

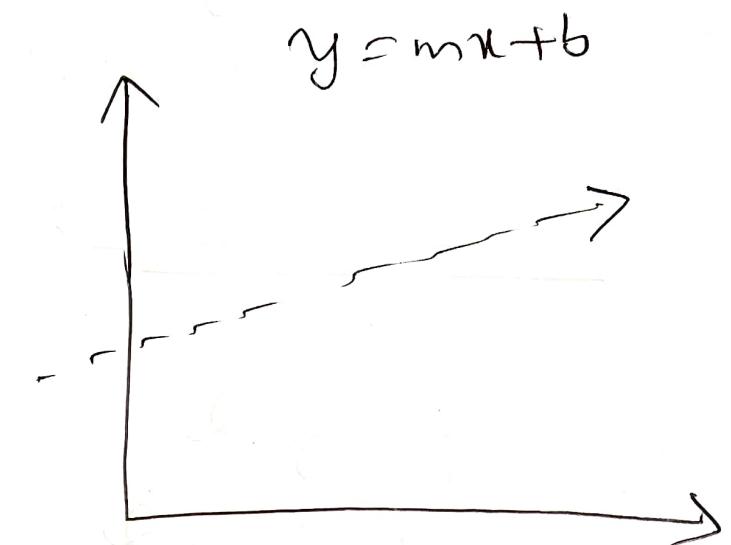
LR with Gradient Descent

①

① $x \rightarrow$ IP data

from x I try
to make a guess

$x \rightarrow \text{guess}$



$x \xrightarrow{\text{Mf}} \text{Y} \rightarrow$ original value
 $x \xrightarrow{\text{Avg}} \text{Avg}$

$$\boxed{\text{Error} = y - \text{guess}} \rightarrow \text{Cost func}$$

$$\boxed{\text{COST} = \sum_{i=1}^N (y_i - \text{Avg}_i)^2}$$

↳ Total error for the particular model, ~~not~~

where y_i are count values of m , &
'b' denotes the particular line.

Goal: minimize the loss or cost
(or I want the smallest fnc).

Mimimizing a fnc: error

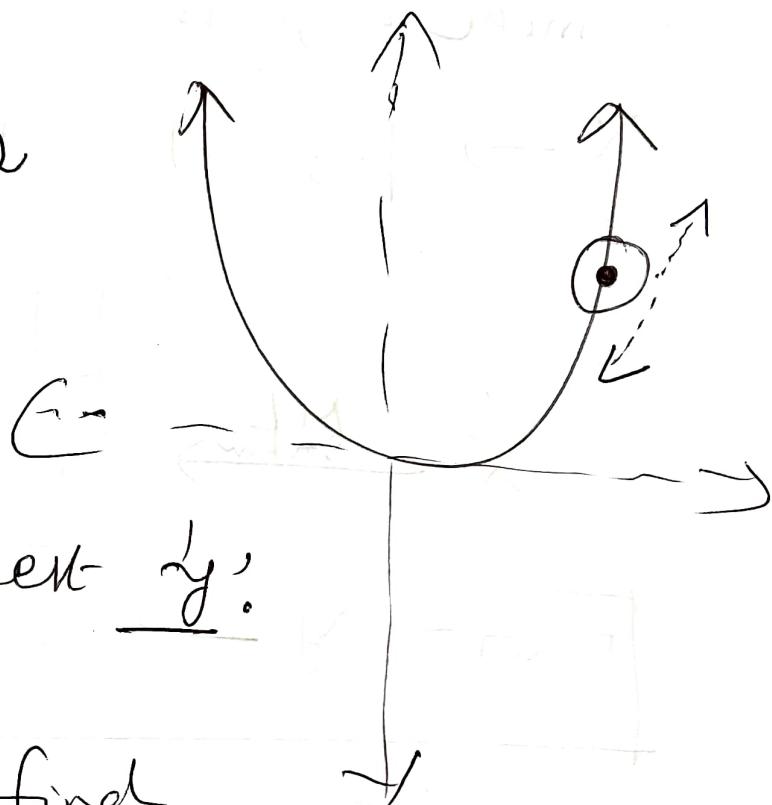
e.g. $y = x^2$

(or) $f(x) = x^2$

i.e.

Find a 'x'
value that

gives the lowest y:



How could I find

(the minimum)

(1) Direction of movement from given point (current or desired)

(2) Step-size = big or small

Noisy Signal

@

Using CI to find the optimal value
of the intercept

Sol) In this case the $R^2 = 0$

$$SS \text{ Residuals} = (\text{Observed} - \text{Predicted})^2$$

$$\begin{aligned} SSRY &= (1.4 - (\text{intercept} + 0.64 \cancel{\times \text{height}}))^2 \\ &= (1.4 - (\text{intercept} + 6.64 + 0.5 \cancel{\times}))^2 \end{aligned}$$

Now we can plug in any value
for intercept and get a new
predicted height.

$$\begin{aligned} \text{Total } SSR &= (1.4 - (\text{intercept} + 0.64 \cancel{\times \text{height}}))^2 \\ &\quad + (1.9 - (\text{intercept} + 0.64 + 2.3))^2 \\ &\quad + (3.2 - (\text{intercept} + 0.64 + 2.9))^2 \end{aligned}$$

This gives an eqn. In this case

Step 01: Take the derivative of SSR's w.r.t. the

'Intercept Term'.

$$\frac{d}{d \text{intercept}} (\text{SSR}^{\text{total}}) = d(\text{SSR}_1) + d(\text{SSR}_2) + d(\text{SSR}_3)$$

Apply chain rule

$$\frac{d}{d \text{intercept}} (\text{SSR}_1) = \frac{d}{d \text{intercept}} (1.4 + (-1) \text{intercept} - 0.64 \times 10^{-5})^2$$

~~$$= \frac{d}{d \text{intercept}} (0 - 2 \text{intercept}) = -2$$~~

~~$$= -2(1.4 - \text{intercept} + 0.64 \times 10^{-5})$$~~

$$\therefore \frac{d}{d \text{intercept}} (\text{SSR}_1)$$

$$= -2(1.4 - \text{intercept} + 0.64 \times 10^{-5})$$

Step: Take the derivative of SSR's w.r.t. the intercept term.

$$\text{Residual} = (\text{Observed} - \text{Predicted})^2$$

$$\frac{d}{d\text{intercept}} \text{ (Total SSR's)}$$

$$= d(\text{SSR}_1) + d(\text{SSR}_2) + d(\text{SSR}_3)$$

$$\frac{d}{d\text{intercept}} (\text{SSR}_1) = \frac{d}{d\text{intercept}} (1.4 + (-1) \text{ intercept} - 0.64 \times 0.5)$$

$$= -2 (1.4 - (\text{intercept} + 0.64 \times 0.5))$$

$$\frac{d}{d\text{intercept}} (\text{SSR}_2) = -2 (1.9 \cancel{-} (\text{intercept} + 0.64 \times 2.3))$$

$$\therefore (\text{SSR}_3) = -2 (3.2 - (\text{intercept} + 0.64 \times 2.9))$$

$$\text{Now } \frac{d}{d\text{intercept}} \text{ (Total SSR's)} = d(\text{SSR}_1) + d(\text{SSR}_2) + d(\text{SSR}_3)$$

now with this derivative value, A.D will use it to find where the SSR is the lowest.

Note: If we were using Least Squares to solve for the optimal value of the intercept, we would simply find where the slope of the curve = 0.

In contrast, AD finds the min. value by taking steps from an initial guess until it reaches the best value.

This makes AD very useful when it is not possible to solve for where the derivative = 0.

Our initial value of intercept = 0.

$$\therefore \underset{\text{d}_{\text{intercept}}}{d} (\text{SSR}^{\text{L}}) =$$

$$-2(1.4 - (0 + 0.64 * \frac{0.32}{1.4 + 2})) = -2.16$$

$$+ -2(1.9 - (0 + 0.64 * \frac{2.3}{1.856})) = -0.856$$

$$+ -2(3.2 - (0 + 0.64 * 2.9)) = -2.688$$

$$= -5.704$$

$$\Rightarrow \frac{d}{d\text{intercept}} (\text{SSR}'s) = -5.7$$

\Rightarrow When intercept = 0 ; slope = -5.7 of the curve \rightarrow

Note: The closer we set to the optimal value for the intercept, the closer the slope of the curve gets to zero.

\Rightarrow we need to take "Baby Steps" since we are close to the optimal value.

And we should take big steps when we are far from optimal value.

\Rightarrow Size of the step is related to the slope; as it tells whether we should take a baby step or big step

Note: We should ensure that the big step is not too big.

Gradient descent determines the step-size by multiplying the slope by a small no. called the "Learning Rate".

$$\begin{aligned}\text{Step-Size} &= -5.7 * 0.1 \\ &= -0.57\end{aligned}$$

With step-size we can calculate a new intercept

$$\begin{aligned}\text{New Intercept} &= \text{old Intercept} \\ &\quad - \text{Step-Size}\end{aligned}$$

17

Now the

$$\begin{aligned}
 \text{New-intercept} &= 0 - (-0.57) \\
 &= 0.57
 \end{aligned}$$

w.r.t. the original line & original intercept, we see that the residuals shrink when the intercept = 0.57.

To take another step, we go back to the derivative & plug in the new intercept ($= 0.57$)

$$\begin{aligned}
 \frac{d}{d\text{intercept}} (\text{SSRD}) &= \\
 &-2(1.4 - (0.57 + 0.64 \times 0.5)) \\
 &+ -2(1.9 - (0.57 + 0.64 \times 2.3)) \\
 &+ -2(3.2 - (0.57 + 0.64 \times 2.9)) \\
 &= -2.3 \quad (= \text{slope of the curve})
 \end{aligned}$$

Now calculate the step-size:

Step-size = slope \times learning rate

$$= -2.3 \times 0.1 \quad \cancel{-2.3}$$

Step-size = -0.23

$$\therefore \text{New-intercept} = 0.57 - (-2.3)$$

$$= 0.57 + \cancel{2.3} 0.23$$

New-intercept = 0.8

Now we can compare the residuals when the intercept = 0.57 & new-intercept = 0.8 ;

Overall the SSQ is getting smaller.

1st step-size = 0 to 0.57 (relative range)
2nd = 0.57 to 0.8

$$\frac{d}{\text{dintercept}} (\text{SSR}/n) = -2 \cdot 2.88 \approx -2.3 \quad (= -2.3)$$

Now calculate the derivative at
the new intercept:

$$\begin{aligned} \frac{d}{\text{dintercept}} &(\text{Sum of squared residuals}) \\ &= (-2)(1.4 - (0.8 + 0.64 \cdot 2.3)) \\ &\quad + (-2)(1.9 - (0.8 + 0.64 \cdot 2.3)) \\ &\quad + (-2)(3.2 - (0.8 + 0.64 \cdot 2.3)) \end{aligned}$$

$$= -0.9 \rightarrow \text{New-intercept}$$

$$\text{Step-size} = -0.9 * 0.1 = -0.09$$

$$\text{New-intercept} = 0.8 - (-0.09)$$

$$= 0.89$$

Now we increase the intercept from 0.8 to 0.89.

Then we take another step 4

$$\text{new-intercept} = 0.92$$

Another check New-interv oft = 0.94

100

$$= 0.95$$

Notice: Each step gets smaller & smaller the closer we get to the bottom of the curve.

After 6 steps, the AD estimate for the intercept is 0.95

Note: The least squares estimate for the intercept is also 0.95.

$$\frac{d}{d\text{intercept}} (\text{SSR} / n) = -2.288 \approx -2.3 \quad (= -2.3)$$

Now calculate the derivative at
the new intercept:

$$\begin{aligned} & \frac{d}{d\text{intercept}} (\text{Sum of Squared residuals}) \\ &= f(2)(1.4 - (0.8 + 0.64 * 2^0.5)) \\ &+ (-2)(1.9 - (0.8 + 0.64 * 2^0.3)) \\ &+ (-2)(3.2 - (0.8 + 0.64 * 2^0.9)) \end{aligned}$$

$$= -0.9 \rightarrow \text{New-intercept}$$

$$\text{Step-size} = -0.9 * 0.1 = -0.09$$

$$\begin{aligned} \text{New-intercept} &= 0.8 - (-0.09) \\ &= 0.89 \end{aligned}$$

Now we increase the intercept from
0.8 to 0.89.

Then we take another step &

$$\text{New-intercept} = 0.92$$

Another step: New-intercept = 0.94

$$= 0.95$$

n

:

a

Notice: Each step gets smaller &
smaller the closer we get
to the bottom of the curve.

After 6 steps, the AD estimate
for the intercept is 0.95

Note: The least squares estimate
for the intercept is also 0.95.

Now see anyone that GD has done its job.

But without comparing its role to a "Gold Std.", how does

GD knows to stop taking steps?

GD Termination Steps

- (1) when Step-size is very close to 0. (= when slope is very close to zero)
- (2) In practice, the minimum step-size = 0.001 (or) smaller

$$\text{Step-size} = \text{slope} * \text{Learning Rate}$$

If slope = 0.009;

then

$$\begin{aligned}\text{Step-size} &= 0.009 \times 0.1 \quad (\text{learning rate}) \\ &= 0.0009 \\ &(< 0.001)\end{aligned}$$

Now the AD stops.

(3) AD also includes a limit on the no. of steps it will take before giving up.

In practice, the user may

steps = 1000,

or greater

Even if step-size is large, if the user no. of steps > max. steps, AD will stop.

Summary:

(13)

(1) Use of SSR as Loss func),
to evaluate how well a line
fits the data

(2) we take the derivative of
SSR's i.e. the derivative
of the Loss Function.

$$\frac{d}{d\text{intercept}} (\text{SSR})$$

(3) Pick a RV for intercept
 $\Rightarrow 0$ in our case.

(4) we calculate the derivative
when $\text{intercept} = 0$. ($0/p = \text{slope of}$
 the curve)

(5) Plug the resultant slope value
into Step-size calculation.

Step-Size = Slope or learning Rate

(6) Calculate NewIntercept value
as:

$$\text{NewIntercept} = \frac{\text{old-intercept}}{-\text{Step-size}}$$

(7) Plug the new-intercept into
the derivative & repeat the steps
until the size goes close
to 0

$$\text{Step-size} = \text{slope} * \text{learning Rate}$$

$$\text{New-Intercept} = \frac{\text{old-intercept}}{-\text{Step-size.}}$$

We understood how CI can estimate the Intercept;

105

Let us understand how to "estimate the Intercept & Slope".

(1) Use SSR as before as the Loss Function.

(2) we have a 3-D graph:

X-axis = Slope-values

Y-axis = SSR

Z-axis = Intercept values.

(3) We need to find the values in the Intercept & Slope that gives us the min. SSR value.

$$\begin{aligned}
 (4) \quad SSR = & \\
 & (1.4 - (\text{Intercept} + \text{Slope } \times 0.5))^2 \\
 & + (1.9 - (\text{Intercept} + \text{Slope } \times 2.3))^2 \\
 & + (3.2 - (\text{Intercept} + \text{Slope } \times 2.9))^2
 \end{aligned}$$

Take the derivative of this func.

$\frac{d}{d\text{Intercept}} \text{SSR}$ w.r.t. Intercept

$$\frac{d}{d\text{Intercept}} \text{SSR} \underset{\text{w.r.t. Intercept}}{\underset{11}{\underset{15}{=}}}$$

+ we also take the derivative
of (err func) (SSR) w.r.t.

"Slope"

(5) Start taking the derivative w.r.t. the intercept. (Now slope is constant)

$$\frac{d}{d\text{intercept}} (\text{SSR}) =$$

$$(-2)(1.4 - (\text{intercept} + \text{slope} * 0.5)) \\ + (-2)(1.9 - (\text{intercept} + \text{slope} * 2.3)) \\ + (-2)(3.2 - (\text{intercept} + \text{slope} * 2.9))$$

Note:

Derivative of Intercept = 0

(6) Take derivative of Loss Fnc (=SSR)

w.r.t. the slope:

\rightarrow Intercept = Constant

$$\frac{d}{ds\text{slope}} (\text{SSR}) = -0.5 * 2(1.4 - (\text{intercept} + (\text{slope} * 0.5)))$$

$$+ (-2) * 2.3(1.9 - (\text{intercept} + (\text{slope} * 2.3)))$$

$$+ (-2) * (2 \cdot 9) \left(3 \cdot 2 - (\text{intercept} + \text{slope} \cdot 2 \cdot 9) \right)$$

Note: When u have 2 or more derivatives of the same fnc, they are called a "gradient".

We use this gradient to descend to lowest point in the Loss Fnc, i.e. the SSR:

Step 1: i) Pick a random no. fn the intercept. (=0)

ii) Pick a Random no. fn the slope. (Pick slope = 1)

Starting point:

19

$$\text{Intercept} = 0$$

$$\text{Slope} = 1.$$

Let the
Learning Rate
 $= 0.01$

Pleeg uit voor values

$$\Rightarrow \frac{d}{d\text{intercept}} (\text{SSR}) = -1.6 \quad \&$$

$$\frac{d}{dslope} (\text{SSR}) = -0.8$$

$$\text{Step-Size}_{\text{intercept}} = -1.6 + 0.01$$

$$\text{Step-Size}_{\text{slope}} = -0.8 + 0.01$$

Note: The layer learning rate
of the perceptrons e.g. ($= \cancel{0.01} 0.1$)
doesn't work this time.

Even after a bunch of steps,
AD doesn't arrive at the
correct answer.

\Rightarrow AD is very sensitive to
"learning Rate".

Note: In practice, a reasonable
learning rate can be determined
automatically by starting
large & getting smaller with
each step.

$$\begin{aligned}
 &= -2(1.4 - (0.57 \frac{0.89}{2.042} + 0.32)) \\
 &\quad + -2(1.9 - (0.57 \frac{2.426}{1.472} + 1.472)) \\
 &\quad + -2(3.2 - (0.57 \frac{2.426}{1.856} + 1.856))
 \end{aligned}$$

$$\begin{aligned}
 &= (-2)(0.51) + (-2)(-0.142) \\
 &\quad + (-2)(0.776)
 \end{aligned}$$

$$= -1.02 + 0.284 - 1.552$$

$$= -2.572 + 0.284$$

$$= -2.288 \approx 2.3$$