$$F = \{f_1, f_2, f_3, f_4 - \cdots - f_m\} \text{ set of base learners}$$

Final prediction: $\hat{y}_i = \sum\limits_{t=1}^{m} f_t(x_i)$

$$D = \{x_1, x_2, x_3 - \cdots - x_n\} \text{ data points}$$

$$L^{<t>} = \sum\limits_{i=1}^{n} \ell(y_i, \hat{y}_i^{<t-1>} + f_t(x_i)) + \Omega(f_t) \quad\text{——}①$$

Taylor Expansion

$$f(a+h) = f(a) + f'(a)h + \frac{1}{2}f''(a)h^2 + \cdots\cdots f^n(a)\frac{h^n}{n!}$$

Here $a = \hat{y}_i^{<t-1>}$

$h = f_t(x_i)$

$$f(a) = \ell(y_i, \hat{y}_i^{<t-1>})$$

$$\therefore L^{<t>} = \sum\limits_{i=1}^{n} \ell(y_i, \hat{y}_i^{<t-1>}) + \left(\frac{\partial \ell(y_i, \hat{y}_i^{<t-1>})}{\partial \hat{y}_i^{<t-1>}}\right) f_t(x_i) + \left(\frac{\partial^2 \ell(y_i, \hat{y}_i^{<t-1>})}{\partial \hat{y}_i^{<t-1>2}}\right) f_t(x_i)^2$$

$\ell(y_i, \hat{y}_i^{<t-1>})$ is constant irrespective of any function

$$\therefore L^{<t>} = \sum\limits_{i=1}^{n} C + g_i f_t(x_i) + h_i f_t(x_i)) + \Omega(f_t)$$

Pick $f_t(x_i) \ni L^{<t>}$ is minimum. Removing constant as it is

equal for any function.

$$L^{<t>} = \sum\limits_{i=1}^{n} (g_i f_t(x_i) + h_i f_t(x_i)) + \Omega(f_t) \quad\longrightarrow②$$

Let $f_t(x)$ has $K$ leaf nodes. $I_j$ be the set of instances

belonging to node 'j'. '$w_j$' be the prediction for node 'j'.

$$\Omega(f_t) = \gamma K + \frac{1}{2}\lambda \sum\limits_{j=1}^{K} w_j^2$$

$$L^{<t>} = \sum\limits_{j=1}^{K}\left[\left(\sum\limits_{i\in I_j} g_i\right)w_j + \frac{1}{2}\left(\sum\limits_{i\in I_j} h_i + \lambda\right)w_j^2\right] + \gamma K \longrightarrow③$$

For each leaf 'j', $\dfrac{dL^{<t>}}{dw_j*} = 0$

$$0 = \sum\limits_{i\in I_j} g_i + \frac{1}{2}\left(\sum\limits_{i\in I_j} h_i + \lambda\right) 2\times w_j*$$

$$\boxed{w_j^* = -\dfrac{\sum\limits_{i\in I_j} g_i}{\sum\limits_{i\in I_j} h_i + \lambda}}$$

Substituting weights into ③

$$L^{<t>} = -\frac{1}{2}\sum_{j=1}^{K}\frac{\left(\sum_{i\in I_j}g_i\right)^2}{\sum_{i\in I_j}h_i+\lambda} + \gamma K \qquad ——④$$

This is the best loss for a fixed base learner with 'K' nodes. There will be several hundreds of possible tree structures. It is impossible to explore all of them.



node (I)

$I_L$          $I_R$

leaf node $I_L$          leaf node $I_R$

loss according to ③

$$\frac{-\frac{1}{2}\left(\sum_{i\in I}g_i\right)^2}{\sum_{i\in I}h_i+\lambda} + \gamma(1)$$

$\downarrow$ reduction in loss

$$\frac{-\frac{1}{2}\left(\sum_{i\in I_L}g_i\right)^2}{\sum_{i\in I_L}h_i+\lambda} + \gamma(1)$$

$$\frac{-\frac{1}{2}\left(\sum_{i\in I_R}g_i\right)^2}{\sum_{i\in I_R}h_i+\lambda} + \gamma(1)$$

$$L_{split} = \frac{-\frac{1}{2}\left(\sum_{i\in I}g_i\right)^2}{\sum_{i\in I}h_i+\lambda} - \left[-\frac{1}{2}\left(\frac{\left(\sum_{i\in I_L}g_i\right)^2}{\sum_{i\in I_L}h_i+\lambda} + \frac{\left(\sum_{i\in I_R}g_i\right)^2}{\sum_{i\in I_R}h_i+\lambda}\right)\right]+\gamma-\gamma-\gamma$$

$$= \boxed{\frac{1}{2}\left[\frac{\left(\sum_{i\in I_L}g_i\right)^2}{\sum_{i\in I_L}h_i+\lambda} + \frac{\left(\sum_{i\in I_R}g_i\right)^2}{\sum_{i\in I_R}h_i+\lambda} - \frac{\left(\sum_{i\in I}g_i\right)^2}{\sum_{i\in I}h_i+\lambda}\right] - \gamma}$$