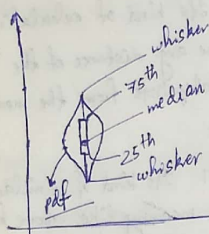


violin-plots \rightarrow (box plot + pdf)



* Always perform analysis oriented towards solving the object.
Don't deviate.

Pseudo Residuals Revisited

* We have seen that pseudo-residuals $\left(\frac{-\partial L}{\partial F_k(x)} \right)$ are similar to residuals/errors for squared loss function.

* What about other loss functions?

\rightarrow consider binary log-loss.

$$L = y_i \cdot \log(p_i) + (1 - y_i) \cdot \log(1 - p_i) \quad \text{where}$$

$$p_i = \frac{1}{1 + e^{-\hat{y}_i}}$$

$$\text{So, } L = y_i \cdot \log\left(\frac{1}{1 + e^{-\hat{y}_i}}\right) + (1 - y_i) \log\left(\frac{e^{-\hat{y}_i}}{1 + e^{-\hat{y}_i}}\right)$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{y_i - (1 - y_i) \cdot e^{\hat{y}_i}}{(1 + e^{\hat{y}_i})} \quad \text{Multiplying by } \frac{e^{-\hat{y}_i}}{e^{-\hat{y}_i}}$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{y_i \cdot e^{-\hat{y}_i} + y_i - 1}{(1 + e^{-\hat{y}_i})} = \frac{y_i (e^{-\hat{y}_i} + 1)}{(e^{-\hat{y}_i} + 1)} - \frac{1}{(e^{-\hat{y}_i} + 1)}$$

$$= \underline{\underline{y_i - p_i}}$$

\hat{y}_i is same as $F_k(x)$.

and $(y_i - p_i)$ is nothing but residue.

$$\text{So, } \frac{-\partial L}{\partial F_k(x)} \propto \text{residual}$$

Hence, even for binary log-loss, the pseudo residuals are similar to the actual residuals.

sq-loss \rightarrow yes.
log-loss \rightarrow yes.
any-loss \rightarrow no

most common losses

* So, if pseudo residuals aren't equivalent to residuals always, why are we even using them?

* It has a more fundamental mathematical reason behind this rather than pseudo residuals just being similar to residuals.

* We have a dataset $D = \{x_i, y_i\}_{i=1}^n$

Loss function = L ,

* we have to find a function F that maps $x_i \rightarrow y_i$ minimizing the loss function.

* In boosting we find a sequence of functions.

* So, in accordance to 'empirical risk minimization' principle, we try to find $F(x)$ that minimizes average value of the loss function on the training set.

It does so by starting with a ~~constant~~ model, consisting of a constant function $F_0(x)$ and incrementally expanding in greedy fashion:-

$$F_1(x) = \arg\min_s \sum_{i=1}^n L(y_i, s)$$

Optimal Problem

$$F_m(x) = F_{m-1}(x) + \arg\min_{h_m \in H} \sum_{i=1}^n L(y_i, F_{m-1}(x) + h_m(x_i))$$

H = set of all base learners possible

Due to the infinite size of H , training the model isn't feasible.

so we try to find some approximations that are close to this optimization problem

We solved the optimization problem of weights using gradient descent.

$$w_{\text{new}} = w_{\text{old}} - \eta \left(\frac{\partial L}{\partial w} \right)_{\text{old}}$$

variable learning rate gradient update step.

Similarly we can generalize this for learning functions as well.

$$F_m(x) = F_{m-1}(x) - \eta_m \sum_{i=1}^n \nabla_{F_{m-1}} L(y_i, F_{m-1}(x_i))$$

learning rate \rightarrow calculated by simple line search.

this is nothing but pseudo-residual (including -2η sign).

Refer to prev notebook for actual algo of boosting using pseudo-residuals.

* So, this is the reason behind using pseudo-residuals.