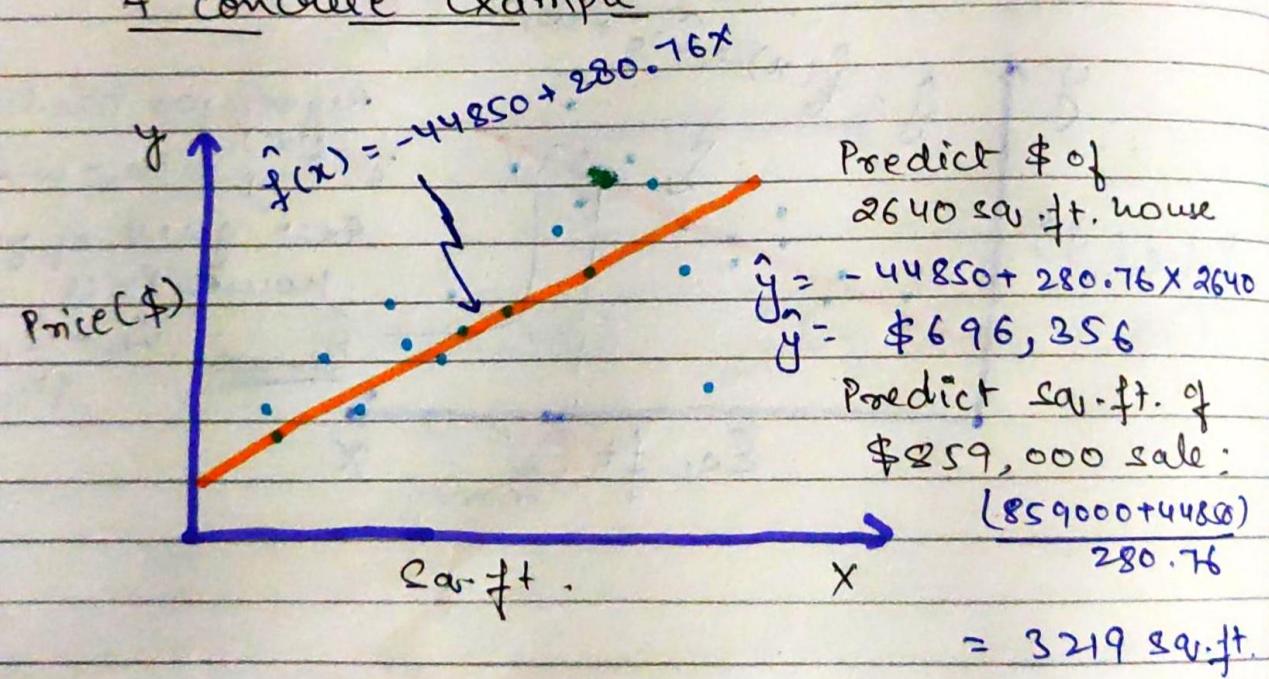
+ Seller: Predicting your house Price



① Buyer: Predicting size of house



4 concrete example



## Interpreting the coefficients

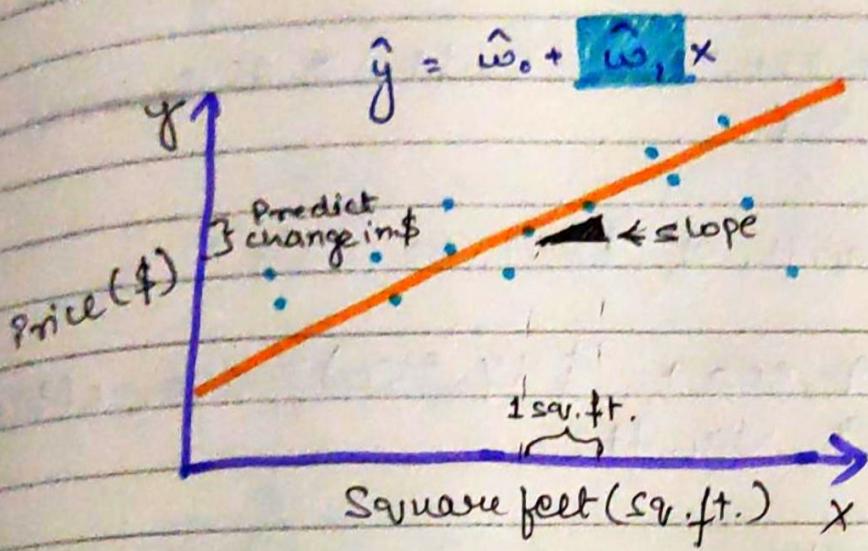
$$\hat{y} = \hat{w}_0 + \hat{w}_1 x$$



$$\hat{y} = \hat{w}_0$$

when  $x=0$

↓  
This is not very meaningful.

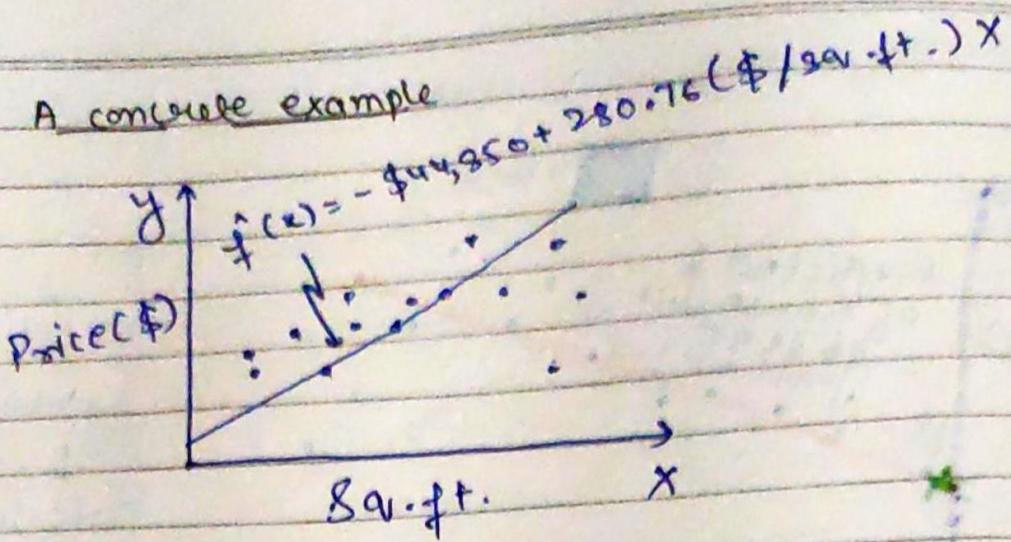


$$\begin{aligned} \$1001 \text{ sq.ft.} - \$1000 \text{ sq.ft.} \\ &= \hat{w}_0 + \hat{w}_1 \cdot 1001 \text{ sq.ft.} \\ &\quad - (\hat{w}_0 + \hat{w}_1 \cdot 1000 \text{ sq.ft.}) \\ &= \hat{w}_1 \end{aligned}$$

↓  
Predict change in the O/P per unit change in I/P.

warning: Magnitude depends on unit of both feature & observations

A concrete example



Predict \$ of 2,640 sq.ft. house:

$$\begin{aligned} & -\$44,850 + 280.76 (\$/\text{sq-ft.}) * 2640 \text{ sq-ft.} \\ &= \$696,356 \end{aligned}$$

Predict sq.ft. of \$859,000 sale:

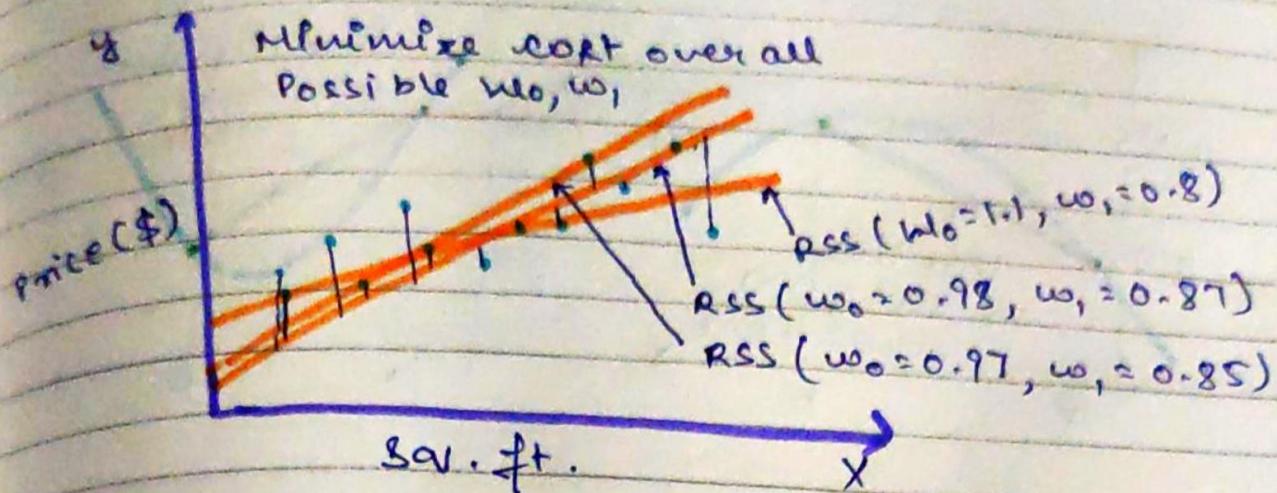
$$\begin{aligned} &= (\$859,000 + \$44,850) / 280.76 (\$/\text{sq-ft.}) \\ &= 3,219 \text{ sq-ft.} \end{aligned}$$

• But what if:

- House was measured in sq. metres?
- Price was measured in RMB?

• Algorithm for fitting the model

• find the "best" line



$$RSS(w_0, w_1) = \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

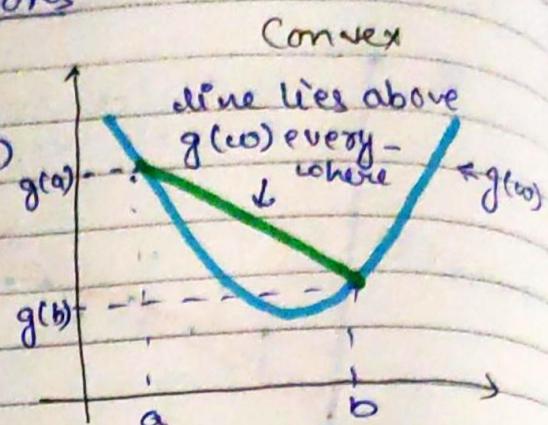
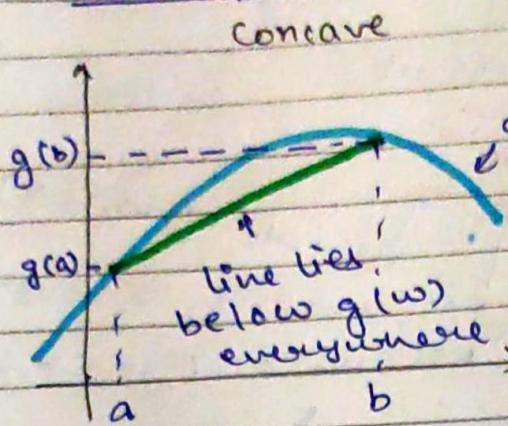
→ minimize function overall possible  $w_0, w_1$

$$\min_{w_0, w_1} \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

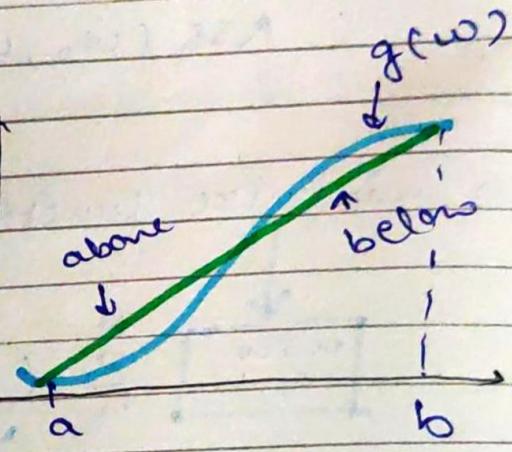
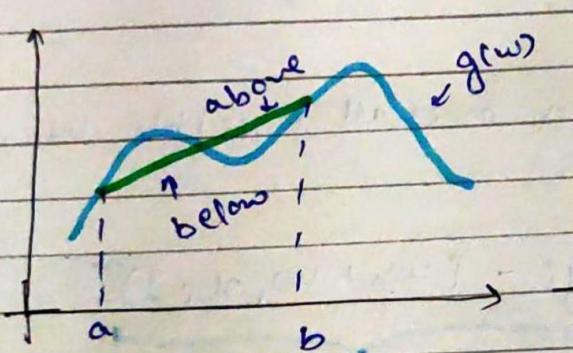
RSS ( $w_0, w_1$ ) is a function of 2 variables =  $g(w_0, w_1)$

→ An aside on optimization :

### Convex / Concave functions



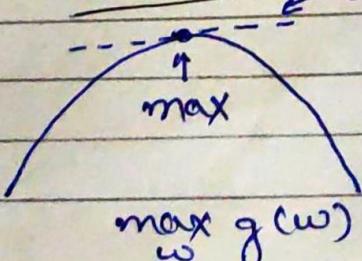
Neither



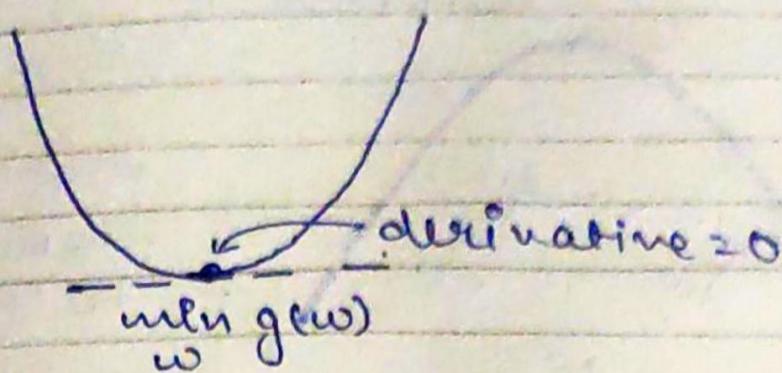
Finding the max or min analytically

concave

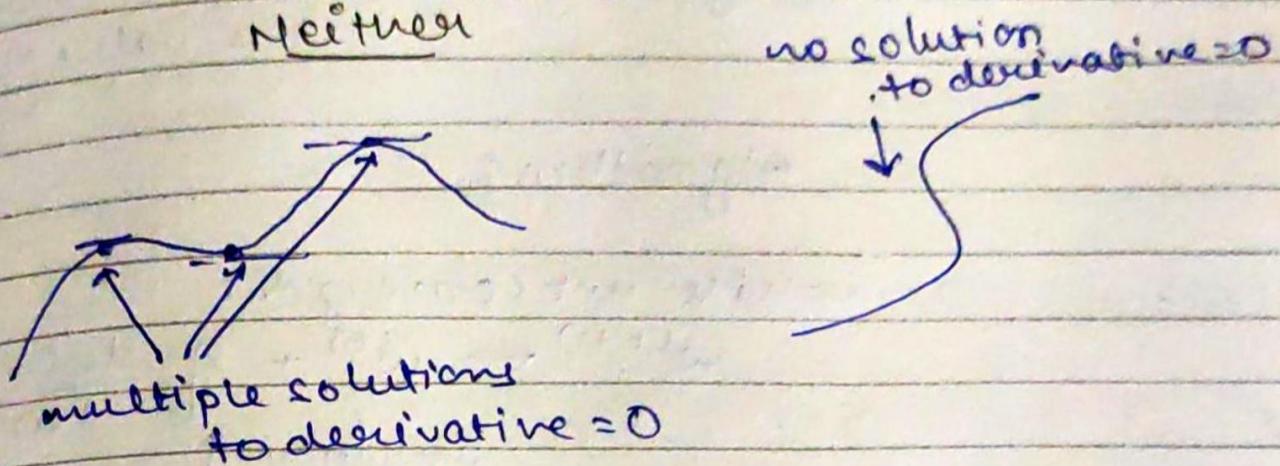
← derivative  $\geq 0$



convex



Neither



Example

$$g(w) = 5 - (w-10)^2$$

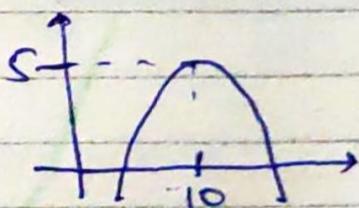
$$\frac{d(g(w))}{dw} = 0 - 2 \times (w-10)' \times 1$$

$$= -2w + 20$$

Set derivative = 0

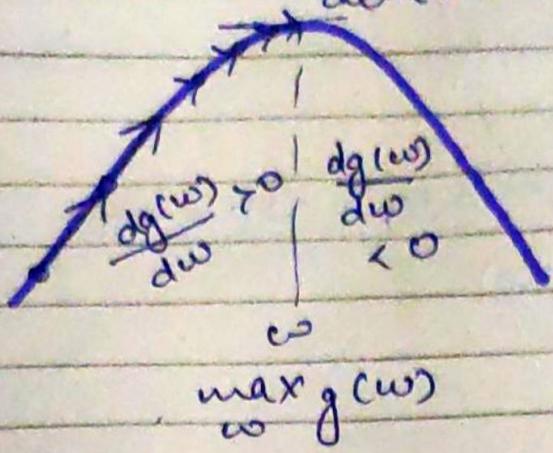
$$-2w + 20 = 0$$

$$\boxed{w = 10}$$



Finding the max via hill climbing

derivative = 0



How does know  
whether to move  
w to right or left?  
↓

Increase or decrease  
the value of w?

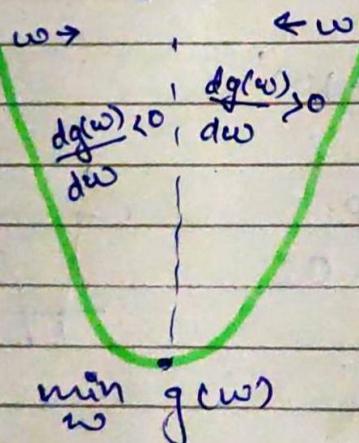
Algorithm:

while not converged :

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \frac{dg(w)}{dw}$$

iteration ↑  
step size ↑

Finding the min via hill descent



when derivative is positive, we want to decrease w.  
and when derivative is negative we want to increase w.

Algorithm:

while not converged

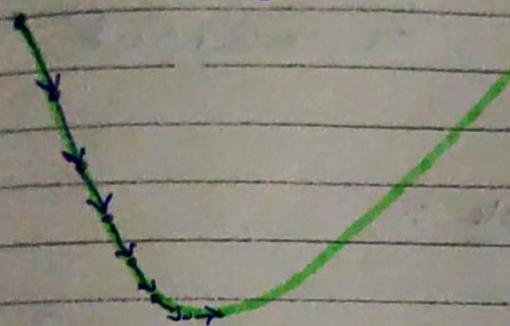
$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{dg}{dw} |_{w^{(t)}}$$

choosing the stepsize -

1) Fixed stepsize



2) Decreasing stepsize or stepsize schedule



common choices:

$$\eta_t = \frac{\alpha}{t}$$

$$\eta_t = \frac{\alpha}{\sqrt{t}}$$

$\downarrow \eta_t$   
iterations

convergence criteria:

- for convex functions, optimum occurs when  $\frac{dg(\omega)}{d\omega} = 0$

- In practice stop when

$$\left| \frac{dg(\omega)}{d\omega} \right| < \epsilon$$

↑  
threshold  
to be set

Algorithm:  
while not converged  
 $\omega^{(t+1)} \leftarrow \omega^{(t)} - \eta \frac{dg}{d\omega} \Big|_{\omega^{(t)}}$

2. Moving to multiple dimensions: Gradient

gradient  $[w_0, w_1, \dots, w_p]$

$$\nabla g(w) = \left[ \begin{array}{c} \frac{\partial g}{\partial w_0} \\ \frac{\partial g}{\partial w_1} \\ \vdots \\ \frac{\partial g}{\partial w_p} \end{array} \right] \quad \text{ } \quad (p+1) - \text{dimensional vector}$$

partial derivative is like  
a derivative w.r.t.  $w_i$ ,  
treating all other variables  
as constant.

Gradient example :

$$g(w) = 5w_0 + 10w_0w_1 + 2w_1^2$$

$$\frac{\partial g}{\partial w_0} = 5 + 10w_1$$

$$\frac{\partial g}{\partial w_1} = 10w_0 + 4w_1$$

$$\nabla g(w) = \left[ \begin{array}{c} 5 + 10w_1 \\ 10w_0 + 4w_1 \end{array} \right]$$

NOTE: Contour Plots is a bird's eye view of 3D plot of RSS with tangent plane at minimum.

### - Gradient Descent Algorithm:

while not converged

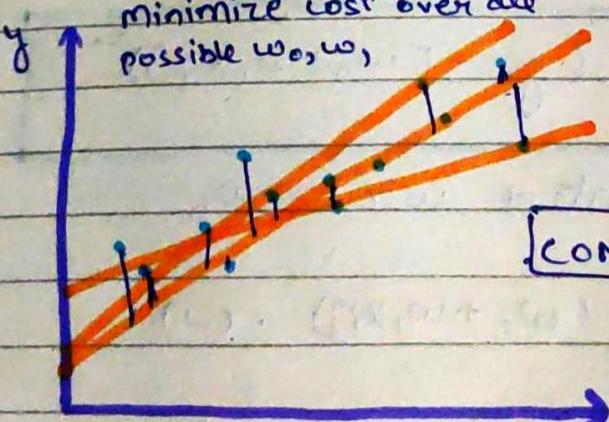
$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla g(\mathbf{w}^{(t)})$$
$$[\vdots] \leftarrow [\vdots] - \eta [\vdots]$$

convergence:  
 $\|\nabla g(\mathbf{w})\| < \epsilon$

### • finding the least square lines:

#### Find "best" line

minimize cost over all possible  $w_0, w_1$



$$\min_{w_0, w_1} \sum_{i=1}^N (y_i - [w_0 + w_1 x_i])^2$$

**CONVEX**

⇒ solution is unique  
+ gradient descent algorithm will converge to minimum

compute the gradient

$$\rightarrow \text{RSS}(\omega_0, \omega_1) = \sum_{i=1}^N (y_i - [\omega_0 + \omega_1 x_i])^2$$

$$\begin{aligned}\rightarrow \frac{d}{d\omega} \cdot \sum_{i=1}^N g_i(\omega) &= \frac{d}{d\omega} (g_1(\omega) + g_2(\omega) + \dots + g_N(\omega)) \\ &= \frac{d}{d\omega} g_1(\omega) + \frac{d}{d\omega} g_2(\omega) + \dots + \frac{d}{d\omega} g_N(\omega) \\ &= \sum_{i=1}^N \frac{d}{d\omega} g_i(\omega)\end{aligned}$$

$\rightarrow$  In our case

$$g_i(\omega) = (y_i - [\omega_0 + \omega_1 x_i])^2$$

$$\frac{\partial \text{RSS}(\omega)}{\partial \omega_0} = \sum_{i=1}^N \frac{\partial}{\partial \omega_0} (y_i - [\omega_0 + \omega_1 x_i])^2$$

same for  $\omega_1$

$$\rightarrow \boxed{\text{RSS}(\omega_0, \omega_1) = \sum_{i=1}^N (y_i - [\omega_0 + \omega_1 x_i])^2}$$

$\Rightarrow$  Taking the derivative w.r.t.  $\omega_0$

$$= \sum_{i=1}^N 2(y_i - [\omega_0 + \omega_1 x_i])' \cdot (-1)$$

$$= -2 \sum_{i=1}^N (y_i - [\omega_0 + \omega_1 x_i])$$

⇒ Taking the derivative w.r.t.  $w_0$ ,

$$= \sum_{i=1}^N 2(y_i - [w_0 + w_1 x_i]) \cdot (-x_i)$$

$$= -2 \sum_{i=1}^N (y_i - [w_0 + w_1 x_i]) x_i$$

⇒ Putting it together:

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$

Approach 1: Set gradient = 0

$$\nabla \text{RSS}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (w_0 + w_1 x_i)] x_i \end{bmatrix}$$

Top term:  $N \leftarrow$  average house sales price estimate of the slope

$$\hat{w}_0 = \frac{\sum_{i=1}^N y_i}{N} - \hat{w}_1 \frac{\sum_{i=1}^N x_i}{N} \leftarrow$$
 average sq.ft.

Bottom term:

$$\sum y_i x_i - \hat{w}_0 \sum x_i - \hat{w}_1 \sum x_i^2 = 0$$

$$\hat{w}_1 = \frac{\sum y_i x_i - \frac{\sum y_i x_i}{N}}{\sum x_i^2 - \frac{\sum x_i^2}{N}}$$

NOTE:

You have to find these 4 terms!

$$1) \sum_{i=1}^N y_i \quad 2) \sum_{i=1}^N x_i$$

$$3) \sum_{i=1}^N y_i x_i \quad 4) \sum_{i=1}^N x_i^2$$

## Approach 2: Gradient Descent

Interpreting the gradient: actual value taken observation  
predicted value  $\hat{y}_i(w_0, w_1)$

$$\nabla_{\text{RSS}}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - (\omega_0 + \omega_1 x_i)] \\ -2 \sum_{i=1}^N [y_i - (\omega_0 + \omega_1 x_i)] x_i \end{bmatrix}$$

$$\nabla_{\text{RSS}}(w_0, w_1) = \begin{bmatrix} -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] \\ -2 \sum_{i=1}^N [y_i - \hat{y}_i(w_0, w_1)] x_i \end{bmatrix}$$

while not converged

$$\begin{bmatrix} w_0^{(t+1)} \\ w_1^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} w_0^{(t)} \\ w_1^{(t)} \end{bmatrix} + \alpha \eta \begin{bmatrix} \sum_{i=1}^N [y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})] \\ \sum_{i=1}^N [y_i - \hat{y}_i(w_0^{(t)}, w_1^{(t)})] x_i \end{bmatrix}$$

- If overall, underpredicting  $\hat{y}_i$ , then  $\sum_i [y_i - \hat{y}_i]$  is positive
  - $w_0$  is going to increase
  - similar intuition for  $w_1$ , but multiply by  $x_i$

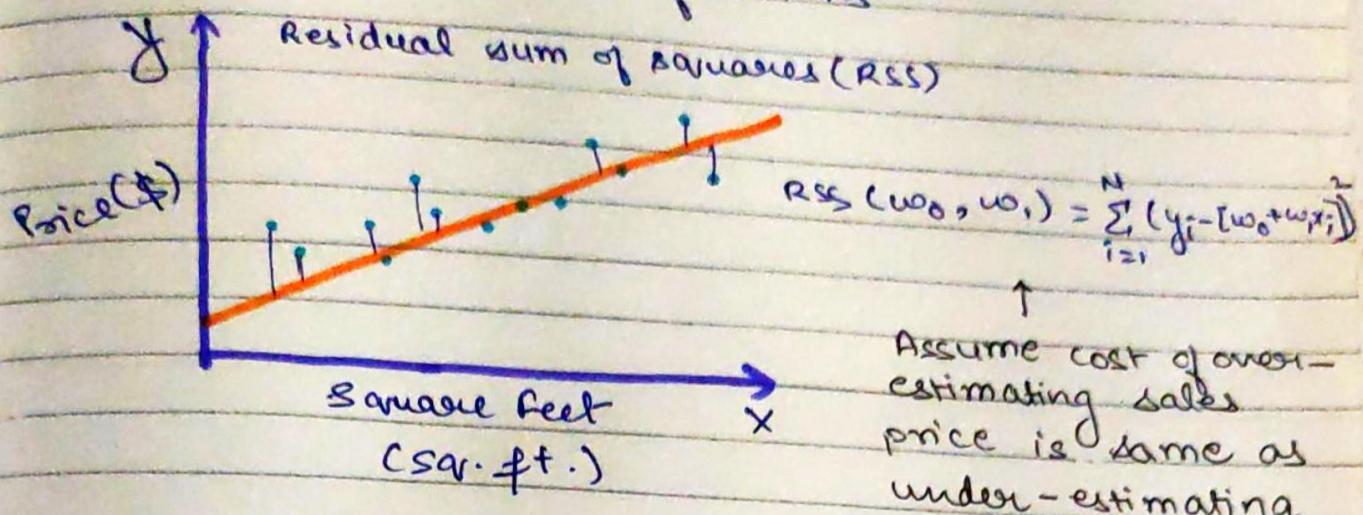
## Comparing the approaches

- For most ML problems, cannot solve  $\text{gradient} = 0$
- Even if solving  $\text{gradient} = 0$  is feasible, gradient descent can be more efficient.
- Gradient Descent relies on choosing step size & convergence criteria.

→ Influence of high leverage points

Asymmetric errors

→ Symmetric cost functions



→ Asymmetric cost functions

