

الف) لزوما درست نمیباشد

در اینجا مسئله bias/variance tradeoff مطرح است ، افزایش پیچیدگی مدل منجر به کاهش خطای آموزش میشود و زمانی که مدل ما بیش از حد پیچیده باشد معضل high variance و overfitting داریم که باعث افزایش خطای ازمون می شود و برعکس زمانی که مدل ما به شدت ساده باشد معضل high bias و underfitting داریم که باعث میشود خطای آموزش و ازمون افزایش یابد زیرا مدل به سختی نظر خود را عوض نمیکند.

ب) درست

هرچه داده های کمتری داشته باشیم، مدل ما بهتر می تواند استثنایها را در مجموعه آموزشی شما به خاطر بسپارد که منجر به دقت بالا در آموزش اما دقت پایین در مجموعه ازمون می شود زیرا مدل ما چیزهایی را که از مجموعه آموزشی کوچک آموخته است تعمیم می دهد.

ج) درست

افزایش پیچیدگی مدل باعث overfitting مدل بر روی داده های آموزش میشود که کاهش خطای آموزش و افزایش خطای ازمون را به همراه دارد (د) MAE برای داده های پرت و نویزی عملکرد بهتری دارد چون از قدر مطلق برای محاسبات استفاده میکند . هرچند در پیاده سازی توابع محاسبه هزینه متریک به صورت پیش فرض از RMSE استفاده می شود ولی به دلیل تفسیر ساده تر MAE (easy to interpretation) الگوریتم MAE بهتر است

۲-

مدل با پیچیدگی بالا نویز داده های آموزش را بخاطر میسپارد و نمیتواند به مجموعه ارزیابی تعمیم دهد یا زمانی که training epochs های زیادی تکرار می شود خطای مجموعه آموزش کاهش می یابد ولی خطای مجموعه ارزیابی شروع به افزایش می کند

- early stopping
- Cross validation
- Regularization

۳-

الف , ج)

$$Y = XW + \epsilon$$

$$\epsilon = Y - XW \rightarrow \|\epsilon\| = \|Y - XW\|$$

$$\|\epsilon\|^2 = \|Y - XW\|^2 = (Y - XW) \cdot (Y - XW) = (Y - XW)^T (Y - XW)$$

$$\begin{aligned} J(W) &= (Y - XW)^T (Y - XW) = ((XW)^T - Y^T)(XW - Y) = \\ &= (XW)^T XW - (XW)^T Y - Y^T XW + Y^T Y = \\ &= (XW)^T XW - 2(XW)^T Y + Y^T Y \end{aligned}$$

$$\frac{dJ}{dW} = 2X^T XW - 2X^T Y = 0 \rightarrow W = (X^T X)^{-1} X^T Y$$

$$J(W) = \sum_{i=1}^m (h(x^i) - y^i)^2 + \lambda \sum_{j=1}^n W_j^2 \quad (ع)$$

$$W = (X^T X + \lambda \cdot I)^{-1} X^T Y$$

ب) زمانی که ماتریس مربعی $X^T X$ وارون نداشته باشد نمیتوانیم از این رابطه به طور مستقیم استفاده کنیم که در دو حالت وارون نداریم :

- ویژگی ها مستقل از هم نباشند و ستون وابسته بهم داشته باشیم
- $\#m \ll \#n$

وقتی ماتریس وارون ندارد میتوان از $sodu\ inverse$ استفاده کرد ولی جواب تقریبی و غیر منحصرفرد است یا میتوان از $regularization$ استفاده کرد

برای کار با تعداد $feature$ کم گزینه خوبی است اگر تعداد $feature$ زیاد باشد محاسبات پیچیده میشود

د) چون قصد داریم $cost\ function$ را $minimum$ کنیم جمله نرمال ساز که اضافه کردیم نیز min میشود و باعث می شود W_j ها که ضرایب چندجمله ای هایمان هستند نیز کم شوند در نتیجه جلوی زیاد شدن X ها با توان بالا گرفته میشود.

-۴-

الف) SGD یا $MSGD$ بهترین گزینه هستند زیرا نیاز به بارگذاری کامل مجموعه آموزش برای انجام یک مرحله از $Gradient\ descent$ نداریم و $normal\ equation$ گزینه مناسبی نیست زیرا پیچیدگی محاسباتی در بدست آوردن معکوس ماتریس داریم

ب) روش $normal\ equation$ نیازمند $feature\ normalizing$ نیست به همین دلیل برای مجموعه داده با مقیاس های پراکنده عملکرد بهتری دارد ولی در عوض روش $gradient\ descent$ برای همگرایی سریع تر نیازمند $feature\ scaling$ می باشد

-۵-

افزودن داده آموزشی منجر به کاهش $variance$ می شود. اما تاثیری روی $bias$ ندارد

-۶-

خط چین مشکی خطای آموزش

خط چین زرد مربع بایاس

خط قرمز خطای ناشی از نویز

خط چین آبی واریانس

خط سبز خطای آزمون



