

$$X = [x_1, x_2, x_3 \dots x_n]$$

Learning from Data

$$\begin{cases} \theta_0 = \bar{y} - \theta_1 \bar{x} \\ \theta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \end{cases}$$

Linear Regression

- Closed form solution

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

→ vector notation.

[Learning from data]

- Independent variable, Input, features

- dependent variable, Output, target (label)

Examples: stock price prediction,

- Quantitative variables, Qualitative variables

# Define the above terms with examples

→ Categorical.

Caesarian Delivery Prediction  
Input: Number of delivery, 1, 2, 3  
Pregnant Vitals = 1

proceed with Caesarian delivery  
Output: True, False  
= 1 = 0 (-1)  
→ {cat / Dog.}  
Qualitative Variable.

⊕ ⊖ ⊗ ⊘ ⊙

Dummy variable (one-hot encoding)

$x = \{cat, dog, zebra\}$

one hot encoding  $[x = 0, 1, 2]$

$x = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$  cat  
dog  
zebra.

Convert to some meaningful numbers.

Computer Algorithm.  $\left\{ \text{look for data } D = \{(x_1, y_1), \dots, (x_n, y_n)\} \right.$   
 $y = \theta_0 + \theta_1 x = \text{parameters}$   $y = mx + c$

# Linear Regression: Worked out Example

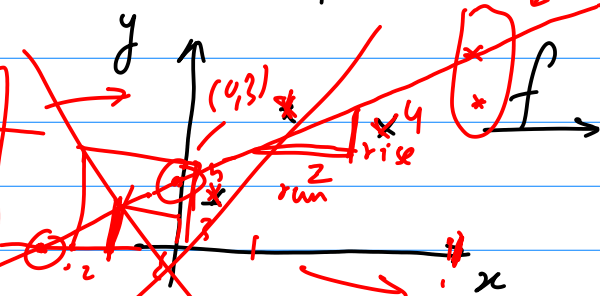
single column

Univariate

colinear

$x_n$	$y_n$
1	4.8
3	11.3
5	17.2
6	?

24, 23:?



We want to fit a line that closely passes through the given point.

$y = 2x + 3$

$y = 2x + c$

The model has the form  $y = f(x) = \theta_0 + \theta_1 x$

This is an example of Univariate linear regression.

$y = mx - ?$

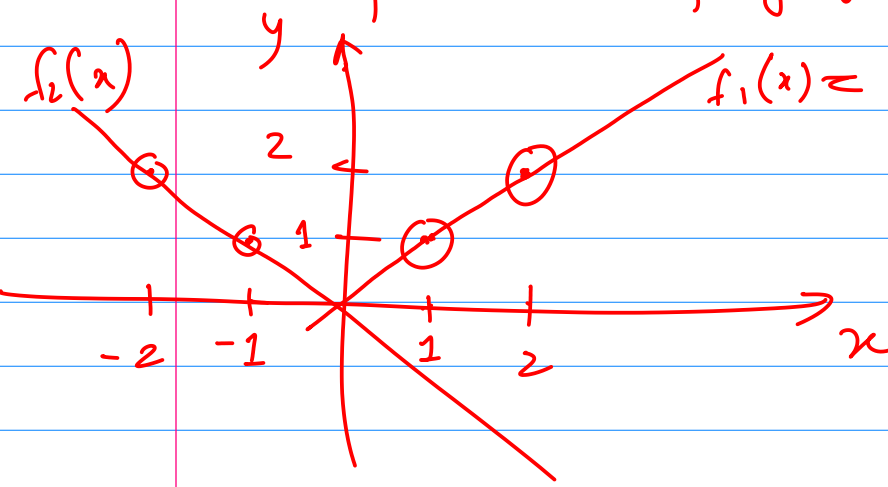
$y =$

Discuss possible solution for  $y = f(x) = \theta_0 + \theta_1 x$  [ $y = mx + c$ ]

$y = mx + c$

$f_1(x) = 1x + 0$

$f_2(x) = -1x + 0$   
 $m$   $c$



$f_1(x) \quad m = \frac{\text{rise}}{\text{run}} = \frac{1}{1} = 1$

$y = mx + c$   
 $\downarrow$   
 $y = \theta_1 x + \theta_0$

$f_2(x) = m$   $x_1 \quad x_2 \quad x_3 \quad \dots \quad x_n \rightarrow y_n$   
 $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$   
multivariate.

Closed form solution for  $y = f(x) = \theta_0 + \theta_1 x$  [Univariate Regression]  
 Under certain criteria, the best possible value of parameters linear model.

Analysis  
 solution  
 (parameters):  $\theta \Rightarrow ?$

$$\left[ \begin{array}{l} \theta_1 = \frac{\sum x_i y_i - \bar{x} \cdot \bar{y}}{\sum x_i^2 - (\bar{x})^2} \quad , \quad \theta_0 = \bar{y} - \theta_1 \cdot \bar{x} \end{array} \right] \text{--- (1)}$$

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

$x_n$	$y_n$	$x_n \cdot y_n$	$x_n^2$
1	4.8	4.8	1
3	11.3	33.9	9
5	17.2	86	25
6	?		
$\bar{x} = 3$	$\bar{y} = 11.1$	$\bar{xy} = 41.57$	$\bar{x^2} = 11.67$

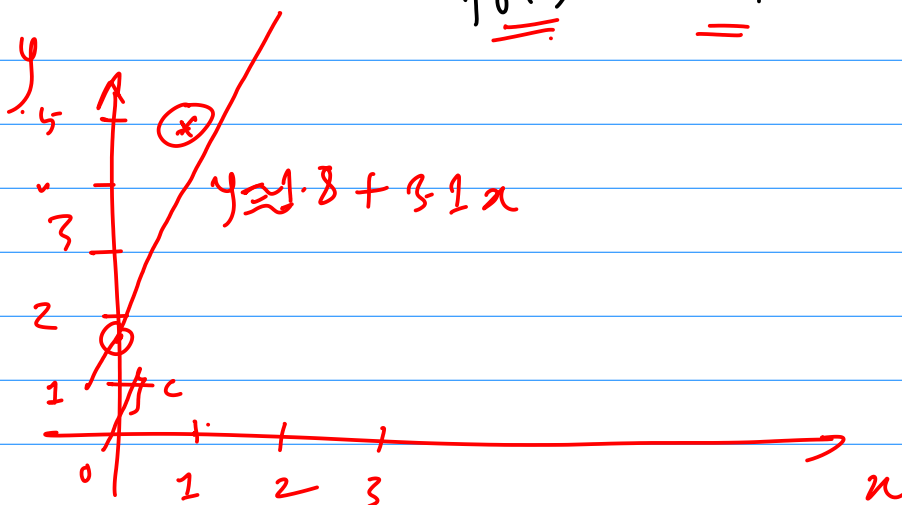
$$\theta_0 = f(x)_{x=0} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$y = mx + c = \theta_0 + \theta_1 x \rightarrow \theta_0$$

Putting values into eqn. (1),  $\theta_1 = \frac{41.57 - 3 \times 11.1}{11.67 - 3 \times 3} = 3.1$ ,  $\theta_0 = 11.1 - 3.1 \times 3 = 1.8$

Hence,  $y = f(x; \theta) = f_\theta(x) = 1.8 + 3.1x$   
 $f_\theta(1) = 20.4$

(intercept)  
 bias.  
 DC value



$$y = f_{\theta}(x) = \theta_0 + \theta_1 x \quad \left[ \begin{aligned} \theta_1 &= \frac{\overline{xy} - \bar{x} \cdot \bar{y}}{\overline{x^2} - (\bar{x})^2} \\ \theta_0 &= \bar{y} - \theta_1 \bar{x} \end{aligned} \right]$$

# Where did this magic formula come from!

Given a dataset  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , find a functional relationship using a linear model.

$$f(x; \theta) = f_{\theta}(x) = \theta_0 + \theta_1 x \quad | \quad (\theta_0, \theta_1) \rightarrow \text{find.}$$

$(-1.7, 2) \leftrightarrow (2.1, 3)$

→ How do we quantify when the solution is good enough?

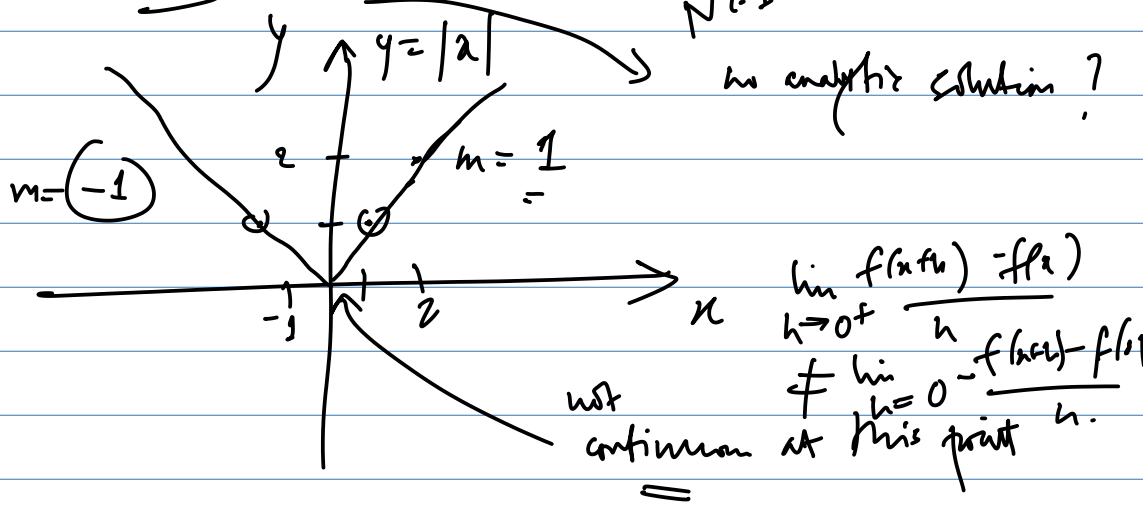
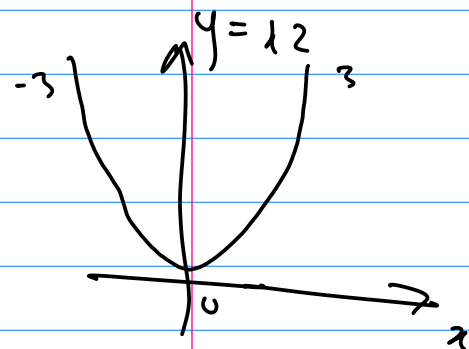
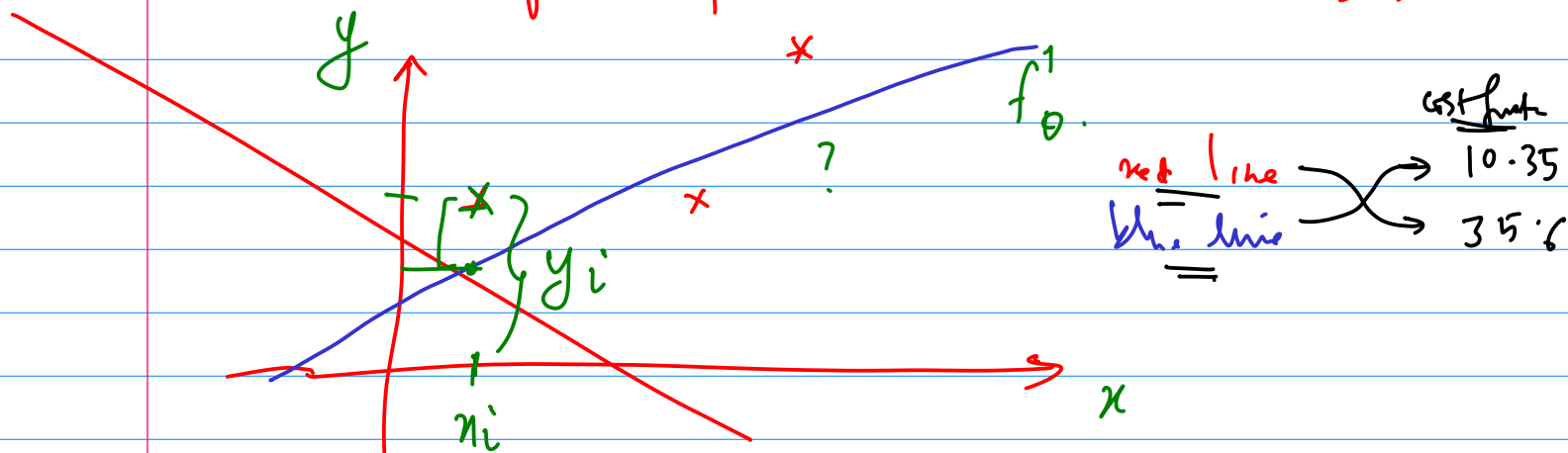
[Cost function] → 0 means this fits the data very well

(loss function) ② Mean Squared Error (MSE) or L2 norm predicted

$$= L(f_{\theta}(x), D) = \sum_{i=1}^N (f_{\theta}(x_i) - y_i)^2 \rightarrow \underline{1}$$

[Smooth, convex, differentiable]

Evaluate this cost function for some  $\theta$  and some dataset  $D$ .



cost function (MSE)  $\frac{dL}{c}$   $\frac{dy}{dx}$   $\frac{dL}{dslope}$   $\frac{dL}{dinh}$   $\frac{\partial L}{\partial \theta_0}$   $\frac{\partial L}{\partial \theta_1}$

# Use the idea that  $\frac{\partial L}{\partial \theta}$  should be 0 if  $\theta$  is the best solution.

$y = x^2$   
 $y = ax^2 + bx + c$   
 $m = 1$   
 $\frac{\partial L}{\partial \theta} \rightarrow$  quadratic eqn.

convex  
vs  
non convex.



$$L = \frac{1}{N} \sum_{i=1}^N L_n(f_{\theta}(x_i), y_i)$$

$$= \frac{1}{N} \sum_{i=1}^N (y_i - (\theta_0 + \theta_1 x_i))^2$$

MSE.

$$= \frac{1}{N} \sum_{i=1}^N (y_i^2 - 2(\theta_0 + \theta_1 x_i)y_i + (\theta_0 + \theta_1 x_i)^2)$$

$$L = \frac{1}{N} \sum_{i=1}^N (y_i^2 - 2\theta_0 y_i - 2\theta_1 x_i y_i + \theta_0^2 + 2\theta_0 \theta_1 x_i + \theta_1^2 x_i^2)$$

Differentiating the cost function,

$$\frac{dL}{d\theta_0}, \frac{dL}{d\theta_1}, \frac{dL}{d\theta_0} = \frac{1}{N} \sum_{i=1}^N (-2y_i + 2\theta_0 + 2\theta_1 x_i) = 0$$

$$y = 3x + 2$$

$$\frac{dy}{dx} = 3$$

$$\frac{dL}{d\theta_1} = \frac{1}{N} \sum_{i=1}^N (-2x_i y_i + 2\theta_0 x_i + 2\theta_1 x_i^2) = 0$$

$$\text{Eqn (i)} \quad \frac{1}{N} \sum_{i=1}^N (-2y_i + 2\theta_0 + 2\theta_1 x_i) = 0$$

$$\text{Eqn (ii)}, \quad -\overline{x_i y_i} + \theta_0 \overline{x_i} + \theta_1 \overline{x_i^2} = 0$$

$$\text{ii}, \quad \frac{1}{N} \sum_{i=1}^N (-y_i + \theta_0 + \theta_1 x_i) = 0$$

$$\text{ii}, \quad \theta_0 \overline{x_i} = \overline{x_i y_i} - \theta_1 \overline{x_i^2}$$

$$\text{ii}, \quad -\overline{y_i} + \theta_0 + \theta_1 \overline{x_i} = 0$$

$$\text{ii}, \quad \theta_0 = \overline{y_i} - \theta_1 \overline{x_i} \quad \text{--- (iv)}$$

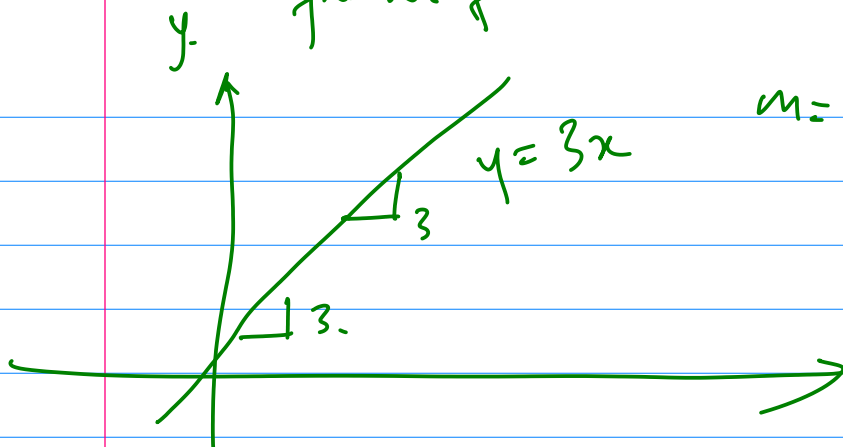
$$\text{ii}, \quad \theta_0 = \frac{\overline{x_i y_i} - \theta_1 \overline{x_i^2}}{\overline{x_i}}$$

$\theta_0, \theta_1$   
 2 unknown,  
 2 eqns.

$$\frac{1}{N} \sum_{i=1}^N \theta_1 x_i = \theta_1 \frac{x_1 + x_2 + \dots + x_N}{N} = \theta_1 \overline{x_i}$$

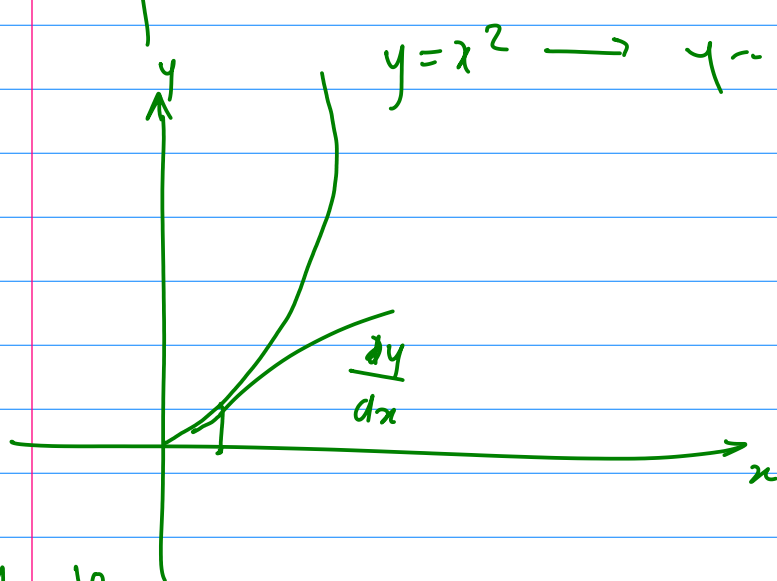
$$\frac{1}{N} \sum_{i=1}^N \theta_0 = \theta_0 \frac{1 + 1 + \dots + 1}{N} = \theta_0$$

# Refresher for Calculus 101: derivative. (1<sup>st</sup>).



$$m = \frac{dy}{dx} = 3 \text{ (slope)} \theta_1$$

general form of slope  
 $\frac{dy}{dx} = f'(x)$



$$y = x^2 \rightarrow y = \text{func.} \\ - y = ?$$

$$\frac{d \sin x}{dx} = \cos x$$

$$\frac{dx^2}{dx} = 2x$$

$$\frac{de^x}{dx} = e^x$$

$$\frac{d \log_{10} x}{dx} = ? \frac{1}{x}$$

$$\frac{d \ln x}{dx} = ? \frac{1}{x}$$

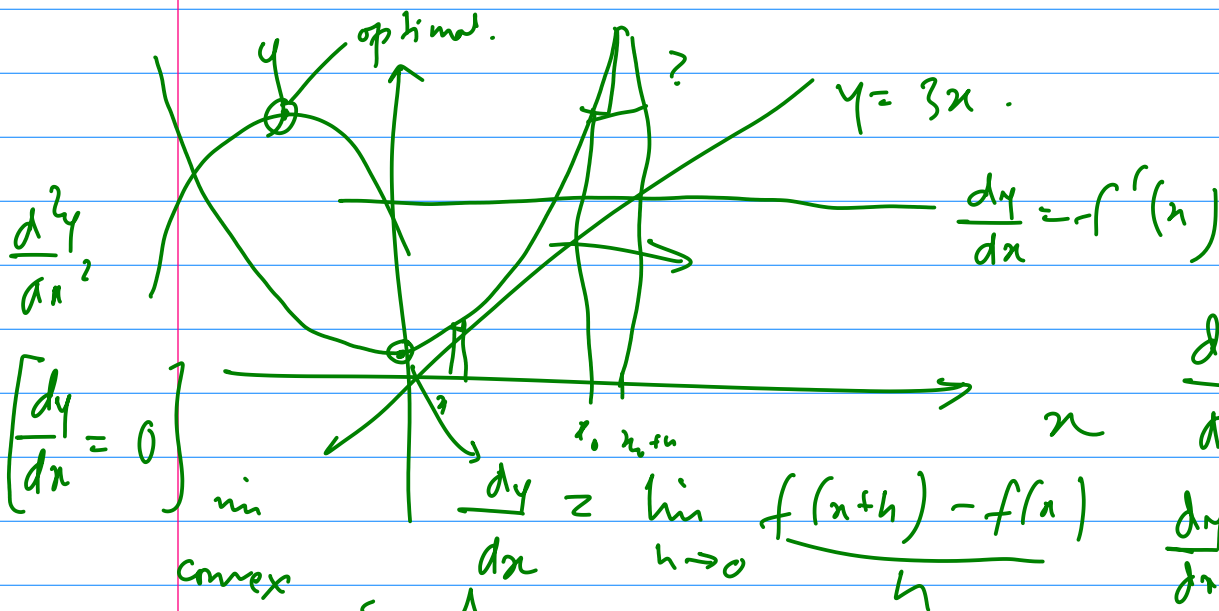
$$\frac{d \log_2 x}{dx} = ? \frac{1}{x}$$

$\log_e 10$

$\log_e 2$

$$\log_{10} x = c \cdot \log_2 x \\ = c' \cdot \log_e x$$

$$\frac{d(c \cdot f(x))}{dx} = c \cdot \frac{df(x)}{dx}$$



$$\frac{dy}{dx} = f'(x)$$

$$\frac{dy}{dx} = \frac{dx^2}{dx} = 2x$$

$$\left. \frac{dy}{dx} \right|_{x=0} = 2 \cdot 0 = 0$$

convex

$$\frac{dy}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$













