

حاشیة مصنوعی - تمرین تدریسی سری چهارم

محاسنر مخاطاری - 98102346

سوال 1  
الف)

$$\{\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}^2\}_{mle} = \underset{\beta_0, \beta_1, \sigma}{\operatorname{argmax}} \prod_{n=1}^N p(\{x_n, y_n\} | \beta_0, \beta_1, \sigma) = \underset{\beta_0, \beta_1, \sigma}{\operatorname{argmax}} \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_n - \beta_0 - \beta_1 x_n)^2}{2\sigma^2}\right)$$

$$\left( \mathcal{N}(y_n | \beta_0 + \beta_1 x_n, \sigma^2) \right) \xrightarrow{\text{MLE}} \underset{\beta_0, \beta_1, \sigma}{\operatorname{argmax}} \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_n - \beta_0 - \beta_1 x_n)^2}{2\sigma^2}\right) \xrightarrow{NLL}$$

دارایان → اید ریاضی توزیع

$$\underset{\beta_0, \beta_1, \sigma}{\operatorname{argmin}} \left( \frac{N}{2} \log(2\pi) + \frac{N}{2} \log(\sigma^2) + \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \beta_0 - \beta_1 x_n)^2 \right) \xrightarrow{NLL} (*)$$

$$\hat{\beta}_0: \frac{\partial}{\partial \beta_0} NLL = 0 = \frac{1}{\sigma^2} \sum_{n=1}^N (y_n - \beta_0 - \beta_1 x_n) \Rightarrow 0 = \sum_{n=1}^N y_n - N\hat{\beta}_0 - (\sum_{n=1}^N x_n) \hat{\beta}_1 \xrightarrow{/N}$$

$$\Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1: \frac{\partial}{\partial \beta_1} NLL = 0 = \frac{1}{\sigma^2} \sum_{n=1}^N (x_n y_n - \hat{\beta}_0 x_n - \beta_1 x_n^2) \rightarrow 0 = \sum_{n=1}^N x_n y_n - \hat{\beta}_0 \sum_{n=1}^N x_n - \hat{\beta}_1 \sum_{n=1}^N x_n^2$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \rightarrow 0 = \sum_{n=1}^N x_n y_n - (\bar{y} - \hat{\beta}_1 \bar{x}) \sum_{n=1}^N x_n - \hat{\beta}_1 \sum_{n=1}^N x_n^2 \rightarrow \hat{\beta}_1 = \frac{\sum_{n=1}^N x_n y_n - \sum_{n=1}^N \bar{x} \bar{y}}{\sum_{n=1}^N x_n^2 - \sum_{n=1}^N \bar{x}^2}$$

(\*) : برای تخمین MLE باید NLL بدست آمده کمینه شود (بر حسب  $\beta_0, \beta_1$ ) که با توجه به NLL بدست آمده معادل کمینه کردن  $\sum_{n=1}^N (y_n - \beta_0 - \beta_1 x_n)^2$  است که همان مجموع مربعات خطا است

$$E\hat{\beta}_1 = E\left[ \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] = E\left[ \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] = E\left[ \frac{\sum_{i=1}^n (x_i - \bar{x})(\beta_0 + \beta_1 x_i + \varepsilon_i)}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

$$* : \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x}) y_i - \bar{y} \sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n (x_i - \bar{x}) y_i$$

جمله  $\sum_{i=1}^n (x_i - \bar{x}) \beta_0$  صفر بودن حذف می شود

$$\rightarrow E[\hat{\beta}_1] = E\left[ \frac{\sum_{i=1}^n (\beta_0 x_i + \varepsilon_i)(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \left( \sum_{i=1}^n \beta_0 x_i + E[\sum_{i=1}^n \varepsilon_i x_i] \right)$$

$$\underline{E[\varepsilon_i] = 0} \quad \frac{\beta_1 \sum_{i=1}^n x_i (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\beta_1 \sum_{i=1}^n (x_i^2 - x_i \bar{x})}{\sum_{i=1}^n (x_i^2 - 2x_i \bar{x} + \bar{x}^2)} = \beta_1 \Rightarrow \hat{\beta}_1 \text{ is unbiased}$$

$$E[\hat{\beta}_0] = E[\bar{y} - \hat{\beta}_1 \bar{x}] = E\bar{y} - E\hat{\beta}_1 \times \bar{x} = E[\beta_0 + \beta_1 x_i + E\varepsilon_i] - \beta_1 \bar{x} = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0 \Rightarrow \hat{\beta}_0 \text{ is unbiased}$$

$$VAR[\hat{\beta}_1] = VAR\left[\frac{\sum (x_i - \bar{x})(\beta_0 + \beta_1 x_i + \varepsilon_i)}{\sum (x_i - \bar{x})^2}\right] \xrightarrow[\text{مقادیر تصادفی نیستند}]{\text{مقادیر تصادفی نیستند}} VAR\left[\frac{\sum (x_i - \bar{x}) \varepsilon_i}{\sum (x_i - \bar{x})^2}\right]$$

$$\xrightarrow{VAR[aX] = a^2 VAR X} VAR[\hat{\beta}_1] = \left\{ \frac{\sum (x_i - \bar{x})^2}{(\sum (x_i - \bar{x})^2)^2} \right\} VAR[\varepsilon_i] = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

$\varepsilon_i$  are independent

$$VAR[\hat{\beta}] = VAR[\bar{y} - \hat{\beta}_1 \bar{x}] = VAR\left[\frac{1}{n}(\beta_0 \sum x_i + n\beta_0 + \sum \varepsilon_i) - \hat{\beta}_1 \sum x_i\right]$$

*deterministic*  $\rightarrow VAR=0$

$$= VAR\left[\frac{1}{n}(\sum \varepsilon_i - (\sum x_i) \hat{\beta}_1)\right] = \frac{1}{n^2} \left( n\sigma^2 + (n\bar{x})^2 \frac{\sigma^2}{\sum (x_i - \bar{x})^2} - 2n\bar{x} \frac{\sum (x_i - \bar{x}) \varepsilon_i}{\sum (x_i - \bar{x})^2} \right)$$

$$= \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{\sum (x_i - \bar{x})^2} - 2\bar{x} \frac{\sum (x_i - \bar{x}) \varepsilon_i}{n \sum (x_i - \bar{x})^2}$$

$$= \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum (x_i - \bar{x})^2} = \frac{\sigma^2 (n\bar{x}^2 + \sum (x_i - \bar{x})^2)}{n \sum (x_i - \bar{x})^2}$$

$$\xrightarrow{\sigma^2 \sum x_i^2} \frac{\sigma^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x}) y_i - \sum (x_i - \bar{x}) \bar{y}}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2}$$

$\sum (x_i - \bar{x}) = 0 \checkmark \Rightarrow Y_i = (x_i - \bar{x})$  به طبق رابطه بالا به هم مشابه رسید

$$E \tilde{\beta}_1 = E \frac{\sum Y_i y_i}{\sum Y_i x_i} = E \frac{\sum Y_i (\beta_0 + \beta_1 x_i + \varepsilon_i)}{\sum Y_i x_i} = E \frac{\beta_0 \sum Y_i + \beta_1 \sum Y_i x_i + \sum Y_i \varepsilon_i}{\sum Y_i x_i}$$

$\sum Y_i = 0$

$$\tilde{\beta}_1 = \frac{\beta_1 \sum Y_i x_i + \sum Y_i \varepsilon_i}{\sum Y_i x_i} = \beta_1 + \frac{\sum Y_i \varepsilon_i}{\sum Y_i x_i}$$

$(VAR(\beta_1, \beta_0) = 0)$   
*deterministic*

$$VAR \tilde{\beta}_1 = VAR\left[\frac{\sum Y_i y_i}{\sum Y_i x_i}\right] = VAR\left[\frac{\sum Y_i (\beta_0 + \beta_1 x_i + \varepsilon_i)}{\sum Y_i x_i}\right] = VAR\left[\frac{\sum Y_i \varepsilon_i}{\sum Y_i x_i}\right]$$

$$= \frac{\sum Y_i^2 VAR \varepsilon_i}{(\sum Y_i x_i)^2} = \sigma^2 \frac{\sum Y_i^2}{(\sum Y_i x_i)^2}$$

$$\frac{VAR \hat{\beta}_1}{VAR \tilde{\beta}_1} = \frac{\frac{\sigma^2}{\sum (x_i - \bar{x})^2}}{\sigma^2 \frac{\sum Y_i^2}{(\sum Y_i x_i)^2}} = \frac{(\sum Y_i x_i)^2}{(\sum (x_i - \bar{x})^2) \sum Y_i^2}$$

$$\frac{(\sum Y_i x_i)^2}{(\sum Y_i^2)(\sum (x_i - \bar{x})^2)} \leq 1 \Rightarrow VAR \hat{\beta} \leq VAR \tilde{\beta}$$

$\sum Y_i = 0 \Rightarrow \begin{cases} \sum Y_i^2 \geq \sum Y_i x_i \\ \sum (x_i - \bar{x})^2 \geq \sum Y_i x_i \end{cases}$



$$\min_w \lambda W^T W + \|XW - Y\|_2^2$$

(الف)

- یک بردار مقدار اولیه برای  $W$  انتخاب می‌کنیم و مقدار  $\eta$  (ضریب یادگیری) را انتخاب می‌کنیم
- تا وقتی که  $F(W)$  از سطح آستانه‌ای بزرگتر است، کار زیر را انجام می‌دهیم:
- از بین  $n$  نمونه‌های  $(x_i, y_i)$  دسته‌ای را به صورت تصادفی انتخاب می‌کنیم (دسته‌ای به اندازه  $\eta$ )
- برای  $n$  نمونه‌های انتخابی در مرحله قبل، کار زیر را انجام می‌دهیم:

$$W \leftarrow W - \frac{\eta}{n} \sum_{i=1}^n \nabla F_i(W)$$

$W$  is vector so,  $\|W\|_2 = \sqrt{W^T W}$  (ب)  
 argument, large  $\|W\|_2$  would be penalized (Tikhonov regularization) so,  $\|W_2\|_2 \leq \|W\|_2$  adding this factor to argmin

به بیان دیگر اضافه کردن  $\lambda W^T W$  باعث کوچک شدن فضای جواب  $W$  برای مینه کردن  $L(W)$  است که به این معنی است که  $L(W_1) < L(W_2)$ ، حال اگر  $\|W_1\|_2 < \|W_2\|_2$  باشد به این معنی است که  $W_2$  پاسخ مناسبی برای مسئله دوم نیست (چرا که  $W_1$  پاسخ بهینه‌تری است) که خلاف فرض است  $W_2$  پاسخ مسئله دوم باشد پس  $\|W_1\|_2 \leq \|W_2\|_2$  است

## سوال 3

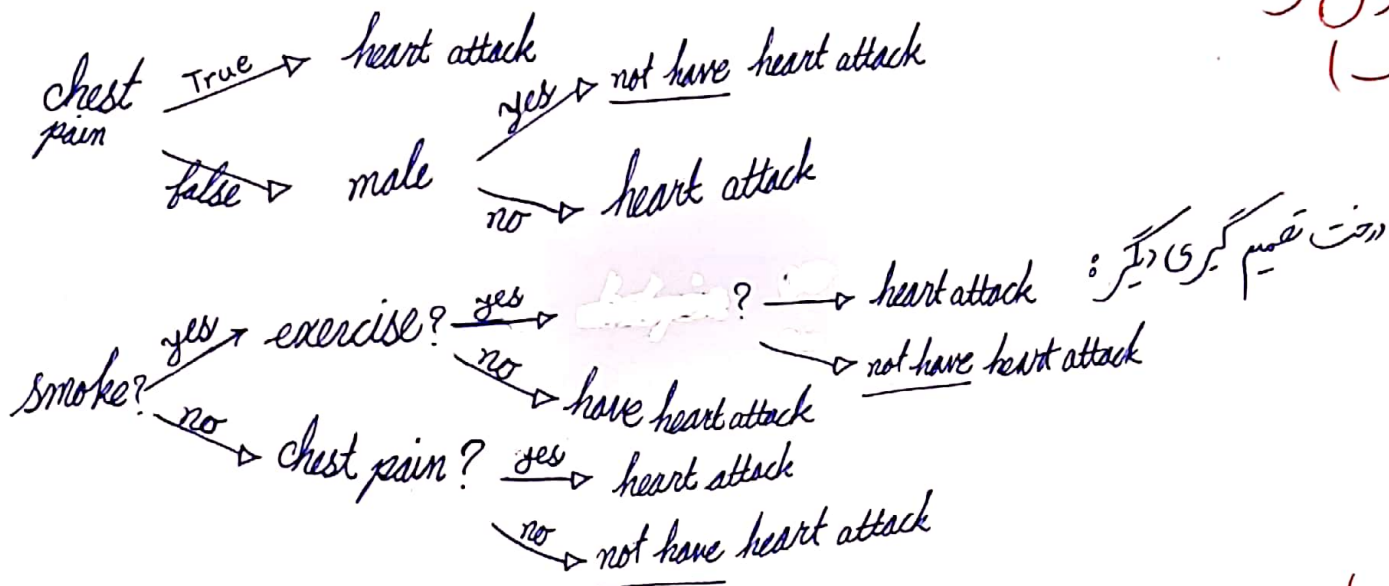
این ویژگی در نگاه اول تاثیر بسیار زیادی روی مدل نداشته و باید حفظ شود چرا که  $n$  را به شدت تحت تاثیر قرار می‌دهد و  $n$  اثر پذیری زیادی از این ویژگی دارد؛ اما «حقیقت بزرگی ضرایب ویژگی» به تنهایی موجب اثرگذاری خیلی زیاد این پارامتر نیست چرا که این ضریب بسیار بزرگ می‌تواند به موجب واحد کوچک آن ویژگی باشد (مثلاً حسب ملل بودن جای بر حسب  $k$  بودن) و می‌توان  $\sigma$  بدون دست داشتن اطلاعات بیشتر «مورد این ویژگی» نظر داد (این اطلاعات بیشتری ترانه اندازه standard deviation باشد و پارامتر مناسب می‌تواند حاصل ضرب  $\sigma x_i$  در ضرب محاسبه شده باشد

$\sigma_i$   
 $\rho$  نمونه‌ها  
 $\hookrightarrow$  ویژگی  $i$ ام

#### سوال 4

- غلط؛ گاهی در محاسبه  $\text{confusion}$  به رابطه ای بر حسب تعداد غونه ای رسم که نشان دهنده اثربخشی  $\text{f1}$  از تعداد غونه ای است.
- غلط؛ بدلیل وجود نویز در داده ای آموزش همچنین عدم تطابق رفتار مدل با رفتار واقعی فرایند، کم کردن زیاد خطا منجر به  $\text{overfit}$  می شود که می تواند منجر به این شود داده ای تست خطای زیادی داشته باشند (غونه ای از  $\text{overfit}$  در اسلاید 6 صفحه 8 وجود دارد)
- به طور کلی غلط است؛ افزایش پیچیدگی تا چندین مرحله می تواند باعث بهبود عملکرد مدل هم در  $\text{training set}$  هم در داده تست شود اما افزایش پیچیدگی بیش از اندازه موجب افزایش خطا در داده تست شود (اصولاً علت وجود داده تست برانداز کردن های پارامترها (مانند پیچیدگی مدل) است)

#### سوال 5 الف)



ب)

از دخت بالایی استفاده می کنیم

- اگر درد سینه دارند، احتمال حمله قلبی وجود دارد

- اگر درد سینه ندارند و مرد هستند، احتمال حمله قلبی وجود ندارد

- و زن هستند، احتمال حمله قلبی وجود دارد

- #### سوال 6
- می توان هر دسته بند در دوی را بارش باینری به طول  $d$  ساخت، همچنین با پرسش  $0 \leq i \leq d$  که کدام شخص می شود که  $d$  صفر است یا نه؛ به این روش در بزرگ های این دخت در دوی کامل همه رشته ای باینری ممکن به طول  $d$  است ساخته می شود که هر رشته در تناظر یک به یک با یک یا عضو مجموعه دسته بند ای  $\{0, 1\}^d$  خواهد بود، پس دخت تصمیم گیری حداکثر تا عمق  $d$  تا گره دارد (که یعنی عمق  $d+1$  دارد) (در صورتی که  $d$  تایی مقابل ساده شدن باشد این امکان وجود دارد که حداکثر عمق این دخت کمتر باشد)