

Mahan Veisi – ML HW1 – 400243081

1- میدانیم که مشکل اورفیت زمانی رخ میدهد که مدل ما به دلیل وابستگی زیاد به داده های موجود خودش (به دلیل کمبود دیتا یا حتی مشکل ساختاری مدل) نتواند روی داده های دیده نشده از قبل (تست دیتا) پردیکشن درستی داشته باشد.

با توجه به تعریف ارائه شده، میدانیم در رندوم فورست اول از همه به نوعی از یک زیرمجموعه از کل داده های موجود استفاده میکنیم زیرا داده ها برای train کردن هر درخت میتوانند تکراری (with replacement) و با شانس برابر باشند و پس از پایان ساختن درخت ها و ساخت جنگل، امکان داشته باشد که هر درخت داده هایی را بررسی کند که دیده نشده باشند (bagging). از طرفی میدانیم معمولاً یکی از نواقص ساخت تک درخت ها (درخت تصمیم)، اورفیت شدن میباشد.

پس اگر ما در رندوم فورست بتوانیم دقت train را بالا ببریم، نسبت به گذشته (که از درخت تصمیم استفاده میکردیم) اطمینان بیشتری داریم که در مقابل داده های دیده نشده، به overfit دچار نشویم (اختلاف بین دقت train و test کمتر است). در نهایت برای پیشبینی کردن از همه درخت ها استفاده میکنیم که این خود نیز یک نوع تکنیک برای regularization است.

2- الف)

$$H(Y|X) = \sum_{x \in X} p(x) H(Y|x) \quad \text{سؤال 2 الف)$$

$Y = \{\text{شكارة}\}$ ، $X = \{\text{نبون و كوكو ديل بون}\}$

$$H(Y|X) = \frac{44}{100} H(Y| \text{کردگودیل})$$

$$+ \frac{34}{100} H(Y| \text{نبودن کردگودیل})$$

$$H(Y| \text{کردگودیل}) = -P(Y| \text{کردگودیل}) \cdot \log_2(P(Y| \text{کردگودیل}))$$

$$-P(\sim Y| \text{کردگودیل}) \cdot \log_2(P(\sim Y| \text{کردگودیل}))$$

$$P(Y| \text{کردگودیل}) = \frac{P(Y, \text{کردگودیل})}{P(\text{کردگودیل})} = \frac{\frac{31}{100}}{\frac{44}{100}} = \frac{31}{44}$$

$$P(\sim Y| \text{کردگودیل}) = \frac{P(\sim Y, \text{کردگودیل})}{P(\text{کردگودیل})} = \frac{\frac{13}{100}}{\frac{44}{100}} = \frac{13}{44}$$

$$H(Y| \text{نکر}) = -P(Y| \text{نکر}) \cdot \log_2(P(Y| \text{نکر}))$$

$$-P(\sim Y| \text{نکر}) \cdot \log_2(P(\sim Y| \text{نکر}))$$

$$P(Y| \text{نکر}) = \frac{P(Y, \text{نکر})}{P(\text{نکر})} = \frac{\frac{14}{100}}{\frac{34}{100}} = \frac{14}{34}$$

$$P(\sim Y| \text{نکر}) = \frac{P(\sim Y, \text{نکر})}{P(\text{نکر})} = \frac{\frac{20}{100}}{\frac{34}{100}} = \frac{20}{34}$$

$$\Rightarrow \text{final} = \frac{48}{100} \left(\begin{array}{c} -0,702 \\ -1,34 \end{array} \left(-\frac{38}{48} \log_2 \frac{38}{48} - \frac{10}{48} \log_2 \frac{10}{48} \right) \right. \\ \left. + \frac{104}{100} \left(\begin{array}{c} -1,34 \\ -0,17 \end{array} \left(-\frac{18}{104} \log_2 \frac{18}{104} - \frac{22}{104} \log_2 \frac{22}{104} \right) \right) \right)$$

$$\Rightarrow \frac{48}{100} \left(0,88 + 0,87 \right) + \frac{104}{100} \left(0,87 + 0,87 \right)$$

0,43 0,38

$$\Rightarrow \text{final} = 0,97$$

)

$$IG(Y|X) = H(Y) - H(Y|X)$$

$$\Rightarrow IG(Y|X) = H(Y) - H(Y|J)$$

شماره کاغذ
↓
J
Previously = 0.98
(0.98 + 0.02)

$$H(Y) = - \sum_{y \in Y} p(y) \log_2 p(y) = - \left(\underbrace{\frac{0.98}{1} \log_2 \frac{0.98}{1}}_{0.014} + \underbrace{\frac{0.02}{1} \log_2 \frac{0.02}{1}}_{0.014} \right)$$

-1

$$\Rightarrow IG(Y|X) = 1 - 0.98 = 0.02$$

boofi

(ج)

$$H(\text{نک} | \sim Y)$$

↓
شماره

(2)

$$= -P(\text{نک} | \sim Y) \log_2 P(\text{نک} | \sim Y)$$

$$- P(\text{ک} | \sim Y) \log_2 P(\text{ک} | \sim Y)$$

$$P(\text{هر} | \sim Y), \frac{P(\text{نک}, \sim Y)}{P(\sim Y)}, \frac{\frac{22}{100}}{\frac{48}{100}} = \frac{22}{48}$$

$$P(\text{ک} | \sim Y), \frac{P(\text{ک}, \sim Y)}{P(\sim Y)}, \frac{\frac{24}{100}}{\frac{48}{100}} = \frac{24}{48}$$

$$\Rightarrow H(\text{نک} | \sim Y) = \underbrace{-\frac{22}{48} \log_2 \frac{22}{48}}_{0.541} - \underbrace{\frac{24}{48} \log_2 \frac{24}{48}}_{0.41}$$

$$\Rightarrow H(\text{نک} | \sim Y) = \underline{0.99}$$