# Linear Regression
## (Multivariate)

i/p. output

| $x_1$ | $x_2$ | ... | $x_n$ | $y$ |
|---|---|---|---|---|
| | | | | |

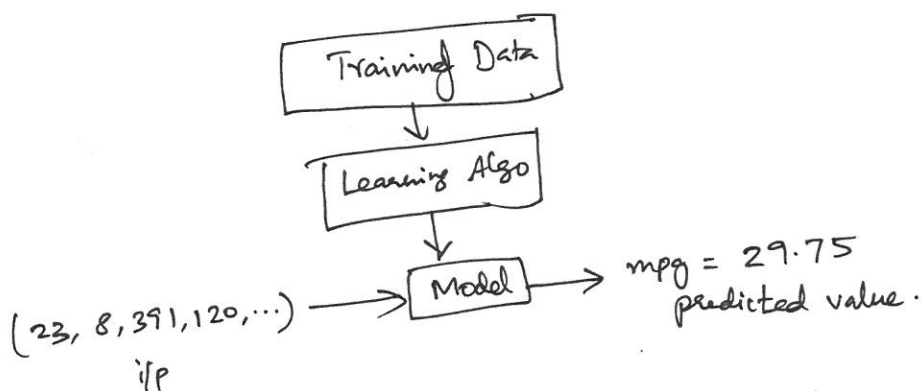**Task:** Given $(x_1, x_2, \ldots x_n)$ we need to predict the value of $y$.

Real values.

**Eg:**

- predict mpg value given engine parameters.

-

| cyclinder | Horse Power | Height | Weight | ... | MPG |
|---|---|---|---|---|---|
| 22.0 | 8 | 390 | 190 | ... | ? |
| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_d$ | |

$N = \# No.$ of instances $= 398$ (rows).

predict mpg of a vehicle.

Training Data → Learning Algo → Model → mpg = 29.75 predicted value.

$(23, 8, 391, 120, \ldots)$ ⟶ i/p

**Hypothesis fn of Regression fn.:**

$$\hat{y}_i(x) = b_0 + b_1 x_{i1} + b_2 x_{i2} + \cdots + b_d x_{id}$$

| $x_1$ | $x_2$ | ... | $x_d$ | $y$ |
|---|---|---|---|---|
| | | | | |

Data set.

$x_1 = (x_{11}, x_{12}, \ldots x_{1d})$.

$x_2 = (x_{21}, x_{22}, \ldots x_{2d})$

- Need to find $(b_0, b_1, b_2, \ldots, b_d)$ parameters?
- Use training data and (learning algorithm) to compute these parameters.
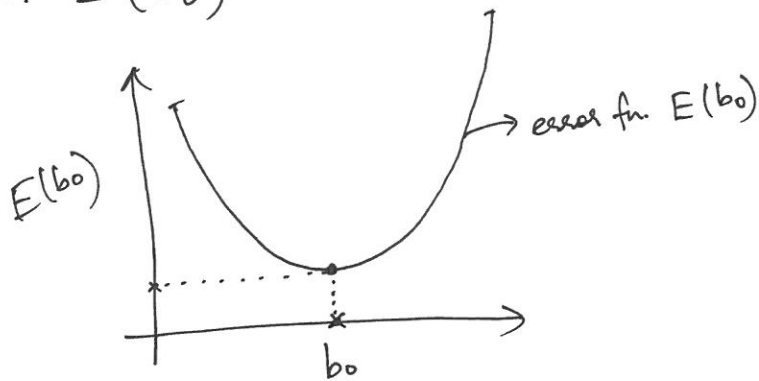
cost fn.?

$$E(b_0, b_1, \ldots b_d) = \frac{1}{2N} \sum_{i=1}^{N} (\hat{y}_i - y_i)^2$$
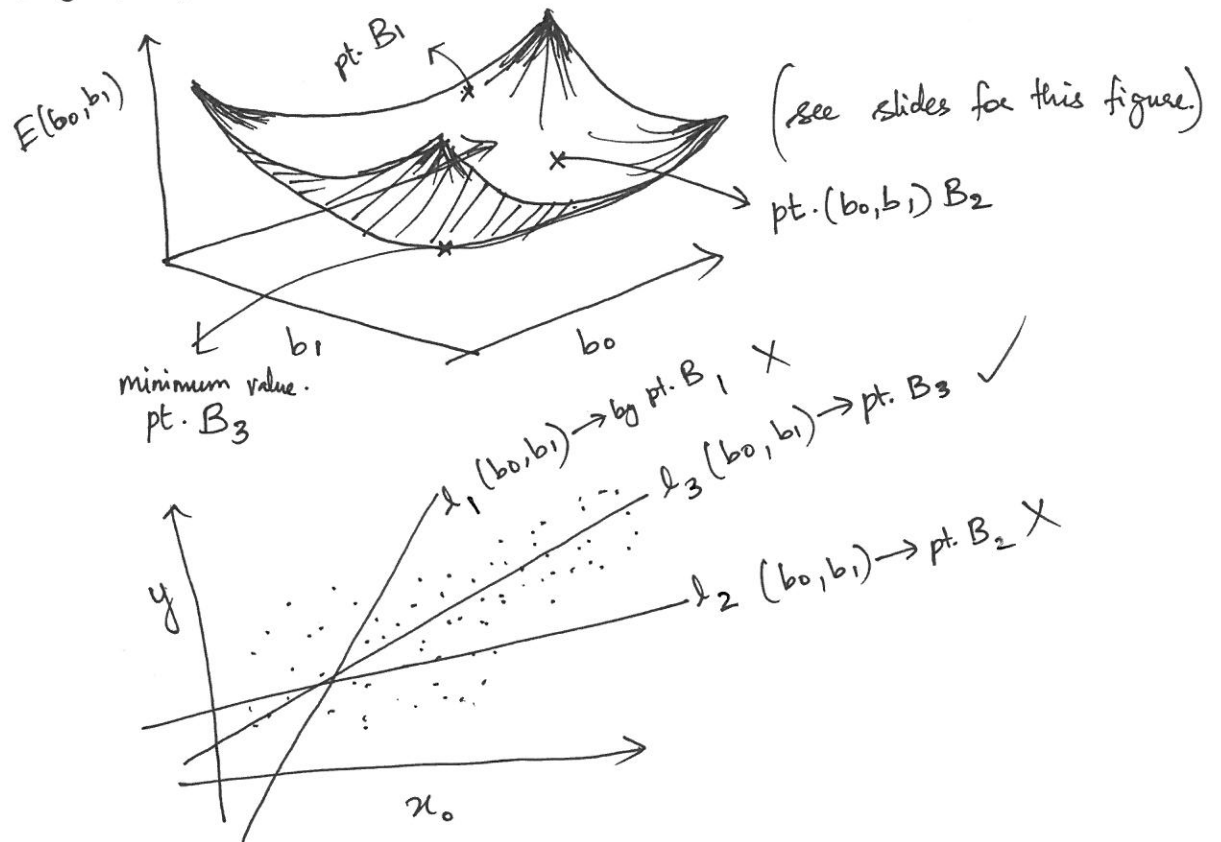
**Goal** minimize $E(\cdots)$ for different values of $(b_0, b_1, \ldots b_d)$

- squared error b/w $y_i$ (from training data and $\hat{y}_i$ (from predicted data using param $b_0, \ldots, b_d$).
- squared error cost fn.

:) Cost / Error fn: $E(b_0) \longrightarrow$ 1 param.



$E(b_0)$

$\longrightarrow$ error fn. $E(b_0)$

$b_0$

$E(b_0, b_1) \longrightarrow$ 2 param.



$E(b_0, b_1)$

pt. $B_1$

(see slides for this figure.)

pt. $(b_0, b_1)$ $B_2$

$b_1$

$b_0$

minimum value.
pt. $B_3$



$\ell_1 (b_0, b_1) \longrightarrow$ by pt. $B_1$ $\times$

$\ell_3 (b_0, b_1) \longrightarrow$ pt. $B_3$ $\checkmark$

$\ell_2 (b_0, b_1) \longrightarrow$ pt. $B_2$ $\times$

$y$

$x_0$

$E(b_0, b_1, \dots b_d) \longrightarrow$ d-param.

Goal $\longrightarrow$ minimize $E(b_0, b, \dots b_d)$
$(b_0, b_1, \dots b_d)$

$\downarrow$

learning Algorithm?
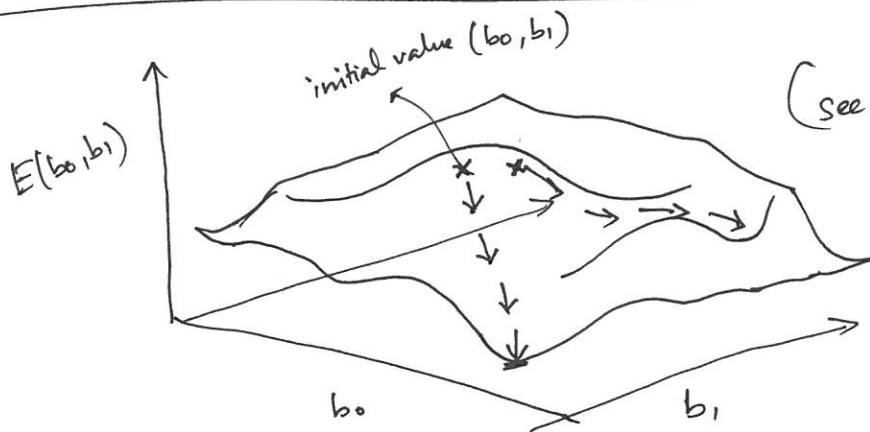
$\boxed{\text{Gradiant Decent}}$

③

# Gradiant Descent

Outline :

1. Start with some $(b_0, b_1)$

or $b_0 = 0$
$b_1 = 0$

2. Keep changing $(b_0, b_1)$ value to reduce $E(b_0, b_1)$, until we end up in min. value

given $E(b_0, b_1)$

Goal $\rightarrow \underset{(b_0, b_1)}{\min} E(b_0, b_1)$

you can Extend it to $(b_0, b_1, \ldots b_d)$

initial value $(b_0, b_1)$

$E(b_0, b_1)$

$b_0$          $b_1$

(see slides for this figure)

---

Math :

no change in $b_j$ values.

repeat until <u>convergence</u> $\{$

$a := b$
$\downarrow$
assignment.

$$b_j := b_j - \alpha \frac{\partial}{\partial b_j} E(b_0, b_1)$$

(for $j = 0$ and $j = 1$)

$\}$

learning rate
or
update rate

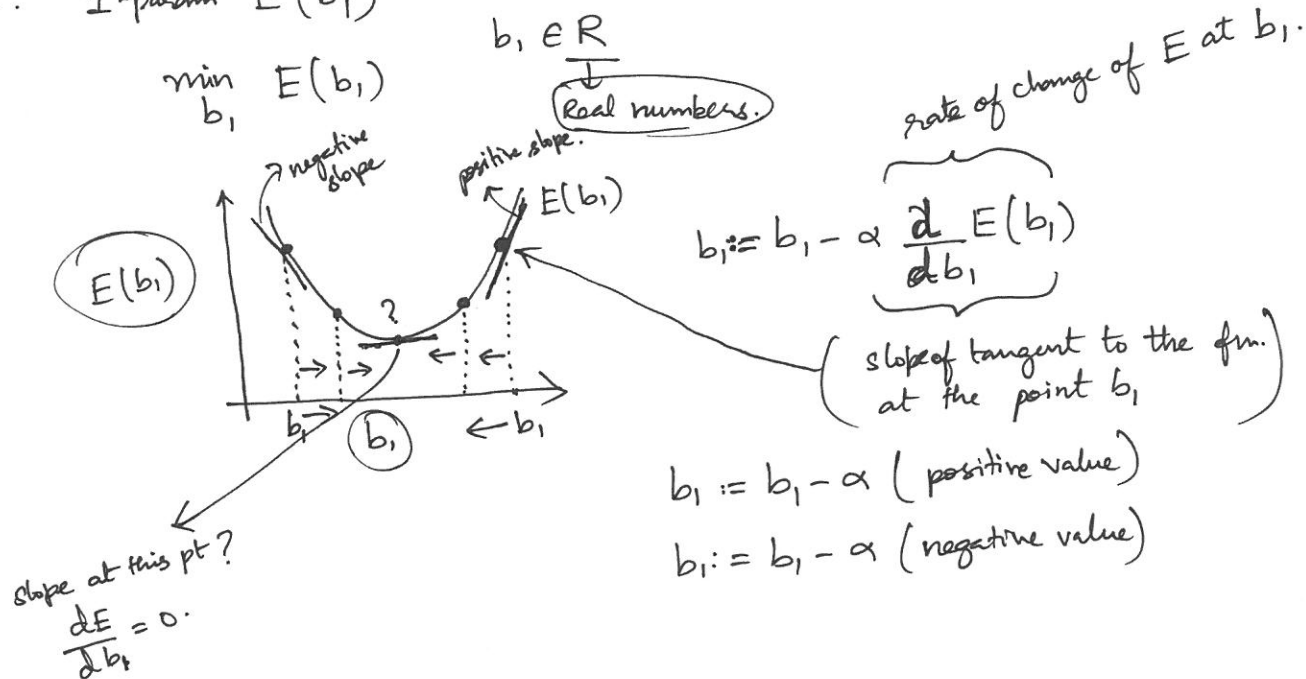partial derivative of $E$ w.r.t $b_j$

simultaneous update

$\{$ $tmp_0 := b_0 - \alpha \frac{d}{db_0} E(b_0, b_1)$

$tmp_1 := b_1 - \alpha \frac{\partial}{\partial b_1} E(b_0, b_1)$

$b_0 := tmp_0$
$b_1 := tmp_1$

$\}$

④

Example: 1-param $E(b_1)$

$$\min_{b_1} E(b_1)$$
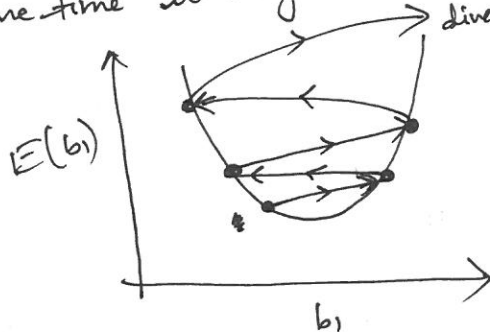
$b_1 \in \mathbb{R}$ — Real numbers.

rate of change of $E$ at $b_1$.

$$b_1 := b_1 - \alpha \underbrace{\frac{d}{db_1} E(b_1)}$$

( slope of tangent to the fn. at the point $b_1$ )

$b_1 := b_1 - \alpha \,(\text{positive value})$

$b_1 := b_1 - \alpha \,(\text{negative value})$



$E(b_1)$ — negative slope — positive slope. — $E(b_1)$

$b_1 \rightarrow$   (b_1)   $\leftarrow b_1$

slope at this pt?
$$\frac{dE}{db_1} = 0.$$

---

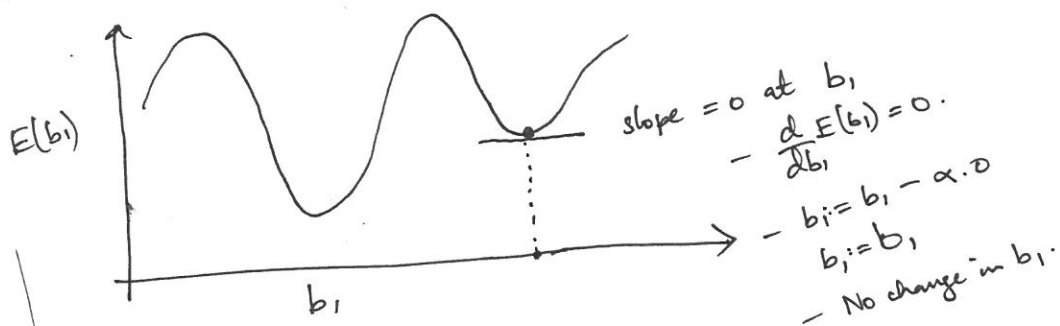$\alpha$-value: (0 to 1) eg: 0.001 → ~~begin~~ trial-and-error.

if $\alpha$ is too small, GD can be slow.



$E(b_1)$ … $b_1$

- if $\alpha$ is too large, GD can overshoot minimum.
- some time it may cause GD to fail to converge.
  → diverge.



$E(b_1)$ … $b_1$

Local Minimum:



$E(b_1)$ … $b_1$

slope = 0 at $b_1$
- $\frac{d}{db_1} E(b_1) = 0.$
- $b_1 := b_1 - \alpha \cdot 0$
  $b_1 := b_1$
- No change in $b_1$.

$E(b_1)$

$b_1$

- As we approach local minimum, GD will automatically takes smaller step. decrease
- So no need to adjust $\alpha$ value over time.

**Apply GD algo in Linear Regression**

**GD**

**L.R (Multivariate)**

repeat until convergence

$\{$

$\quad b_j := b_j - \alpha \dfrac{\partial}{\partial b_j} E(b_0, b_1)$

$\quad$ (for $j=0$ and $j=1$)

$\}$

$y(x) = b_0 + b_1 x$

$E(b_0, b_1) = \dfrac{1}{2N} \sum\limits_{i=1}^{N} (\hat{y}_i - y_i)^2$

minimize $E(b_0, b_1)$
$(b_0, b_1)$

$\dfrac{\partial}{\partial b_j} E(b_0, b_1) = \dfrac{\partial}{\partial b_j} \left( \dfrac{1}{2N} \sum\limits_{i=1}^{N} (\hat{y}_i - y_i)^2 \right.$

$= \dfrac{\partial}{\partial b_j} \dfrac{1}{2N} \sum\limits_{i=1}^{N} \left( (b_0 + b_1 x) - y_i \right)^2$

$j=0 \rightarrow \dfrac{\partial E(b_0, b_1)}{\partial b_0} = \dfrac{1}{N} \sum\limits_{i=1}^{N} (\hat{y}_i - y_i)$

$j=1 \rightarrow \dfrac{\partial E(b_0, b_1)}{\partial b_1} = \dfrac{1}{N} \sum\limits_{i=1}^{N} (\hat{y}_i - y_i)(x_i)$
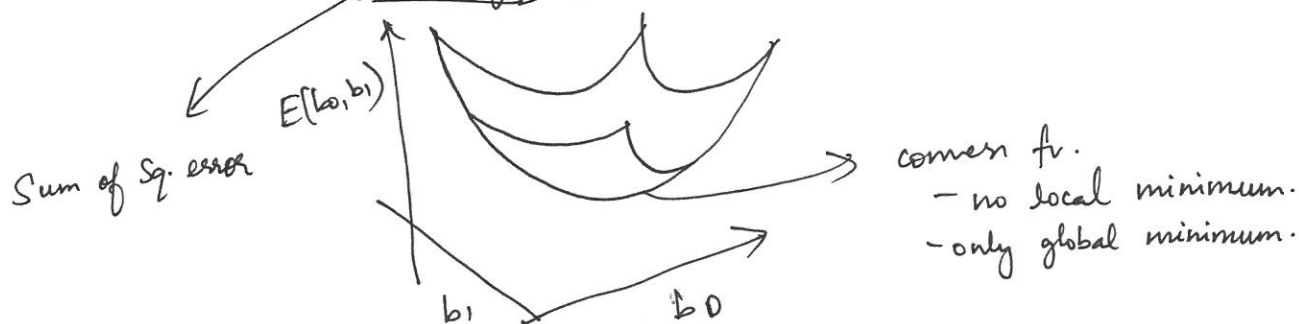
Eq. 1

GD. for LR:

repeat until convergence

$\{ \quad b_0 := b_0 - \alpha \dfrac{1}{N} \sum\limits_{i=1}^{N} (\hat{y}_i - y_i) \qquad \| $ simultaneous update.

$\quad b_1 := b_1 - \alpha \dfrac{1}{N} \sum\limits_{i=1}^{N} (\hat{y}_i - y_i)(x_i)$

$\}$

(6)

Note:
- GD suffers from local minimum solutions.
- The **Error fn** for Linear Regression is always <u>CONVEX</u> fn.

$E(b_0, b_1)$

Sum of Sq. error

convex fn.
- no local minimum.
- only global minimum.

$b_1$        $b_0$

---

## GD for Multivariate Variables (Linear Regression)

- The example we saw is for 1-variable $x$, we have $b_0$ and $b_1$ param to be estimated

- if more than 1-variable, then (Eq.1) does not hold.

$\therefore$ for $X = (x_0, x_1, x_2, \ldots x_d)^T$, $\boxed{y = b_0 x_0 + b_1 x_1 + \cdots + b_d x_d}$

always $x_0 = 1$

| | $x_1$ | $x_2$ | $\ldots x_d$ | $y$ |
|---|---|---|---|---|
| 1 | $x_{11}$ | $x_{12}$ | $\ldots x_{1d}$ | $y_1$ |
| 2 | | | | |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $n$ | $x_{n1}$ | $x_{n2}$ | $\ldots x_{nd}$ | $y_n$ |

$\}$ training samples

#$n$ samples.

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \qquad B = \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_d \end{bmatrix} \qquad \begin{array}{l} X \in R^{d+1} \\ B \in R^{d+1} \end{array}$$

$$y = b_0 x_0 + b_1 x_1 + \cdots + b_d x_d$$

$$= B^T \cdot X$$

$$= \underbrace{\left[ b_0 \ b_1 \ \ldots \ b_d \right]^T}_{\substack{(d+1) \times 1 \\ \text{matrix}}} \cdot \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}}_{\substack{1 \times (d+1) \\ \text{matrix.}}}$$

$\boxed{y = B^T \cdot X}$ Regression Eq. for Multiple Variables.

Cost / Error fn:

$$E(B) = \frac{1}{2N} \sum_{i=1}^{N} (\hat{y}_i - y_i)^2$$

$\downarrow$

$(b_0, b_1, \ldots b_d)$

Repeat until convergence {

$$b_j := b_j - \alpha \frac{1}{N} \sum_{i=1}^{N} (\hat{y}_i - y_i) x_{ij}$$

(simultaneously update $b_j$ for $j = 0, 1, \ldots, d$)

}

Example:

$$b_0 := b_0 - \alpha \frac{1}{N} \sum_{i=1}^{N} (\hat{y}_i - y_i) \underbrace{x_{i0}}_{} \rightarrow \text{always 1.}$$

$$b_1 := b_1 - \alpha \frac{1}{N} \sum_{i=1}^{N} (\hat{y}_i - y_i) x_{i1}$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots$$

$$b_d := b_d - \alpha \frac{1}{N} \sum_{i=1}^{N} (\hat{y}_i - y_i) x_{id}$$

---

GD for Linear Reg (Practical ideas to improve in Real time)
(performance)

1. Make sure features are on a similar scale.

Eg: $x_1 = \text{size} (0 - 2000 \text{ sq.ft})$ $\rbrace$ → or standardize → need to normalize to 0 to 1.

$x_2 = \# \text{ of rooms } (1-5)$

– Make every features are in range $(-1 \le x_i \le 1)$ roughly.

$0 \le x_1 \le 3$ ✓

$-2 \le x_2 \le 0.5$ ✓

$-100 \le x_3 \le 100$ ✗

$-0.0001 \le x_4 \le 0.0001$ ✗



2. Plot $E(B)$ graph

– $E(B)$ should decrease after every iterations.

(8.)

## Gradient Descent – Not Working.



$J(\beta)$
$E(\beta)$

#iteration.

$\longrightarrow$ check No bugs in code.

$\longrightarrow$ change $\alpha$ value.
(Smaller $\alpha$ value)

$E(\beta)$

#iterations

— if $\alpha$ is too small, then GD will be very slow to converge.
— if $\alpha$ is too large, $E(\beta)$ will not decrease for every iterations.
                                                                    $\downarrow$
                                                              may not converge.

— try $\alpha = 0.001, 0.01, 0.1, 1, \ldots$ then vary
                      $\downarrow$                        the $\alpha$-value in this interval
            0.003, 0.002