# Machine Learning Flow (Any Algorithm)

*refers to Linear Regression -v2.ipynb

**A** — Find Data

**B** — Load Data

1.1 Load data

**C** — Data Preprocessing

1.2 Explore data

C1

1.3 Selecting variable

1.3 Prepare data

**D** — Split data

- Train data
- Testing data

**E** — Train model on data

- Algorithm

1.4 Train

maths here

**F** — Evaluate model

- 1.5 Evaluate

C2 — Exploratory Data Analysis

C3 — Feature Engineering → missing values, Encoding

One HOT ENCODING

C4 — Feature Scaling

# C3. Feature Engineering
## One Hot Encoding

| | age | sex | bmi | children | smoker | region | charges |
|---|---|---|---|---|---|---|---|
| 0 | 19 | female | 27.90 | 0 | yes | southwest | 16884.9240 |
| 11 | 62 | female | 26.29 | 0 | yes | southeast | 27808.7251 |
| 14 | 27 | male | 42.13 | 0 | yes | southeast | 39611.7577 |
| 19 | 30 | male | 35.30 | 0 | yes | southwest | 36837.4670 |
| 23 | 34 | female | 31.92 | 1 | yes | northeast | 37701.8768 |

Sex Encoding $\begin{cases} 0 : \text{Female} \\ 1 : \text{Male} \end{cases}$

Children Encoding

| 0 | 1 | 2 | 3 | 4 | 5 | Number of child |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 | 0 | 0 | 2 |
| 0 | 0 | 0 | 1 | 0 | 0 | 3 |
| 0 | 0 | 0 | 0 | 1 | 0 | 4 |
| 0 | 0 | 0 | 0 | 0 | 1 | 5 |

| age | sex | bmi | OHE_1 | OHE_2 | OHE_3 | OHE_4 | OHE_5 |
|---|---|---|---|---|---|---|---|
| 19 | 0 | 27.90 | 0 | 0 | 0 | 0 | 0 |
| 34 | 0 | 31.92 | 1 | 0 | 0 | 0 | 0 |

# D.) Split the data

1. First shuffle the rows of the data

| | Bmi | Charges |
|---|---|---|
| 1 | 27 | 10,000 |
| ⋮ | | |
| 9337 | | |
| 9338 | 30 | 50,000 |
| 338 | | |

Train data

2. X_train is $0.7 \times 1338 = 9337$ rows

Y_train is $= 9337$ rows

Test data

3. X_test is $0.3 \times 1338 = 401$ rows

Y_test is $0.3 \times = 401$ rows

Smoker only has 274 rows
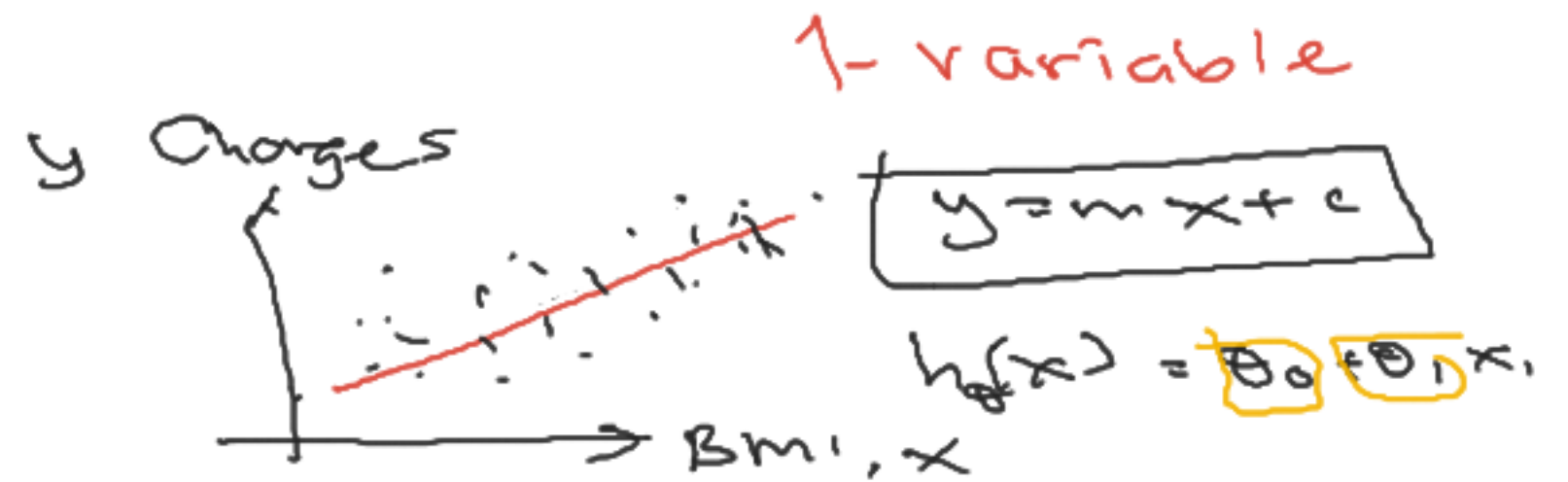
$$X\_train = 0.7 \times 274 = 191$$

$$X\_test = 0.3 \times 274 = 82$$

# Linear Regression (LR)

└→ 1) One variable LR

2) Multi variable LR



1-variable

$$y = mx + c$$

$$h_\theta(x) = \theta_0 + \theta_1 x_1$$

$$y = m_1 x_1 + m_2 x_2 + c$$

2 variable

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

Q How about 11 variables? What does the eq^n looks like?

A) $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \ldots + \theta_{11} x$

where $m_1, m_2, m_3, m_i$

# Linear Regression (LR)

- We are given this insurance.csv dataset.

- We want to do some LR on this dataset.

- for example. Can we predict how much a person has to pay for his/her insurance based on their BMI ?

| | age | sex | bmi | children | smoker | region | charges |
|---|---|---|---|---|---|---|---|
| 0 | 19 | female | 27.900 | 0 | yes | southwest | 16884.92400 |
| 1 | 18 | male | 33.770 | 1 | no | southeast | 1725.55230 |
| 2 | 28 | male | 33.000 | 3 | no | southeast | 4449.46200 |
| 3 | 33 | male | 22.705 | 0 | no | northwest | 21984.47061 |
| 4 | 32 | male | 28.880 | 0 | no | northwest | 3866.85520 |

# How to look at dataset

(input) Features, variable, dimension, attribute    target output,

| | age | sex | bmi | children | smoker | region | charges |
|---|-----|-----|-----|----------|--------|--------|---------|
| 0 | 19 | female | 27.900 | 0 | yes | southwest | 16884.92400 |
| 1 | 18 | male | 33.770 | 1 | no | southeast | 1725.55230 |
| 2 | 28 | male | 33.000 | 3 | no | southeast | 4449.46200 |
| 3 | 33 | male | 22.705 | 0 | no | northwest | 21984.47061 |
| 4 | 32 | male | 28.880 | 0 | no | northwest | 3866.85520 |

rows
= observations/
Samples

$i^{th}$ observation

## Continuous variable

1. Age
2. Bmi
3. Charges

## Categorical variable

1. Sex
2. Children
3. Smoker
4. region

# Revision of Linear Algebra

$x$, $x = 5$ : scalar

$x$, $x = [1 \quad 2 \quad 3]$ $\rightarrow$ $x$ is a $1 \times 3$ row vector

(3 col, 1 row)

row × col

$$x = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$ (column) vector

3 rows, 1 col

$x$ is a $3 \times 1$ column vector

row × col

$$X = \begin{bmatrix} 1 & 10 & 100 \\ 2 & 20 & 200 \\ 3 & 30 & 300 \end{bmatrix}$$

3 row, 3 col

$X$ is a $3 \times 3$ matrix

Continue

$$x_1 = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \qquad x_2 = \begin{bmatrix} 10 \\ 20 \\ 30 \end{bmatrix} \qquad x_3 = \begin{bmatrix} 100 \\ 200 \\ 300 \end{bmatrix}$$

$$X = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 10 & 100 \\ 2 & 20 & 200 \\ 3 & 30 & 300 \end{bmatrix}$$

Now when we see a dataset, lets imagine it as vectors and matrices

| | $x_1$ age | $x_2$ sex | $x_3$ bmi | $x_4$ children | $x_5$ smoker | $x_6$ region | $y$ charges |
|---|---|---|---|---|---|---|---|
| 0 | 19 | female | 27.900 | 0 | yes | southwest | 16884.92400 |
| 1 | 18 | male | 33.770 | 1 | no | southeast | 1725.55230 |
| 2 | 28 | male | 33.000 | 3 | no | southeast | 4449.46200 |
| 3 | 33 | male | 22.705 | 0 | no | northwest | 21984.47061 |
| 4 | 32 | male | 28.880 | 0 | no | northwest | 3866.85520 |

$x_i$ : the $i^{th}$ feature/ attributes of the dataset . $x_i$ is __a col vector__

$y$ : a col vector of the target output

$$X = \begin{bmatrix} x_1 & x_2 & \cdots & x_6 \\ 19 & & & \\ 18 & & & \\ \vdots & \vdots & & \vdots \\ & & & \\ & & & \end{bmatrix} \begin{matrix} 0 \\ \\ \\ \\ \\ 1337 \end{matrix}$$

- we have 6 features. But how many rows of observation/samples do we have?

cols = 6
rows = 1338

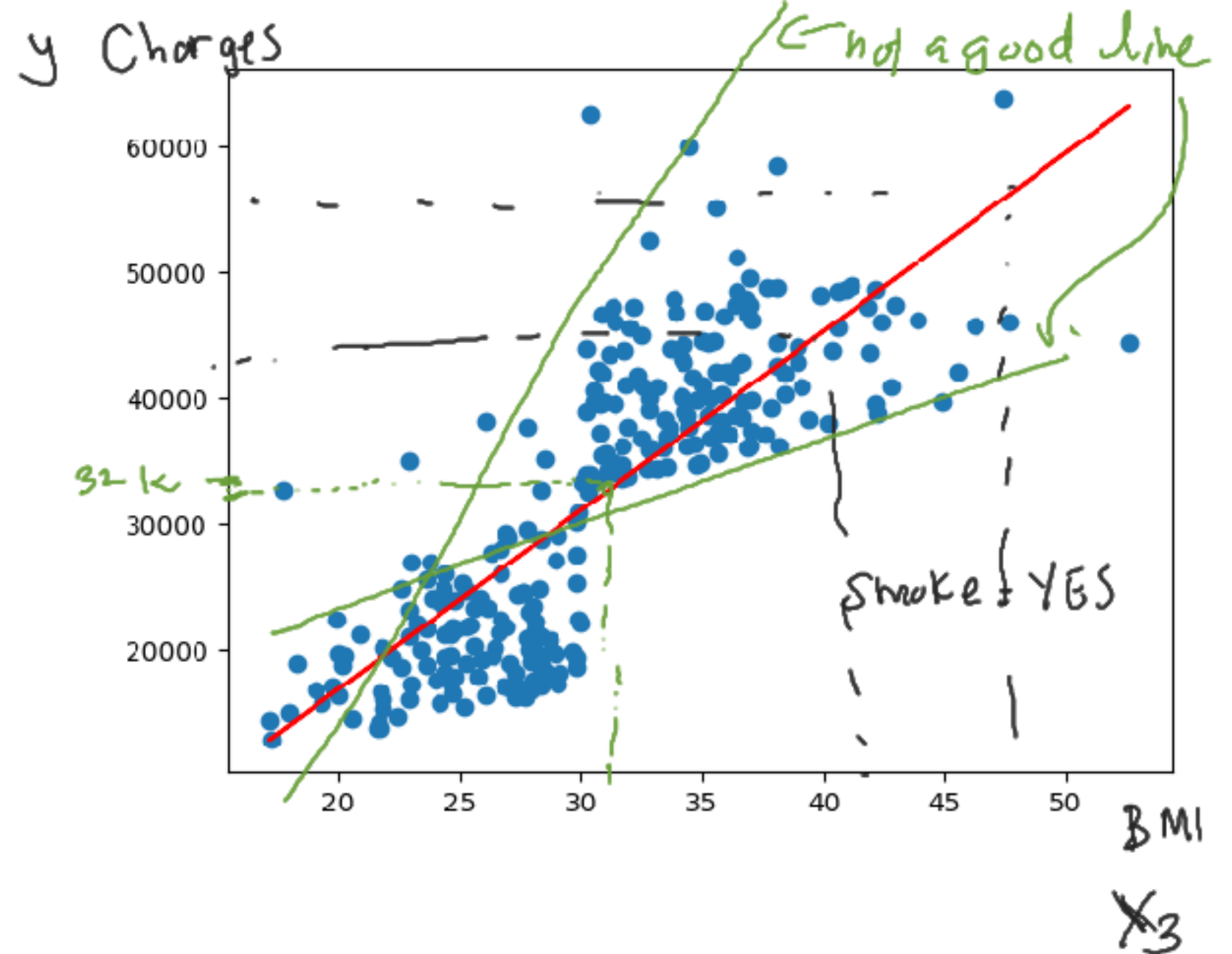$X$ is a $1338 \times 6$ matrix #

Our 'intuition'



| | age | sex | bmi | children | smoker | region | charges |
|---|---|---|---|---|---|---|---|
| 0 | 19 | female | 27.900 | 0 | yes | southwest | 16884.92400 |
| 1 | 18 | male | 33.770 | 1 | no | southeast | 1725.55230 |
| 2 | 28 | male | 33.000 | 3 | no | southeast | 4449.46200 |
| 3 | 33 | male | 22.705 | 0 | no | northwest | 21984.47061 |
| 4 | 32 | male | 28.880 | 0 | no | northwest | 3866.85520 |

$X_1$   $X_2$   $X_3$   $X_4$   $X_5$   $X_6$   $y$

$X$: features / attributes, *independent* variable

$y$: target output, *dependant* variable

# Linear Regression with one variable

| | $X_3$ BMI | $y$ Charges |
|---|---|---|
| 1 | 29 | 50000 |
| 2 | -30 | 60000 |
| ⋮ | | |
| 1338 | 25 | 70000 |
| 1337 | 31 | ? |



- Linear regression is to find <u>a linear line</u> that fits the dataset.

- We can make predictions for unseen obscruction. For example what is the insurance charges for a person with BMI = 31?

How do we find the best line (red)
that fits the data?

$$y = mx + c$$

$$h(x) = \theta_0 + \theta_1 x$$

C: intercept    m: slope

~ So our job (for one variable) is to solve
for m and c.



y
Charges

60000

50000

40000

30000

20000

y=mx+c

m

Smoker: YES

C

20   25   30   35   40   45   50

X    BMI
     Xs

How do we find the 'best' m and c ?

$$y = mx + c$$

Two methods

A) Solve using equation

B) Use gradient descent

# A) Solve using equation

## Ordinary Least Square method (from Stats)

$$c, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x};$$

$$m, \hat{\beta}_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2},$$

$$y = mx + c$$

$$y = \beta_0 + \beta_1 x$$

← what is this?



y, charges

BMI
x

## From our dataset

| index | BMI $x_3$ | Charges $y$ |
|-------|-----------|-------------|
| 1 | 19 | |
| 2 | 23 | |
| $i$ | $x_i = 25$ | $y_i = 20000$ |

$n = 1338$

$$\bar{x} = \text{averge. of } x = \frac{1}{N}\sum_{i=1}^{N} x_i$$

$$\bar{y} = \text{'' of } y = \frac{1}{N}\sum_{i=1}^{N} y_i$$

# B.) Use gradient Descent

<span style="color:green">(gradient of the cost function)</span>

We need to define a cost function
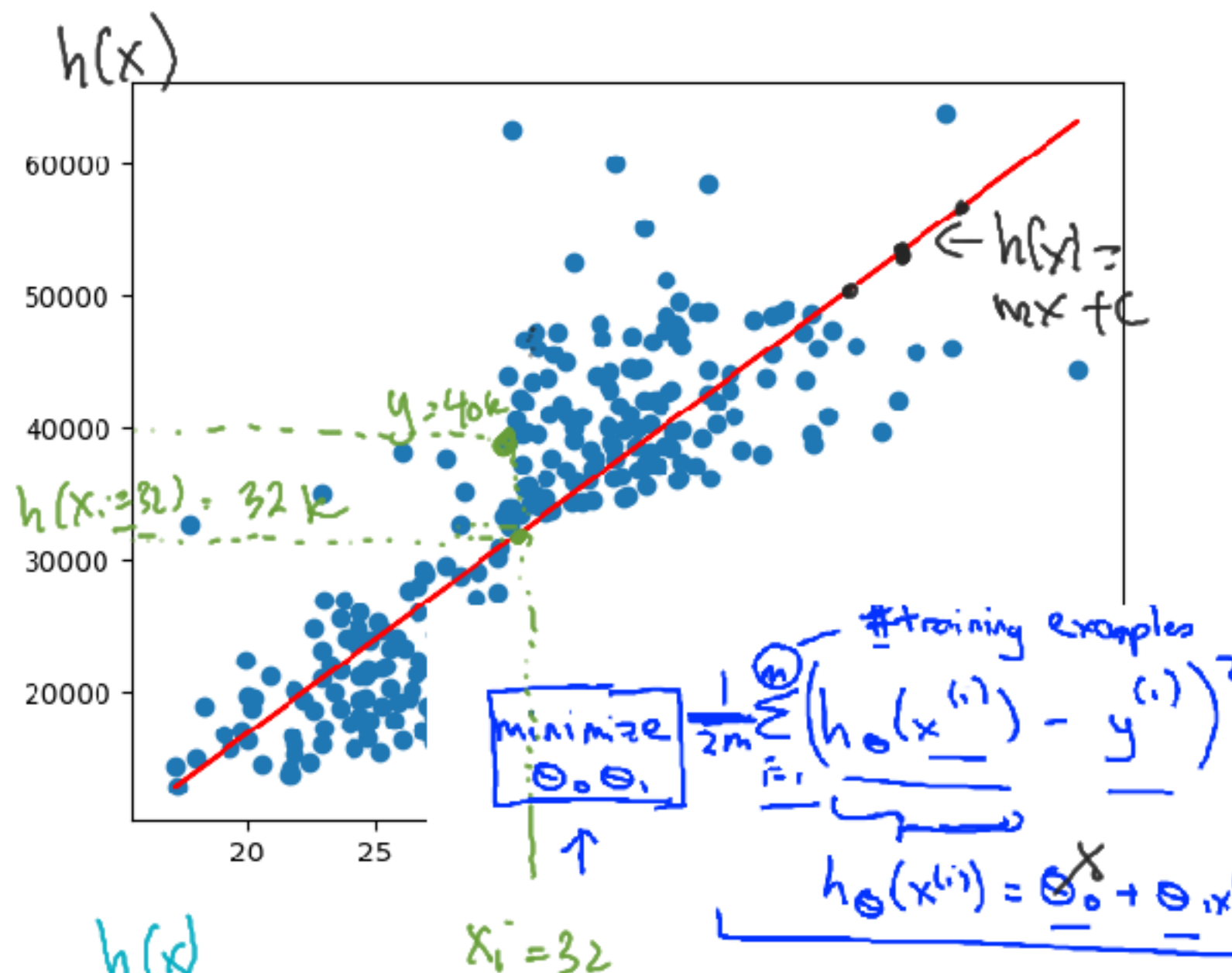
$h(x) = mx + c$ ← our red line

● : our data   ● : our prediction (line)

$$\left(h(x^i) - y^i\right)^2$$

∴ So our cost function is

$$J(c,m) = \frac{1}{2N} \sum_{i=1}^{N} \left(h(x^i) - y^i\right)^2$$



h(x)

$y = 40k$

$h(x_i = 32) = 32k$

minimize $\frac{1}{2m} \sum_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)^2$   #training examples
$\theta_0, \theta_1$

$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x$

$x_i = 32$

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

minimize $J(\theta_0, \theta_1)$
$\theta_0, \theta_1$

Cost function
Squared error function

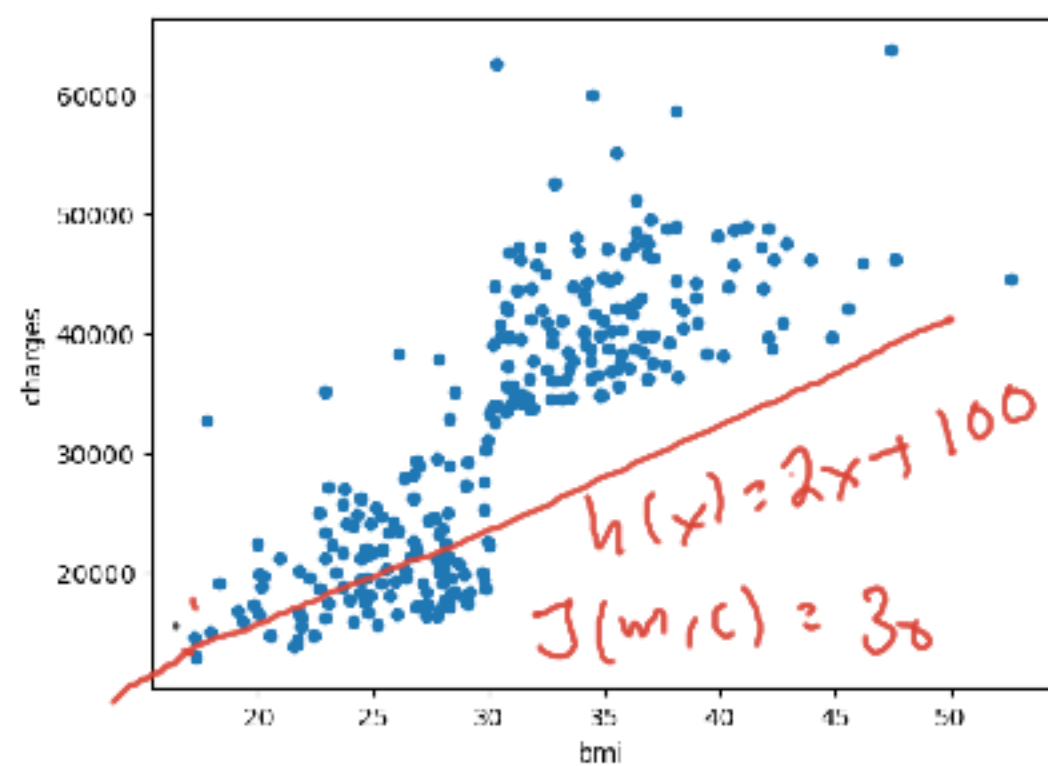| indx | X BMI | y Charges | h(x) |
|------|-------|-----------|------|
| 1 | 27 | 49000 | |
| 2 | 31 | 45,000 | |
| i | $x^i = 32$ | $y_i = 40,000$ | 32,000 |
| N = 1338 | 25 | 38 000 | |

3.) Use gradient Descent
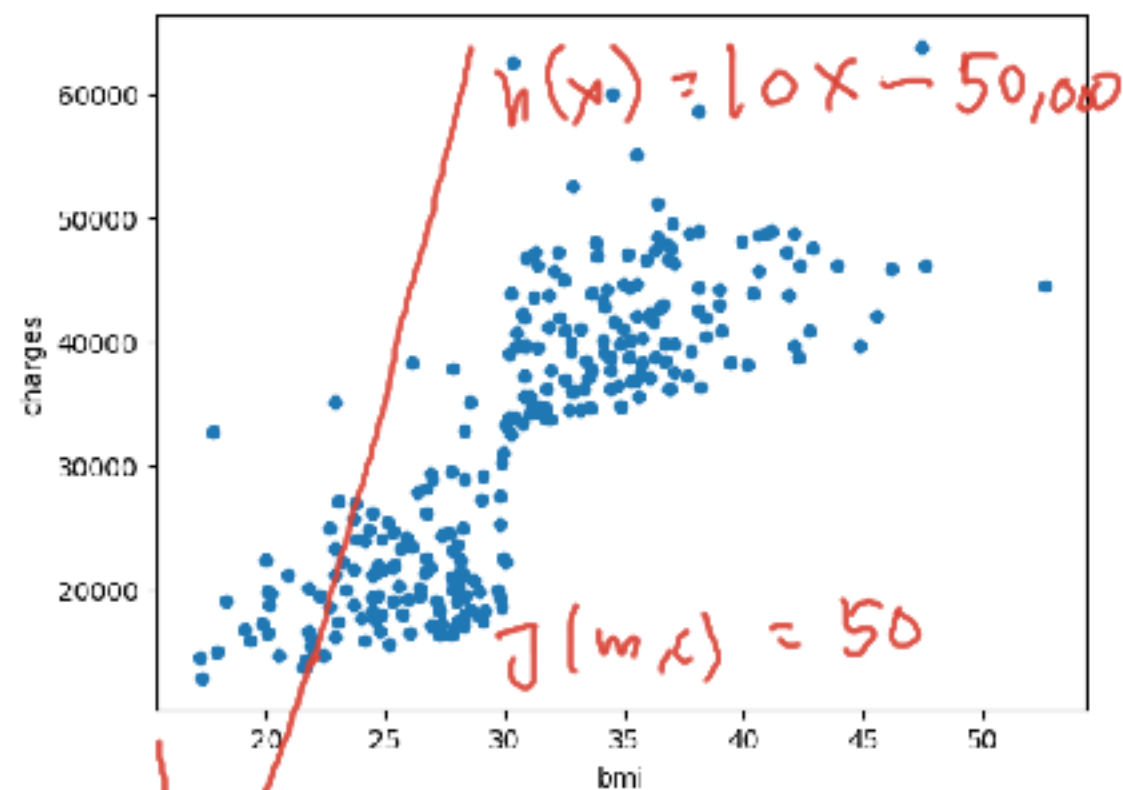
Cost function    $h(x) = mx + c$

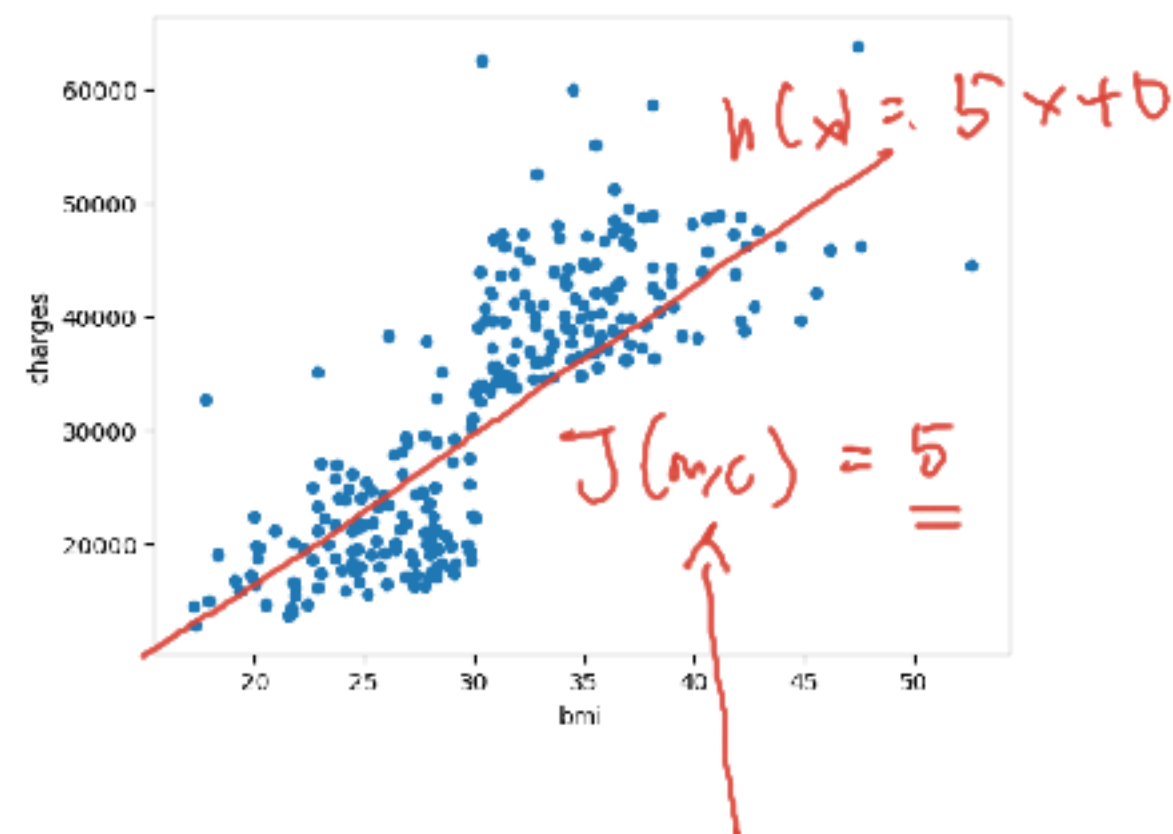$$J(m,c) = \frac{1}{2N} \sum_{i=1}^{N=1338} \left( h(x^i) - y^i \right)^2$$



- The cost function calculates the distance between our predictions $h(x)$ and the original data.

- So then we need to find the values of m and c that minimizes the cost function $J$

$h(x) = 2x + 100$

$J(m,c) = 30$

bad

$R^2 = 0.1$

$h(x) = 10x - 50,000$

$J(m,c) = 50$

really bad

$R^2 = 0.01$

$h(x) = 5x + 0$

$J(m,c) = \underline{5}$

$m = 5, c = 0$
yields the lowest
error/cost function

good

$R^2 = 0.7$ ✓

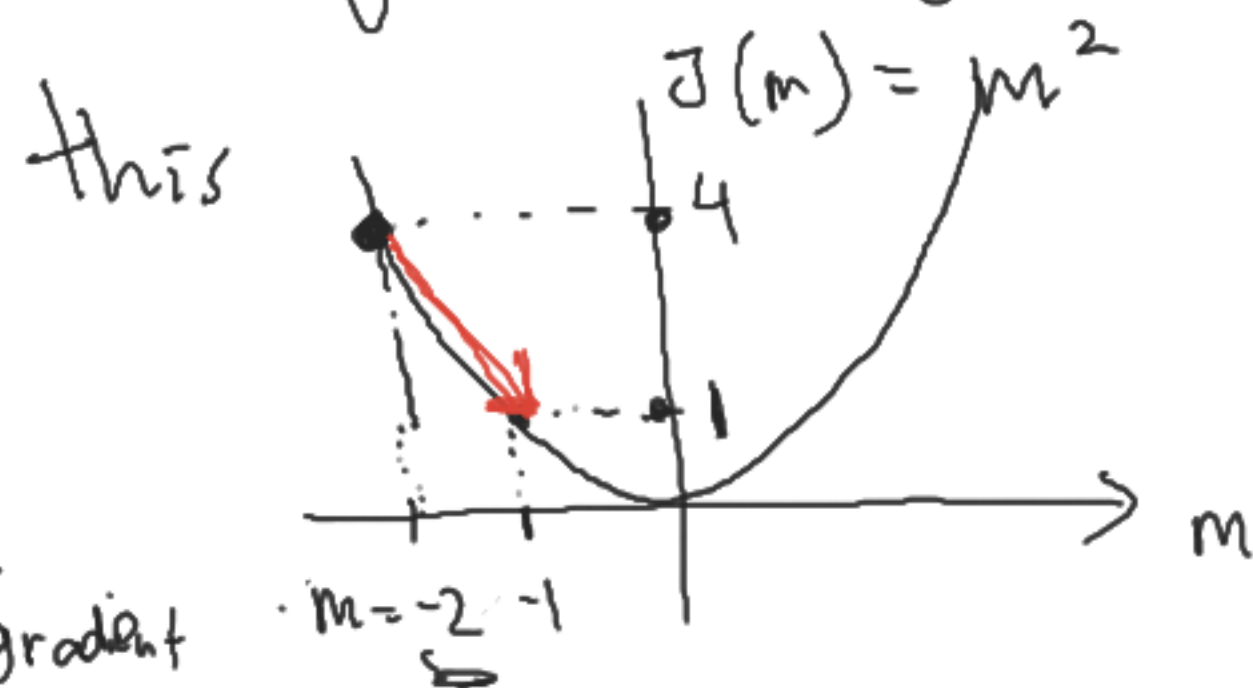How do we find the best m and c that minimizes the cost function?

$$J(m,c) = \frac{1}{2N} \sum_{j=1}^{N=1338} (h(x^i) - y^i)^2$$

what does this cost function looks like? (I don't know)

How do we find the best m and c that minimizes the cost function?

Let's just imagine our cost function looks like this



$J(m) = m^2$

$4$

$1$

$m = -2 \quad -1$

$\alpha$: learning rate/step $= 1/4$

Update using gradient descent

Is it new

$$M_{new} = M_{old} - \alpha \frac{\partial J(m)}{\partial m}$$

$$= -2 - \frac{1}{4} \cdot 2(-2)$$

$$= -1$$

$Y_{new} = M_{new}^2 = 1$

Derivative/gradient

$$J(m) = m^2$$
$$\frac{\partial J(m)}{\partial m} = 2m$$

THIS IS HOW WE GO DOWN THE GRADIENT OF THE COST FUNCTION

We define the cost function

$$\text{(1)} \quad J(m,c) = \frac{1}{2N} \sum_{j=1}^{N=1338} (h(x^i) - y^i)^2$$

y Charges

$$h(x) = mx + c$$

BMI

$J(m,c)$

m~5

C=0

c

m

Cnew =

(2) we find the best value of m and c which minimizes cost function J.

(3) This is done by going down the derivative/gradient of J

$$M_{new} = M_{old} - \alpha \frac{\partial J(m,c)}{\partial m}$$