

# Gradient-Descent.

## applied to LINEAR Regression

A quick Refresher!

Recall Simple Linear Regression.

$$\hat{y} = mx + b$$

predicted

Slope  $\downarrow$  Constant  $\downarrow$   
 $y = mx + b + \epsilon$   
 $\downarrow$   
error or Residual

$$\text{Error} = y_{\text{actual}} - \hat{y}_{\text{predicted}}$$

For 'n' observations.

The loss term is <sup>mean</sup> Sum of Squares of errors.

$$\text{Thus } L(m, b) = \frac{1}{N} \sum_{i=1}^N (y_i - (mx_i + b))^2$$

Take gradients

$$\frac{\partial L}{\partial m} = \frac{2}{N} \sum_{i=1}^N -x_i (y_i - (mx_i + b)) \quad \text{'m' gradient}$$

$$\frac{\partial L}{\partial b} = \frac{2}{N} \sum_{i=1}^N -(y_i - (mx_i + b)) \quad \text{'b' gradient}$$

Update Rule

$$\text{New-}b = \text{Current-}b - \alpha \times \text{b-gradient}_{\text{Current}}$$

$$\text{New-}m = \text{Current-}m - \alpha \times \text{m gradient}_{\text{Current}}$$

- Compute New line. with  $\text{New-}b$   $\text{New-}m$
- Compute Loss → get gradients
- Iterate. till convergence.

# Gradient Boosting For Regression.

GBM

→ Given  $(x_1, y_1) \dots (x_n, y_n)$   
Fit  $F(x)$ .

[Initial Model  $\rightarrow F(x)$ ]

Say we have the following.

$$\begin{aligned}h(x_1) &= y_1 - F(x_1) \\h(x_2) &= y_2 - F(x_2) \\&\vdots \\h(x_n) &= y_n - F(x_n)\end{aligned}$$

$F(x)$  can be any Model!

Fit a Regression tree to the data.

$$\begin{aligned}&(x_1, y_1 - F(x_1)), \\&(x_2, y_2 - F(x_2)) \\&\vdots \\&(x_n, y_n - F(x_n))\end{aligned}$$

[Boosted Model.  $F(x) + h(x)$ ]

Gradient Descent.

$$\theta_i := \theta_i - \rho \frac{\partial J}{\partial \theta_i}$$

How is this Related to Gradient Descent?

Consider the loss function. for the fitted Model  $F(x)$

$$L(y, F(x)) = \frac{1}{2} [y - F(x)]^2$$

We want to minimize

$$J = \sum_i L(y_i, F(x_i))$$

Consider  $F(x_1) \dots F(x_n)$  as parameters

$$\frac{\partial J}{\partial F(x_i)} = \frac{\partial \sum_i L(y_i, F(x_i))}{\partial F(x_i)} = \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}$$

$$\text{or } \frac{\partial J}{\partial F(x_i)} = F(x_i) - y_i \quad \text{--- (A)}$$

Re arranging (A)

$$\underbrace{y_i - F(x_i)}_{\text{residual}} = \underbrace{-\frac{\partial J}{\partial F(x_i)}}_{\text{-ve gradient}}$$

$\therefore$  We can interpret residuals as negative gradients.  
Thus we have the following:

$$\begin{array}{l} F(x_i)_{\text{New}} := F(x_i)_{\text{Current}} + \underbrace{h(x_i)}_{\text{residuals}} \end{array} \quad \left. \vphantom{\begin{array}{l} F(x_i)_{\text{New}} := F(x_i)_{\text{Current}} + h(x_i) \end{array}} \right\} \text{Concept of Boosting}$$

$$\text{or } F(x_i)_{\text{new}} := F(x_i)_{\text{Current}} + \underbrace{y_i - F(x_i)_{\text{Current}}}_{\text{residual Current}}$$

$$\text{or } F(x_i)_{\text{new}} := F(x_i) - \frac{\partial J}{\partial F(x_i)}$$

Compare with Gradient Descent

$$\rightarrow \theta_i^{\text{New}} = \theta_i^{\text{Current}} - \rho \frac{\partial J}{\partial \theta_i}$$



# [CASE] Regression with Square loss.

[6BM-3]

Negative gradient:

$$-g(x_i) = \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} = y_i - F(x_i)$$

## Algorithm

- Start with an initial model, say  $F(x) = \sum_{i=1}^n \frac{y_i}{n}$
- Iterate till convergence
  - Calculate -ve gradients  $-g(x_i)$
  - Fit a regression tree  $h$  to -ve grad.
  - $F_{\text{new}} = F_{\text{old}} + \rho h$
  - ↑
  - $\rho = 1$

A Fundamental Question arises:

Why -ve gradient over residuals?

Consider the types of Loss Functions shown below.

①  $L = \frac{(y_i - F(x_i))^2}{2}$  ← Square Loss Function.

② Absolute loss.

$L(y, F) = |y - F|$  ← Outliers are heavily punished. More Robust to outliers.

③ Huber loss! (More robust to outliers)

$$L(y, F) = \begin{cases} \frac{1}{2} (y - F)^2 & |y - F| \leq \delta \\ \delta (|y - F| - \frac{\delta}{2}) & |y - F| > \delta \end{cases}$$

## Negative Gradient vs Residual.

[GBM-4]

Example: using Huber loss.

Huber loss:

$$L(y, F) = \begin{cases} \frac{1}{2} (y - F)^2 & |y - F| \leq \delta \\ \delta (|y - F| - \delta/2) & |y - F| > \delta \end{cases}$$

CASE A

→ Update by residual

$$h(x_i) = y_i - F(x_i)$$

CASE B

→ Update by Negative Gradient.

$$h(x_i) = -g(x_i) = \begin{cases} y_i - F(x_i) & |y_i - F(x_i)| \leq \delta \\ \delta \operatorname{sign}(y_i - F(x_i)) & |y_i - F(x_i)| > \delta \end{cases}$$

Negative gradient  $\rightarrow$  outliers get less attention.

Q] Why should we moderate the effect of

outliers?

→ We do not want a 'boost' in the wrong direction !!!



# Gradient Boosting for Classification

[6BM.]  
5

Considers the use case:

Handwritten Capital letter classification

•) 26 classes (A to Z)

•) 16 features]



an 'ON' Pixel

1] horizontal pos of box

2] vertical pos " "

3] width of box

4] height of box

5] Total on Pixels.

6] Mean x of 'ON' Pixels

7] Mean y of 'ON' Pixels

8] Mean x variance.

9] Mean y Vari

10] Mean xy correl.

11] Mean of  $x^2y$

12] Mean of  $xy^2$

13] Mean edge count

L to R.

14] Correl x edge

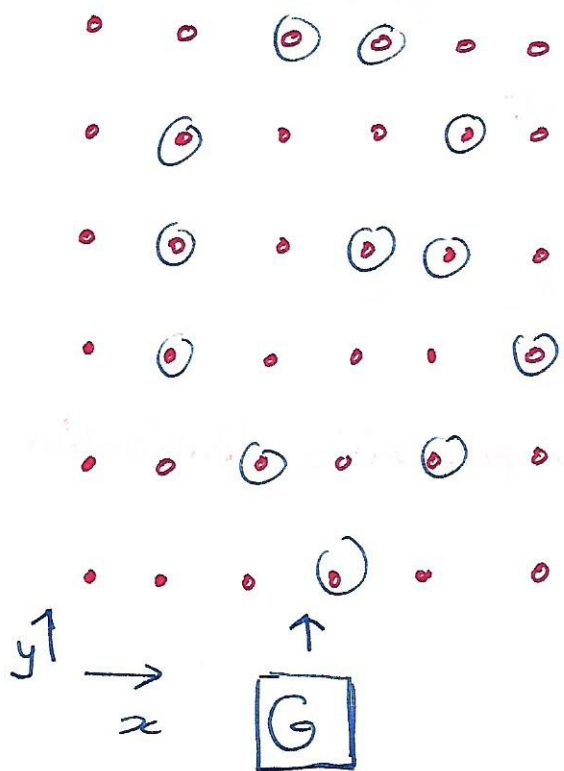
with y

15] Mean edge count

bottom to top

16] Correl. of

'y' edge with



← example of a  
36 (6x6)  
Pixel box.

⊙ denotes an 'ON' Pixel.

Modelling

⊙ Scoring

] 26 score functions, one for each class.

$F_A(x)$  assigns a score for class A.

Probabilities:

$$P_A(x) = \frac{e^{F_A(x)}}{\sum_{c=A}^Z e^{F_c(x)}}$$

$$P_Z(x) = \frac{e^{F_Z(x)}}{\sum_{c=A}^Z e^{F_c(x)}}$$

Consider  $i$ th Row  $X_i [x_1, x_2 \dots x_{16}]$   $\bar{y}_i =$

16 features per observation

For each Row: Eval.  $P_A(X_i)$  to  $P_B(X_i)$

Predicted label = class that has highest probability.

① Loss function at each data point

Say For  $X_5 [x_{1,5}, x_{2,5} \dots x_{16,5}]$   $y_5 = 'G'$

Then we can say.

$$y_A(X_5) = 0$$

$$y_B(X_5) = 0$$

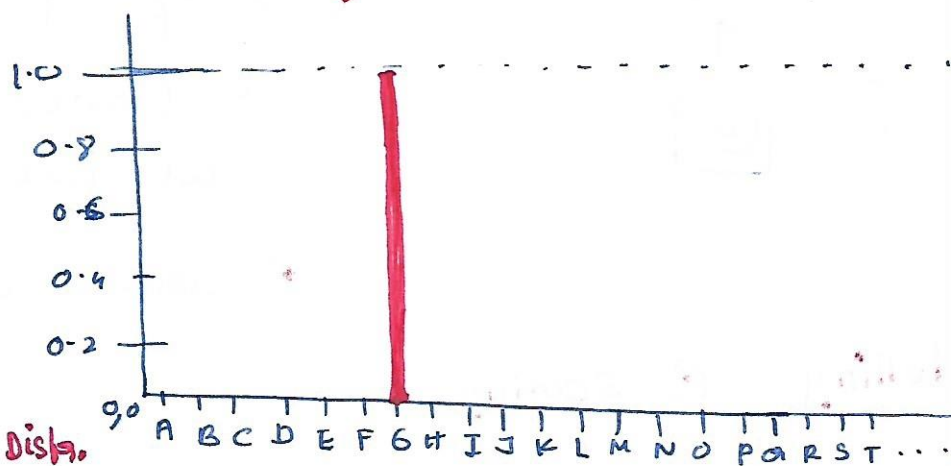
$\vdots$

$$y_G(X_5) = 1$$

$\vdots$

$$y_Z(X_5) = 0$$

Convert this into a true probability distribution



Turn the Label

into a 'TRUE' Probability Dist.

Assume: Predicted Prob. at current Model Iteration is.

$$P_A(X_5) = 0.03$$

$$P_B(X_5) = 0.05$$

$\vdots$

$$P_G(X_5) = 0.3$$

$\vdots$

$$P_Z(X_5) = 0.05$$

Minimize total Loss

what is KL-divergence

How do we go ahead

in this case?

What is our initial guess for  $F_{GS}$ ?

Regression  
with Loss Function.General Procedure .

•) Given any differential Loss function  $L$

Start]  $\rightarrow$  Start with Initial model  $F$

Iterate until Converge!

•) Calculate -ve gradients  $-g(x_i) = -\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}$

Size of the tree

2 Terminal Nodes  
(Stump only)

•) Fit Regression tree ' $h$ ' to -ve gradients  $-g(x_i)$

•) update  $F := F + \rho h$ .

Procedure For alphabet classification case.

$\rightarrow$  Start with initial Models For  
 $F_A = 0, F_B = 0, \dots, F_Z = 0$ .

Iterate till convergence

•)  $-g_A(x_i) = -\frac{\partial L}{\partial F_A(x_i)}$

$-g_B(x_i) = -\frac{\partial L}{\partial F_B(x_i)}$

$-g_Z \dots$

•) Fit Regression tree  $h_A$  to -ve grad  $-g_A(x_i)$

•) update  $h_Z$ .

$F_A := F_A + \rho_A h_A$

$F_Z := F_A + \rho_Z h_Z$



K.L. Divergence  $\rightarrow$  Loss Function.

$$D_{KL}(p||q) = \sum_{i=1}^N p(x_i) \log \left( \frac{p(x_i)}{q(x_i)} \right)$$

$q(x)$  is approximation.

$p(x)$  is true distribution.

Intuitively. This measures: How much a given arbitrary distribution is away from the true distribution.