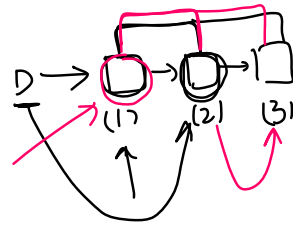


Gradient Boosting

04 August 2023 20:07

Ensemble technique

↳ Boosting



LB LV

stagenwisc

Additive modelling

$$B \propto \frac{1}{V}$$

weak learner
+ HB / LV

$$HB + HB + HB$$

(LB)

Regression.
↓
Classification

Gradient Boosting

mse / log loss

Input: training set $\{(x_i, y_i)\}_{i=1}^n$, a differentiable loss function $L(y, F(x))$, number of iterations M .

1. Initialize $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$.

2. For $m = 1$ to M :

(a) For $i = 1, 2, \dots, N$ compute

$$r_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}.$$

(b) Fit a regression tree to the targets r_{im} giving terminal regions R_{jm} , $j = 1, 2, \dots, J_m$.

(c) For $j = 1, 2, \dots, J_m$ compute

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma).$$

(d) Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$.

3. Output $\hat{f}(x) = f_M(x)$.

8 rows 3 cols

cgpa	iq	is_placed
6.82	118	0
6.36	125	1
5.39	99	1
5.50	106	1
6.39	148	0
9.13	148	1
7.17	147	1
7.72	72	0

mean $\frac{5}{3}$ $\log(\frac{5}{3}) = f_0(x)$ $\rightarrow 1^{st}$ model

3 stages pseudoresidual $y_i - f_0(x_i)$

log odds $\log \rightarrow \text{mean}$ $\text{class} \rightarrow \log \text{ odds}$

$F(x) = f_0(x) + f_1(x) + f_2(x)$

$p = \frac{1}{1 + e^{-\log \text{ odds}}} = \frac{1}{1 + e^{-0.51}} = 0.62$

cgpa	iq	is_placed	prel(log-odds)	prel(probability)
6.82	118	0	0.510826	0.625
6.36	125	1	0.510826	0.625
5.39	99	1	0.510826	0.625
5.50	106	1	0.510826	0.625
6.39	148	0	0.510826	0.625
9.13	148	1	0.510826	0.625
7.17	147	1	0.510826	0.625
7.72	72	0	0.510826	0.625

$\frac{e^x}{1+e^x}$

good stopping point

cgpa	iq	is_placed	prel(log-odds)	prel(probability)	res1
6.82	118	0	0.510826	0.625	-0.625
6.36	125	1	0.510826	0.625	0.375
5.39	99	1	0.510826	0.625	0.375
5.50	106	1	0.510826	0.625	0.375
6.39	148	0	0.510826	0.625	-0.625
9.13	148	1	0.510826	0.625	0.375
7.17	147	1	0.510826	0.625	0.375
7.72	72	0	0.510826	0.625	-0.625

regression tree weak learner \rightarrow leaf nodes = 3

node #0: cgpa <= 6.375, squared_error = 0.234, samples = 8, value = 0.0

node #1: squared_error = 0.0, samples = 3, value = 0.375

node #2: iq <= 132.5, squared_error = 0.24, samples = 5, value = -0.225

node #3: squared_error = 0.0, samples = 2, value = -0.625

node #4: squared_error = 0.222, samples = 3, value = 0.042

diff. prob

$\frac{\sum \text{Residual}}{\sum [\text{Previous Prob} * (1 - \text{Previous Prob})]}$

$= \frac{-0.625 - 0.625}{0.625(1 - 0.625) + 0.625(1 - 0.625)} = \frac{-2 \times 0.625}{2 \times 0.625 \times 0.375} = -2.66$

output \rightarrow log odds

$0.51 + (-2.66) = -2.15$

$p = \frac{1}{1 + e^{-\log \text{ odds}}} = \frac{1}{1 + e^{-2.159}} = 0.10$

cgpa	iq	is_placed	prel(log-odds)	prel(probability)	res1	leaf_entry1	pre2(log-odds)	pre2(probability)
6.82	118	0	0.510826	0.625	-0.625	3	-2.159174	0.103477
6.36	125	1	0.510826	0.625	0.375	1	2.110826	0.891951
5.39	99	1	0.510826	0.625	0.375	1	2.110826	0.891951
5.50	106	1	0.510826	0.625	0.375	1	2.110826	0.891951
6.39	148	0	0.510826	0.625	-0.625	4	0.690826	0.666151
9.13	148	1	0.510826	0.625	0.375	4	0.690826	0.666151
7.17	147	1	0.510826	0.625	0.375	4	0.690826	0.666151
7.72	72	0	0.510826	0.625	-0.625	3	-2.159174	0.103477

node #0: cgpa <= 6.375, squared_error = 0.234, samples = 8, value = 0.0

node #1: squared_error = 0.0, samples = 3, value = 0.375

node #2: iq <= 132.5, squared_error = 0.24, samples = 5, value = -0.225

node #3: squared_error = 0.0, samples = 2, value = -0.625

node #4: squared_error = 0.222, samples = 3, value = 0.042

node #0: cgpa <= 6.995, squared_error = 0.09, samples = 8, value = 0.015

node #1: iq <= 136.5, squared_error = 0.09, samples = 5, value = -0.089

node #2: squared_error = 0.043, samples = 3, value = 0.188

node #3: squared_error = 0.008

node #4: squared_error = 0.0

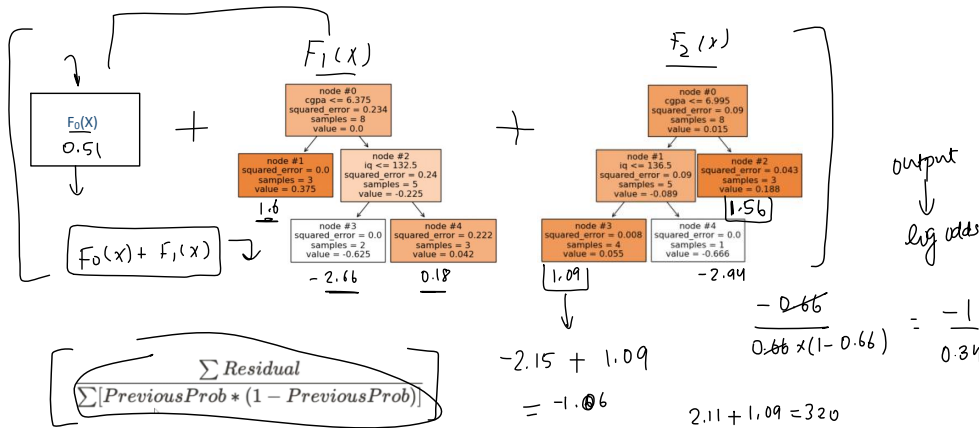
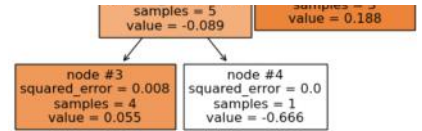
res2

$y - [f_0(x_i) + f_1(x_i)]$

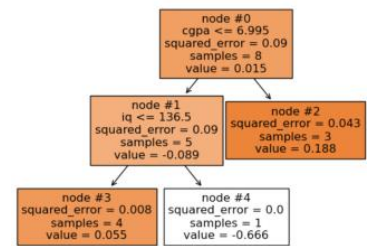
$\log \text{ odds} \rightarrow \text{prob}$

6.39	148	0	0.510826	0.625	-0.625	4	0.690826	0.666151
9.13	148	1	0.510826	0.625	0.375	4	0.690826	0.666151
7.17	147	1	0.510826	0.625	0.375	4	0.690826	0.666151
7.72	72	0	0.510826	0.625	-0.625	3	-2.159174	0.103477

-0.666151
0.333849
0.333849
-0.103477



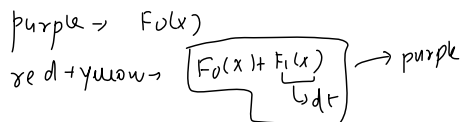
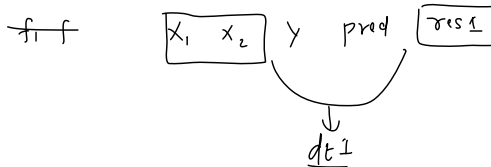
cgpa	iq	is_placed	pre1(log-odds)	pre1(probability)	res1	leaf_entry1	pre2(log-odds)	pre2(probability)	res2	leaf_entry2
6.82	118	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477	3
6.36	125	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3
5.39	99	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3
5.50	106	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3
6.39	148	0	0.510826	0.625	-0.625	4	0.690826	0.666151	-0.666151	4
9.13	148	1	0.510826	0.625	0.375	4	0.690826	0.666151	0.333849	2
7.17	147	1	0.510826	0.625	0.375	4	0.690826	0.666151	0.333849	2
7.72	72	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477	2



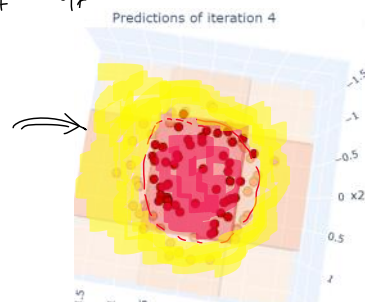
$$\frac{1}{1+e^{1.06}}$$

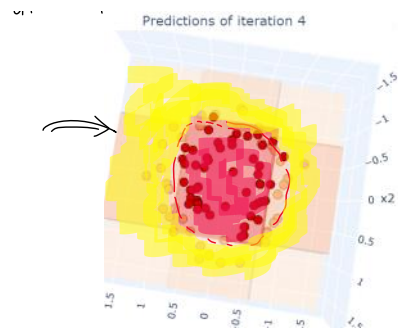
$$\text{out} \rightarrow p = \frac{1}{1+e^{-\log \text{odds}}}$$

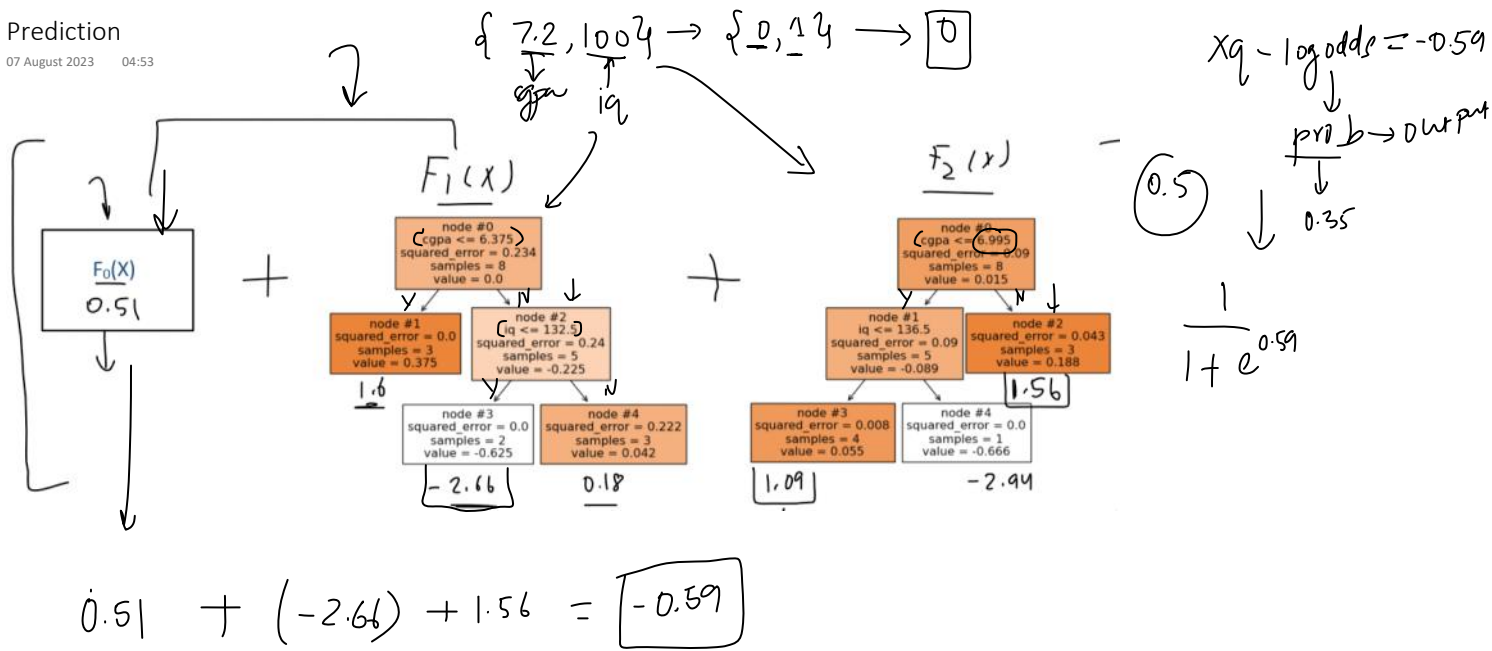
cgpa	iq	is_placed	pre1(log-odds)	pre1(probability)	res1	leaf_entry1	pre2(log-odds)	pre2(probability)	res2	leaf_entry2	pre3(log-odds)	pre3(probability)
6.82	118	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477	3	-1.068349	0.255717
6.36	125	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3	3.201651	0.960896
5.39	99	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3	3.201651	0.960896
5.50	106	1	0.510826	0.625	0.375	1	2.110826	0.891951	0.108049	3	3.201651	0.960896
6.39	148	0	0.510826	0.625	-0.625	4	0.690826	0.666151	-0.666151	4	-1.798349	0.142052
9.13	148	1	0.510826	0.625	0.375	4	0.690826	0.666151	0.333849	2	2.251651	0.904793
7.17	147	1	0.510826	0.625	0.375	4	0.690826	0.666151	0.333849	2	2.251651	0.904793
7.72	72	0	0.510826	0.625	-0.625	3	-2.159174	0.103477	-0.103477	2	-0.598349	0.354722



$$F_0(x) + F_1(x) + F_2(x)$$

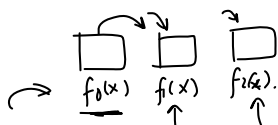






Geometric Intuition

04 August 2023 20:15



Input: training set $\{(x_i, y_i)\}_{i=1}^n$, a differentiable loss function $L(y, F(x))$, number of iterations M .

→ 1. Initialize $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$.

→ 2. For $m = 1$ to M :

→ (a) For $i = 1, 2, \dots, N$ compute

$$r_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$$

$\frac{\partial L}{\partial f_0(x_i)}$

→ (b) Fit a regression tree to the targets r_{im} giving terminal regions R_{jm} , $j = 1, 2, \dots, J_m$.

→ (c) For $j = 1, 2, \dots, J_m$ compute

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma).$$

$$\frac{\sum \text{Residual}}{\sum [\text{PreviousProb} * (1 - \text{PreviousProb})]}$$

(d) Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$.

3. Output $\hat{f}(x) = f_M(x)$.

Step 0 -> The Loss Function

10 August 2023 09:43

$$\text{Log loss} \rightarrow \mathcal{L} = -\frac{1}{n} \sum_{i=1}^n y_i \log p_i + (1-y_i) \log(1-p_i) \rightarrow \mathcal{L}(p_i)$$

$\hat{y}_i \rightarrow$ output prob

cgpa	iq	placement	$\hat{y}(p_i)$
8	80	1	0.62
7	70	0	0.32
6	60	1	0.51

log odds \rightarrow prob

$\mathcal{L}(\log \text{odds})$

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n y_i \log p_i + (1-y_i) \log(1-p_i)$$

log odds

$$\log\left(\frac{p_i}{1-p_i}\right)$$

$$\mathcal{L} = -\frac{1}{n} \left[\sum_{i=1}^n y_i \log p_i + \log(1-p_i) - y_i \log(1-p_i) \right]$$

$$= -\frac{1}{n} \left[\sum_{i=1}^n y_i [\log p_i - \log(1-p_i)] + \log(1-p_i) \right]$$

$$= -\frac{1}{n} \left[\sum_{i=1}^n y_i \log(\text{odds}_i) + \log(1-p_i) \right]$$

$$p_i = \frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}}$$

$$= -\frac{1}{n} \left[\sum_{i=1}^n y_i \log(\text{odds}_i) + \log\left(1 - \frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}}\right) \right]$$

$$= -\frac{1}{n} \left[\sum_{i=1}^n y_i \log(\text{odds}_i) + \log\left(\frac{1}{1 + e^{\log(\text{odds}_i)}}\right) \right]$$

$$= -\frac{1}{n} \left[\sum_{i=1}^n y_i \log(\text{odds}_i) + \log 1 - \log(1 + e^{\log(\text{odds}_i)}) \right]$$

$$\mathcal{L} = -\frac{1}{n} \left[\sum_{i=1}^n y_i \log(\text{odds}_i) - \log(1 + e^{\log(\text{odds}_i)}) \right] \leftarrow \text{data}$$

$$L = -\frac{1}{n} \left[\sum_{i=1}^n x_i \log(\text{odds}_i) - y \dots \right]$$

$$L = -y \log(\text{odds}) + \log(1 + e^{\log(\text{odds})})$$

Step 1 -> Finding $F_0(x)$

10 August 2023 09:44

→ Loss function

1. Initialize $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$.

$$\rightarrow L = -\frac{1}{n} \left[\sum_{i=1}^n y_i \log(\text{odds}_i) - \log(1 + e^{\log(\text{odds}_i)}) \right]$$

$$\downarrow$$

$$L(y_i, \log(\text{odds}_i))$$

$$\downarrow$$

$$L(y_i, \gamma)$$

$$\arg \min_{\gamma} L = -\frac{1}{n} \left[\sum_{i=1}^n y_i \gamma - \log(1 + e^{\gamma}) \right]$$

cg pa | ia | place

1
0
1
0

$$\frac{\partial L}{\partial \gamma} = -\frac{1}{n} \left[\sum_{i=1}^n y_i - \frac{e^{\gamma}}{1 + e^{\gamma}} \right] = 0$$

$$\frac{1+1+0+0}{4}$$

$$= -\frac{1}{n} \sum_{i=1}^n y_i + \frac{1}{n} \sum_{i=1}^n \frac{e^{\gamma}}{1 + e^{\gamma}} = 0$$

$$\boxed{\frac{2}{4}}$$

p(1)

$$= -p_{avg} + \frac{e^{\gamma}}{1 + e^{\gamma}} = 0$$

p_{avg}

$$\frac{e^{\gamma}}{1 + e^{\gamma}} = p_{avg}$$

$$e^{\gamma} = p_{avg} + p_{avg} e^{\gamma} \Rightarrow e^{\gamma} - p_{avg} e^{\gamma} = p_{avg}$$

$$e^{\gamma} (1 - p_{avg}) = p_{avg}$$

$$e^{\gamma} = \frac{p_{avg}}{1 - p_{avg}}$$

$$\gamma = \log\left(\frac{p_{avg}}{1 - p_{avg}}\right)$$

$$f_0(x) = \gamma = \log\left(\frac{p_{avg}}{1 - p_{avg}}\right)$$

$$\log\left(\frac{5/8}{1 - 5/8}\right) = \log\left(\frac{5/8}{3/8}\right) = \log\left(\frac{5}{3}\right)$$

Step 2.a -> Pseudo Residuals

10 August 2023 11:23

n rows

2. For $m = 1$ to M :

(a) For $i = 1, 2, \dots, N$ compute

$$r_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$$

$L = -y \log(\text{odds}) + \log(1 + e^{\log(\text{odds})})$

$\frac{\partial L}{\partial \log(\text{odds})} = -y + \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$

$= -y + p = y - p \rightarrow y_i - p_i$

$f(x) \rightarrow f_0(x) \rightarrow \log \text{odds}$

$\rightarrow \text{prob output}$

$\rightarrow \log \text{odds}$

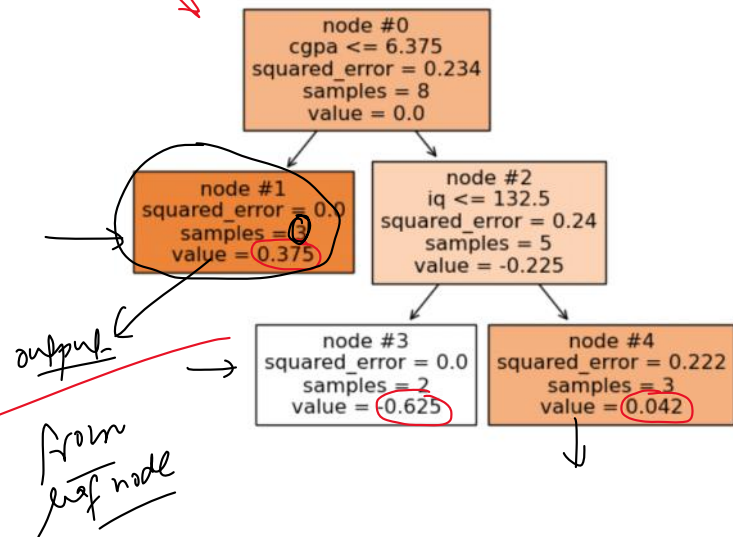
Step 2.b -> Train Regression Tree

10 August 2023 11:41

cgpa	iq	is_placed	pre1(log-odds)	pre1(probability)	res1
6.82	118	0	0.510826	0.625	-0.625
6.36	125	1	0.510826	0.625	0.375
5.39	99	1	0.510826	0.625	0.375
5.50	106	1	0.510826	0.625	0.375
6.39	148	0	0.510826	0.625	-0.625
9.13	148	1	0.510826	0.625	0.375
7.17	147	1	0.510826	0.625	0.375
7.72	72	0	0.510826	0.625	-0.625

fol(x)

regression



$$\frac{\sum \text{Residual}}{\sum [\text{PreviousProb} * (1 - \text{PreviousProb})]}$$

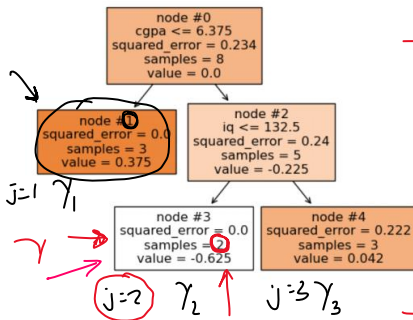
Step 2.c -> Compute Lambda for all leaf nodes

10 August 2023 12:02

$m = 1$

(c) For $j = 1, 2, \dots, J_m$ compute

$$\gamma_j = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma)$$



cgpa	iq	is_placed	prel(log-odds)	prel(probability)	res1	leaf_entry1
6.82	118	0	0.510826	0.625	-0.625	3
6.36	125	1	0.510826	0.625	0.375	1
5.39	99	1	0.510826	0.625	0.375	1
5.50	106	1	0.510826	0.625	0.375	1
6.39	148	0	0.510826	0.625	-0.625	4
9.13	148	1	0.510826	0.625	0.375	4
7.17	147	1	0.510826	0.625	0.375	4
7.72	72	0	0.510826	0.625	-0.625	3

$$\gamma_2 = \arg \min_{\gamma} L(y_1, f_0(x_1) + \gamma) + L(y_8, f_0(x_8) + \gamma)$$

$$\frac{\partial L}{\partial \gamma} = 0$$

minimize $\frac{\text{res1}}{p_i(1-p_i)} + \frac{\text{res8}}{p_i(1-p_i)}$

Taylor Series (w/o 2nd derivative)

$$f(x) \approx f(a) + \frac{f'(a)}{1!} (x-a) + \frac{f''(a)}{2!} (x-a)^2 + \dots$$

$$\gamma_2 = \frac{\sum \text{residual}_i}{\sum p_i(1-p_i)}$$

$$L = L(y_i, f_0(x_i)) + \frac{\partial L}{\partial f_0(x_i)} \gamma + \frac{1}{2} \frac{\partial^2 L}{\partial f_0(x_i)^2} \gamma^2$$

$$\frac{\partial L}{\partial \gamma} = \frac{\partial L}{\partial f_0(x_i)} + \frac{\gamma}{2} \frac{\partial^2 L}{\partial f_0(x_i)^2} = 0$$

$$\gamma \frac{\partial^2 L}{\partial f_0(x_i)^2} = - \frac{\partial L}{\partial f_0(x_i)}$$

$$\gamma = \frac{- \frac{\partial L}{\partial f_0(x_i)}}{\frac{\partial^2 L}{\partial f_0(x_i)^2}} \rightarrow \text{residual w/o row 1 } (y_i - p_i)$$

$$p_i(1-p_i)$$

$$\frac{\partial L}{\partial f_0(x_i)} = - \left(y_i - \frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}} \right) = y_i - p_i \quad \text{residual}$$

$$\left[\frac{\partial \mathcal{L}}{\partial f_0(x_i)} = - \left(-y_i + \frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}} \right) \right] = \dots$$

$$\frac{\partial}{\partial f_0(x_i)} \frac{\partial \mathcal{L}}{\partial f_0(x_i)} = \frac{\partial}{\partial \log(\text{odds}_i)} \left(y_i - \frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}} \right)$$

$$\frac{\partial}{\partial \log(\text{odds}_i)} \left[\underbrace{y_i}_x - \underbrace{\frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}}}_{\frac{y}{x' y + x y'}} \right]$$

$$= - \left[\frac{e^{\log(\text{odds}_i)}}{(1 + e^{\log(\text{odds}_i)})^2} - \frac{e^{\log(\text{odds}_i)}}{(1 + e^{\log(\text{odds}_i)})^2} \right]$$

$$= \frac{e^{\log(\text{odds}_i)}}{(1 + e^{\log(\text{odds}_i)})^2} - \frac{e^{\log(\text{odds}_i)}}{(1 + e^{\log(\text{odds}_i)})^2}$$

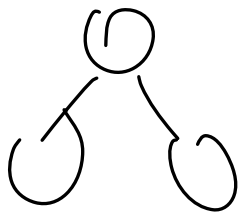
$$= \frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}} \left[\frac{1}{1 + e^{\log(\text{odds}_i)}} - 1 \right] \rightarrow \frac{1}{1 + e^{\log(\text{odds}_i)}}$$

$$= \boxed{\frac{e^{\log(\text{odds}_i)}}{1 + e^{\log(\text{odds}_i)}}} \boxed{\frac{1}{1 + e^{\log(\text{odds}_i)}}} = p_i (1 - p_i)$$

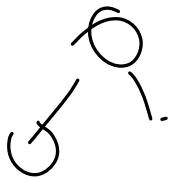
Step 2d - Update the model

10 August 2023 12:10

$$\left[f_1(x) = f_0(x) + \underbrace{\text{output from dt}}_{\substack{\downarrow \\ \text{log odds}}} \left\{ \text{some leaf node} \right\} \right]$$



$f_2(x)$



$f_3(x) \rightarrow M \text{ deis}$

Step 3 - Final Model

10 August 2023 12:10

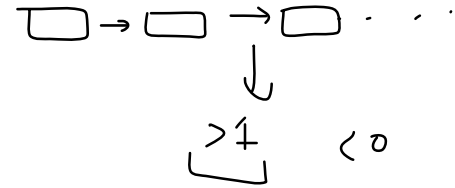
$$f_M(x) = f_{m-1}(x) + \text{last decision output}$$

Boosting model

Log(odds) Vs Probability

10 August 2023 13:06

$-\infty$ $+\infty$ $(0-1)$



- ① **Unconstrained Prediction Space:** Log odds can span the entire real line ($-\infty$ to $+\infty$), while probabilities are constrained between 0 and 1. Algorithms like gradient boosting involve adding corrections (via the weak learners) to the predictions iteratively. If you're working in the log odds space, there's no need to worry about your predictions going out of bounds.

$-\infty - +\infty$

- ② **Better Gradients:** When computing gradients (which guide the addition of new trees in boosting), the gradients can be more informative and have better magnitudes in the log odds space than in the probability space, especially when probabilities are near 0 or 1.

