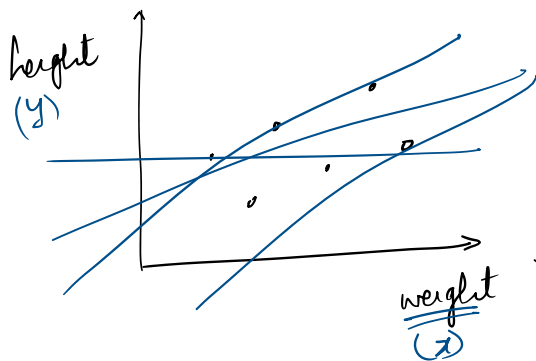


Gradient Descent →



parameters

$$\text{predicted height} = \text{intercept} + \text{slope} \times \text{weight}$$

$$y = c + mx$$

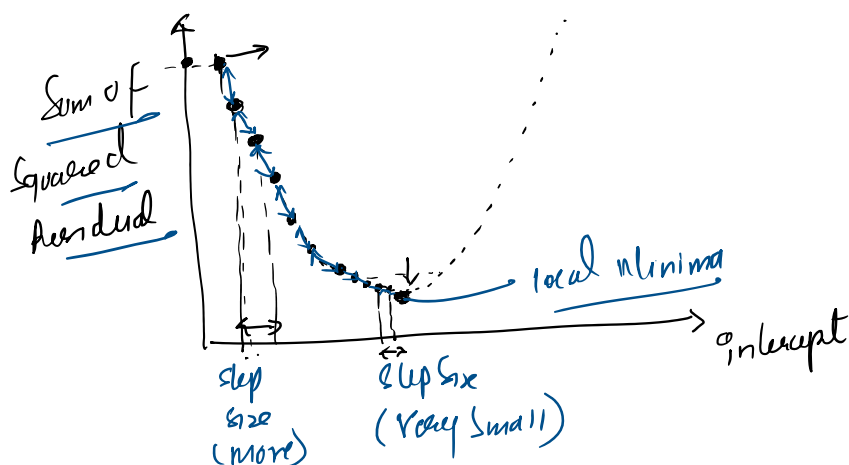
* Here we can have many line with different intercept and slope values

* We want a line with intercept and slope such that the Cost function should give minimum value.

* Let the cost function be Sum of Squared Error (Residual Error)

* Let us Consider only intercept

* Let's plot a graph of intercept v/s Sum of Squared Residual Error.



To find Minimum value for Sum of Squared Residual Error, we will start with very small value of intercept and will increase the value by very small amount.

Ex

	x	y
→	0.5	1.4
→	2.3	1.9
→	2.2	2.2

Let the Cost = Sum of Squared Error (Residual).
Fⁿ be

Ex

→ 0.5	1.4
→ 2.3	1.9
→ 2.9	3.2
↑	↑ Actual

Let the cost = sum of squared residuals.
Fⁿ be

$$\text{Sum of Squared Residual Error} = \left(1.4 - (\text{intercept} + \text{slope} \times 0.5) \right)^2 + \left(1.9 - (\text{intercept} + \text{slope} \times 2.3) \right)^2 + \left(3.2 - (\text{intercept} + \text{slope} \times 2.9) \right)^2$$

Here we want to find value of intercept and slope that give minimum sum of squared error.

Step 1 → Find $\frac{d}{d \text{intercept}}$ (Sum of Squared Residual Error)
 $\frac{d}{d \text{slope}}$ (Sum of Squared Residual Error).

$$\text{Sum of Squared Residual Error} = \left(1.4 - (\text{intercept} + \text{slope} \times 0.5) \right)^2 + \left(1.9 - (\text{intercept} + \text{slope} \times 2.3) \right)^2 + \left(3.2 - (\text{intercept} + \text{slope} \times 2.9) \right)^2$$

(SSE)

Using chain Rule $\left[\frac{d}{dx} x^2 = 2x \frac{d}{dx} x \right]$

$$\begin{aligned} \frac{d}{d \text{intercept}} (\text{Sum of Squared Error}) &= 2 \times \left(1.4 - (\text{intercept} + \text{slope} \times 0.5) \right) \times \frac{d}{d \text{intercept}} \left(1.4 - (\text{intercept} + \text{slope} \times 0.5) \right) \\ &\quad + 2 \times \left(1.9 - (\text{intercept} + \text{slope} \times 2.3) \right) \times \frac{d}{d \text{intercept}} \left(1.9 - (\text{intercept} + \text{slope} \times 2.3) \right) \\ &\quad + 2 \times \left(3.2 - (\text{intercept} + \text{slope} \times 2.9) \right) \times \frac{d}{d \text{intercept}} \left(3.2 - (\text{intercept} + \text{slope} \times 2.9) \right) \end{aligned}$$

(-1)

(-1)

(-1)

$$\frac{d}{d \text{intercept}} (3.2 - (\text{intercept} + \text{slope} \times 2.9)) \quad (-1)$$

$$\frac{d(SS E)}{d \text{intercept}} = (-2)(1.4 - (\text{intercept} + \text{slope} \times 0.5)) + (-2)(1.9 - (\text{intercept} + \text{slope} \times 2.3)) + (-2)(3.2 - (\text{intercept} + \text{slope} \times 2.9))$$

Eq 1.

$$\text{Sum of Squared Residual Error} = \left(1.4 - (\text{intercept} + \text{slope} \times 0.5)\right)^2 + \left(1.9 - (\text{intercept} + \text{slope} \times 2.3)\right)^2 + \left(3.2 - (\text{intercept} + \text{slope} \times 2.9)\right)^2$$

(SS E)

Now

$$\frac{d(SS E)}{d \text{slope}} = \frac{d}{d \text{slope}} \left(1.4 - (\text{intercept} + \text{slope} \times 0.5)\right)^2 + \frac{d}{d \text{slope}} \left(1.9 - (\text{intercept} + \text{slope} \times 2.3)\right)^2 + \frac{d}{d \text{slope}} \left(3.2 - (\text{intercept} + \text{slope} \times 2.9)\right)^2$$

$$\frac{d(SS E)}{d \text{slope}} = \frac{(-2)(0.5)(1.4 - (\text{intercept} + \text{slope} \times 0.5))}{+(-2)(2.3)(1.9 - (\text{intercept} + \text{slope} \times 2.3)) + (-2)(2.9)(3.2 - (\text{intercept} + \text{slope} \times 2.9))}$$

Eq 2.

We got the above 2 equations.

Let start with random values of intercept & slope

$$\text{Let } \begin{cases} \text{intercept} = 0 \\ \text{slope} = 1 \end{cases}$$

Let intercept = 0 & slope = 1 in Eq 1 & Eq 2

Substitute $\phi_{\text{intercept}} = 0$ & $\text{slope} = 1$ in Eq 1 & Eq 2

$$\frac{d(\text{SSE})}{d\text{slope}} = -0.8 \Rightarrow \text{slope of line wot parameter } \underline{\text{slope}}$$

$$\frac{d(\text{SSE})}{d\text{intercept}} = -1.6 \Rightarrow \text{slope of line wot parameter } \underline{\phi_{\text{intercept}}}$$

For step Size

$$\text{New Step Size} = \text{slope of line at point} * \text{Learning Rate}$$

* Learning Rate \Rightarrow takes smaller value [10^{-6} to 1.0]

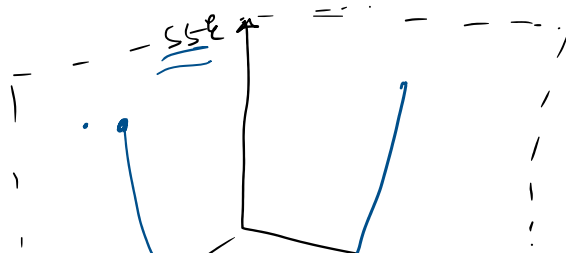
$$\text{Let us take Learning Rate} = 0.01$$

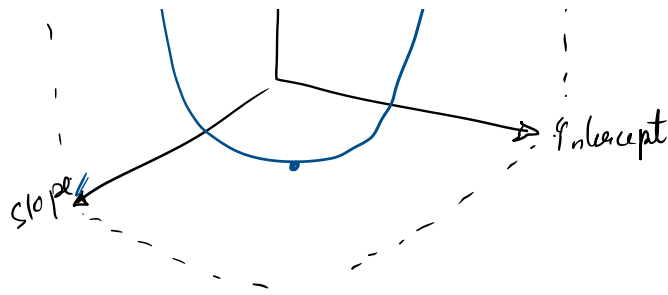
$$\text{stepSize}_{(\phi_{\text{intercept}})} = \frac{d(\text{SSE})}{d\text{intercept}} * 0.01 = -1.6 * 0.01 = \underline{\underline{-0.016}}$$

$$\text{stepSize}_{(\text{slope})} = \frac{d(\text{SSE})}{d\text{slope}} * 0.01 = -0.8 * 0.01 = -0.008$$

$$\text{New Intercept} = \text{old } \phi_{\text{intercept}} - \frac{\text{stepSize}}{(\text{intercept})} = 0 - (-0.016) = \underline{\underline{0.016}}$$

$$\text{New Slope} = \text{old slope} - \frac{\text{stepSize}}{(\text{slope})} = 1 - (-0.008) = \underline{\underline{1.008}}$$





Substitute the new intercept and slope in $\frac{d(SSC)}{d \text{intercept}}$ & $\frac{d(SSC)}{d \text{slope}}$
 and Calculate new step size for intercept & slope
 And again Calculate New Intercept & New slope

* We will repeat until the step size becomes very small
 $\Rightarrow \underline{\underline{0.001}}$

or Some maximum Number of steps is reached Ex: 1000

For above Dataset Best fitting line will have
 $\text{intercept} = 0.95$ & $\text{slope} = 0.64$

This is How Gradient Descent Optimizes Parameters

Summary \rightarrow Gradient Descent

- ① get Eqⁿ of line [Identify Simple linear or Polynomial]
- ② Want intercept & slope value that minimizes cost fⁿ
- ③ Identify cost (loss) function. [we took SSE]
- ④ Find Eq for $\frac{d(SSC)}{d \text{intercept}}$ & $\frac{d(SSC)}{d \text{slope}}$
- \rightarrow ⑤ Start with some initial value of slope & intercept.
- ⑥ Find the value of $\frac{d(SSC)}{d \text{intercept}}$ & $\frac{d(SSC)}{d \text{slope}}$

⑥ Find the value of $\frac{d(SSE)}{d \text{intercept}}$ & $\frac{d(SSE)}{d \text{slope}}$

⑦ Calculate $\text{StepSize}_{(\text{intercept})} = \frac{d(SSE)}{d \text{intercept}} * \text{learning rate}$

$\text{StepSize}_{(\text{slope})} = \frac{d(SSE)}{d \text{slope}} * \text{learning rate}$

⑧ Update intercept = old intercept - $\text{StepSize}_{(\text{intercept})}$

Update slope = old slope - $\text{StepSize}_{(\text{slope})}$

⑨ Repeat step 6 - 8 until stepSize is very small
or Number of iterations ≈ 1000

Gradient Boost for Regression

* Gradient Boost start with single leaf instead of Tree/stump.

* leaf represents initial guess for the weights of the samples.

Consider

Height	Favourite color	Gender	<u>Weight</u> ^(Actual)
1.6	Blue	male	88
1.6	Green	Female	76
1.5	Blue	female	56
1.8	Red	male	73
1.5	Green	male	77
1.4	Blue	Female	57

Step 1 \Rightarrow First let make Initial Guess = weight (Average) = 71.2

Step 2 \Rightarrow Like Adaboost, even Gradient Boost builds tree based on previous Tree

But unlike Adaboost, the Gradient Boost tree are larger than stump

[However the height is still Restricted [No of leaf 8 \rightarrow 32] in Each Tree]

* Gradient Boost will build Another Tree based on Error of previous Tree.

* [Gradient Boost continuous to build Tree in this fashion untill it has made number of Trees you asked for or the Additional tree fails to improve the fit]

⇒ step 3 : Initial Guess = weight (Average) = 71.2

$$\text{Predicted Value} = \boxed{71.2} + \text{Tree 1} + \text{Tree 2} \dots$$

Initial guess

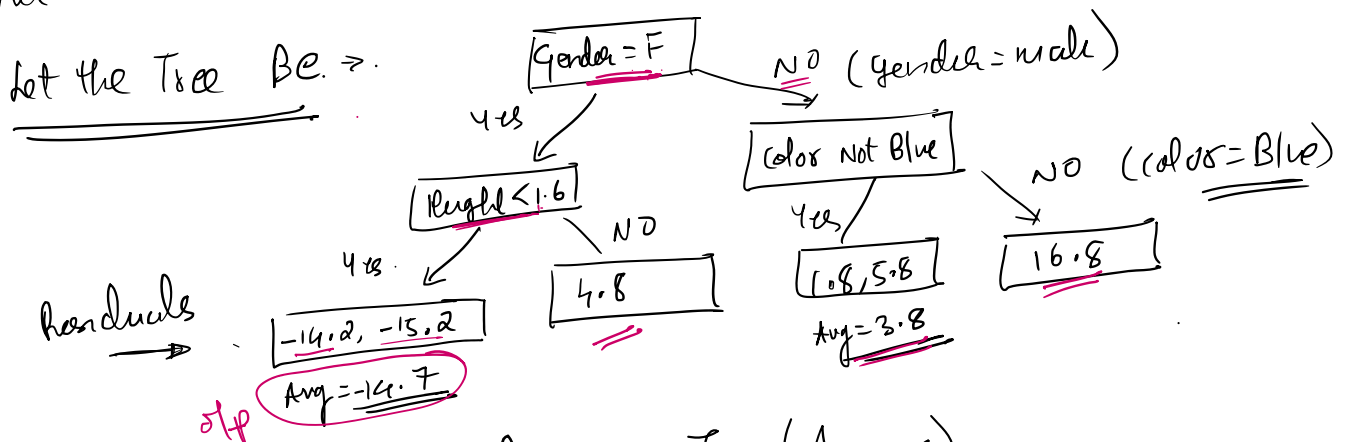
Boosting

Initial 71.2

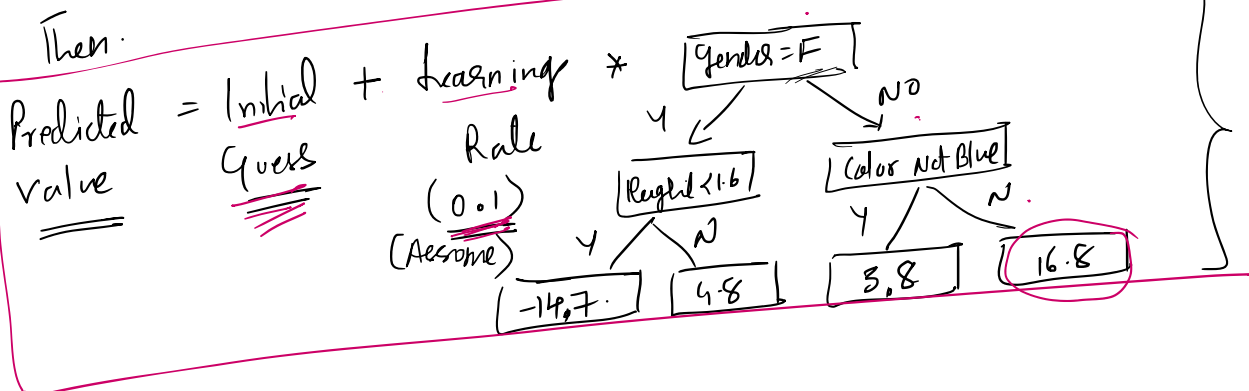
(Observed - Average)

Height	Favorite color	Gender	Weight	<u>Residual</u>
1.6	Blue	male	88	16.8
1.6	Green	Female	76	-4.8
1.5	Blue	Female	58	-15.2
1.8	Red	male	73	1.8
1.5	Green	male	77	5.8
1.4	Blue	Female	57	-14.2

Let Build a Tree using Residual and not weight



if there will be only one Tree (Assume).



Now first Sample =

1.6	Blue	male	?
-----	------	------	---

$$0 \text{ and } 1 \text{ and } 1 - 71.2 + 0.1 * 16.8 = \boxed{72.88}$$

$$\text{Predicted value} = 71.2 + 0.1 * 16.8 = \underline{\underline{72.88}}$$

* Repeat this for every sample & calculate Predicted value for each sample
→ call it as New Prediction.

Height	Favorite color	Gender	Weight	<u>Residual</u>	<u>New Prediction</u>	<u>New Residual</u>
1.6	Blue	male	88	16.8 ✓	72.88 ✓	R ₁ ✓
1.6	Green	Female	76	-4.8	P ₁ ✓	R ₂ ✓
1.5	Blue	female	58	-15.2	P ₂ ✓	R ₃ ✓
1.8	Red	male	73	1.8	P ₃ ✓	R ₄ ✓
1.5	Green	male	77	5.8	P ₄ ✓	R ₅ ✓
1.4	Blue	Female	57	-14.2	P ₅ ✓	R ₆ ✓

* Now find New Residual for New Prediction for all the samples

* Using this new Residual Construct New Tree and so on
Unill the New Tree does not make significant Contribution
to the new Prediction.

* let us say we have 3 Tree for above Dataset →

$$\text{Final Prediction} = \text{Initial value} + 0.1 * \boxed{\text{Tree 1}} + 0.1 * \boxed{\text{Tree 2}} + 0.1 * \boxed{\text{Tree 3}}$$

[Gradient Boost Principle → Taking lot of small steps in
Right Direction will result into
better Predictions]