

1.

$$L_{MSE}(\theta) = \frac{1}{n} \sum (x_i - \theta)^2 \Rightarrow \theta^* = \underset{\theta}{\operatorname{argmin}} \frac{1}{n} \sum (x_i - \theta)^2 = \underset{\theta}{\operatorname{argmin}} \sum (\theta - x_i)^2 \quad (a)$$

$$\Rightarrow \frac{d}{d\theta} \sum (\theta - x_i)^2 = 0 \Rightarrow \sum (\theta - x_i) = 0 \Rightarrow \sum_{i=1}^n \theta = \sum_{i=1}^n x_i$$

$$\Rightarrow n\theta = \sum x_i \Rightarrow \theta^* = \frac{\sum x_i}{n} \rightarrow \underline{\text{میانگین نمونه ما.}}$$

$$L_{MAE}(\theta) = \frac{1}{n} \sum |x_i - \theta| = \frac{1}{n} \left(\sum_{f.i. x_i > \theta^*} (x_i - \theta) + \sum_{f.i. x_j < \theta^*} (\theta - x_j) + \sum_{f.i. x_k = \theta^*} (0) \right)$$

Suppose there are m_1 x_j which are less than θ^* .

Suppose there are m_2 x_j which are greater than θ^* .

$$\frac{d}{d\theta} \frac{1}{n} \left(\sum_{f.i. x_i > \theta^*} (x_i - \theta) + \sum_{f.i. x_j < \theta^*} (\theta - x_j) \right) = 0 \Rightarrow \frac{1}{n} (-m_2 + m_1) = 0$$

$\Rightarrow \boxed{m_1 = m_2} \rightarrow$ پس θ^* باید به گونه ای انتخاب شود که تعداد باری از x_i ها از آن کم‌تر و بیش‌تر باشند. پس θ^* میان نمونه ما است.

(b) نکته اول در رابطه با بحث مشتق پذیری است. از MSE می‌توان مشتق گرفت و محاسبات را ساده‌تر کرد. مقدار بهینه را راحت‌تر می‌دهد اما MAE مشتق پذیر نیست و نیاز به حالت بندی دارد.

نکته دوم در رابطه با انعطاف پذیری آن‌ها در مواجه با داده پرت است. همانگونه که مشاهده شد MSE به مقدار داده‌ها بسیار دالیتی نشان می‌دهد و حتی با داده پرت بسیار بزرگی می‌تواند مقدار بهینه را به سمت سمت تأثیر قرار دهد. اما در MAE میانگین است و باید داده پرت خیلی بزرگ باید داده صرفاً نزدیک‌تر از میانگین تعدادی در میزان تغییر مقدار بهینه ندارد. پس MAE به معنای کمتر حساس است تا MSE که به سمت پرت‌های بیشتر می‌تواند به عنوان مزیت یا ایراد شناخته شود. بهتر است MAE به داده‌های نزدیک حساس‌تر شود و MSE بی‌محابا تر به داده دور.

$$\text{Huber loss} = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta \\ \delta(|a| - \frac{1}{2}\delta) & \text{o.w.} \end{cases}$$

(6)

$$\cosh = f(x) = \ln\left(\frac{e^x + e^{-x}}{2}\right), \quad L_{\log-\cosh}(x, \theta) = \sum_i \log(\cosh(\theta - x_i))$$

$\log-\cosh$ برای مقادیر کوچک شبیه $\frac{x^2}{2}$ و برای مقادیر بزرگ شبیه $\log(x)$ عمل می‌کند.
 بنابراین برای مقادیر کوچک (عملاً نزدیک به میانگین با بیکار) مانند MSE عمل می‌کند اما چون محاسبه مقادیر برای مقادیر دور از میانگین دارد، حساسیت کمتری نسبت به داده‌های پرت نشان می‌دهد.
 در حقیقت مشتق برابر هم هست اما فب e^x و e^{-x} و \ln دارد که از نظر محاسباتی سنگین‌تر از MSE است.
 Huber نیز همان طوری که تعریفش ذکر شده، برای داده‌های نزدیک $\frac{x^2}{2}$ است و عملاً محاسبه MSE را دارد اما برای داده‌های دور محاسبه‌اش دارد. بنابراین بسیار شبیه $\log-\cosh$ عمل می‌کند.
 اما فب حالت نهی دارد برای محاسبه، البته که توابع کمی هستند از نظر محاسباتی.
 در نتیجه می‌توان گفت تا حدودی مشکلات MSE و MAE را حل می‌کند. گویا نسبت به داده‌های نزدیک به میانگین شبیه MSE و نسبت به داده‌های پرت شبیه MAE عمل می‌کند.

پس در ۲ روش پیشنهادی، برای خطای زیاد مانند MAE و برای خطای کم شبیه MSE رفتار می‌کند. (تقریباً!)
 $\log-\cosh$ مشتق برابر هم هست در هر نقطه که از Huber بهتر است.
 Huber محاسبات ریاضی ساده‌تری دارد.

تحلیل کلی : با افزایش تعداد iteration ها، جمله رفتن زاینده یادگیری، خطای train کاهش یافته و به یک مقدار optimal رسیده است. خطای Validation 2 نیز کاهش یافته و به یک مقدار optimal رسیده است که البته از مقدار خطای train پلش ز است که خوب منطقی است چرا که پارامترها بر اساس بهبود train تغییر می کنند. مشاهده می شود که در هر به نمودار، late بهبود در ابتدا سریع و زیاد است و به مرور این late کم می شود. دلیل این این است که پارامترها در ابتدا مقدار خیلی پستی دارند و در هر iteration، گام های بزرگی به سمت حالت بهینه برداشته می شود که با ادامه زاینده این کار سرعت کمتری به خود می گیرد. اما اتفاق عجیب این است که پس از یک نقطه ای، خطای Validation 1 دیگر کم نشده و حتی افزایش می یابد.

توزیع متفاوت : اگر از ۲ توزیع متفاوت استفاده باشند، می توان اینکه تحلیل کرد که توزیع train با Validation 2 یکسان است و با Validation 1 متفاوت است. پارامترها در ابتدا خیلی بد هستند بنابراین احتمال دارد که حرکت به سمت بهر شدن train حتی باعث بهتر شدن Validation 1 هم بشود. اما از یک نقطه به بعد پارامترها خیلی تخصصی و دقیق روی توزیع train متمرکز می شوند و جنب این گرایش با توزیع Validation 1 سازگار نیست و باعث افزایش خطای آن می شود. یعنی میر Validation 1 جای خود را از Validation 2 می دهد برای آن.

توزیع یکسان : اگر توزیع یکی باشد، می توان به نحوه نمونه برداری (sampling) ایراد گرفت. مثلاً شاید متد Validation 1 70% بالای توزیع را تشکیل داده باشند. مثلاً sort شده بعد sample شده! یک احتمال دیگر این است که در زمان training، برای بهبود hyper-parameters به Validation 2 رجوع شده باشد و تغییرات ایجاد شده باعث بهبود خطای Validation 2 می شود اما این کار برای Validation 1 ماسارگام بوده و باعث افزایش خطای آن می شود. مثلاً اگر برای تعیین hyper-parameter به هر دو validation رجوع می شد، احتمالاً خطای هر دو کاهش می یافت اما به مقدار بهینه بزرگی نسبت به مقدار بهینه الان Validation 2 می رسید.

← یک گزینه محتمل دیگر این است که تعداد داده ها خیلی کم است. به همین دلیل اینکه دارد رفتار می کند.

$$E[L_q] = \iint |f(x) - y|^q p(x, y) dx dy$$

سؤال ۳)

در رابطه بالا، مقدار $f(x)$ به ازای هر x مستقل محاسب می شود. پس محادای تراشیم مسئله را به مسئله زیر تبدیل کنیم:

این بهمان باید کمینه شود. $\rightarrow \int |f(x) - y|^q p(y|x) dy$

$$\Rightarrow \frac{d}{df(x)} \int |f(x) - y|^q p(y|x) dy = 0$$

$$\Rightarrow \frac{d}{df(x)} \left\{ \int_{-\infty}^{f(x)} (f(x) - y)^q p(y|x) dy + \int_{f(x)}^{\infty} (y - f(x))^q p(y|x) dy \right\} = 0$$

$$\Rightarrow \int_{-\infty}^{f(x)} q(f(x) - y)^{q-1} p(y|x) dy - \int_{f(x)}^{\infty} q(y - f(x))^{q-1} p(y|x) dy = 0$$

$$\Rightarrow \int_{-\infty}^{f(x)} (f(x) - y)^{q-1} p(y|x) dy = \int_{f(x)}^{\infty} (y - f(x))^{q-1} p(y|x) dy$$

$$\int_{-\infty}^{f(x)} p(y|x) dy = \int_{f(x)}^{\infty} p(y|x) dy : q=1$$

پس مساحت سمت چپ $f(x)$ با مساحت سمت راست آن باید برابر باشند. محادای همان میانه است در توزیع پوی شده.

$q=0$: اگر دقیقاً صفر باشد، مقدار $|f(x) - y|^q$ صفر شده و املا نمی توان مشتق گرفت و $E[L_q] = 1$ می شود.

اگر $q \rightarrow \infty$ میل می کند، باز هم مقدار $|f(x) - y|^q$ بسیار نزدیک به 1 است. البته در حسابهای نزدیک $f(x)$ ، این مقدار به صفر نزدیک است. با در نظر گرفتن این در y نرمال است و می توان گفت بیشترین کاهش در نزدیک $y = f(x)$ است که بیشترین $p(y|x)$ را دارد.

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}, \quad S(x) = \frac{1}{1 + e^{-x}}$$

سوال ۴

$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + e^{-x}} \Rightarrow \forall S(x) - 1 = \frac{e^x}{e^x + e^{-x}} - 1 = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\Rightarrow \tanh(x) = \forall S(x) - 1$$

در ساختار آلفا شده رابط پارامترها برای $f(x, w^{\text{sigmoid}})$ و $f(x, w^{\text{tanh}})$ لینک شده.

$$f(x, w^{\text{sig}}) = w_0^{\text{sig}} + \sum_{j=1}^m w_j^{\text{sig}} S(x)$$

$$f(x, w^{\text{tanh}}) = w_0^{\text{tanh}} + \sum_{j=1}^m w_j^{\text{tanh}} \tanh(x) = w_0^{\text{tanh}} + \sum_{j=1}^m w_j^{\text{tanh}} (\forall S(x) - 1)$$

$$= w_0^{\text{tanh}} + \sum_{j=1}^m \{ \forall w_j^{\text{tanh}} S(x) - w_j^{\text{tanh}} \} = w_0^{\text{tanh}} + \sum_{j=1}^m w_j^{\text{tanh}} + \sum_{j=1}^m \forall w_j^{\text{tanh}} S(x)$$

$$\Rightarrow w_0^{\text{sigmoid}} = w_0^{\text{tanh}} + \sum_{j=1}^m w_j^{\text{tanh}}, \quad \text{for } j=1 \text{ to } m: w_j^{\text{sigmoid}} = \forall w_j^{\text{tanh}}$$

یا به شکل دیگر:

$$w_0^{\text{tanh}} = w_0^{\text{sigmoid}} - \frac{1}{\forall} \sum_{j=1}^m w_j^{\text{sigmoid}}, \quad \text{for } j=1 \text{ to } m: w_j^{\text{tanh}} = \frac{1}{\forall} w_j^{\text{sigmoid}}$$

Q)

$$y(x, w) = w_0 + \sum_{i=1}^D w_i x_i \rightarrow \text{our linear model.}$$

$$E_D(w) = \frac{1}{P} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 \rightarrow \text{sum square error.}$$

$$x = [x_1, x_2, \dots, x_D]^T \xrightarrow{\text{add noise}} x_{\text{noisy}} = [x_1 + \epsilon_1, x_2 + \epsilon_2, \dots, x_D + \epsilon_D]^T$$

$$E_D(w)_{\text{noisy}} = \frac{1}{P} \sum_{n=1}^N \{y(x_{\text{noisy}}, w) - t_n\}^2 = \frac{1}{P} \sum_{n=1}^N \left\{ \left[w_0 + \sum_{i=1}^D w_i (x_i + \epsilon_i) \right] - t_n \right\}^2$$

$$= \frac{1}{P} \sum_{n=1}^N \left\{ \left[w_0 + \sum_{i=1}^D w_i x_i \right] - t_n + \sum_{i=1}^D w_i \epsilon_i \right\}^2 = \frac{1}{P} \sum_{n=1}^N \left\{ y(x_n, w) - t_n + \sum_{i=1}^D w_i \epsilon_i \right\}^2$$

$$= \frac{1}{P} \sum_{n=1}^N \left\{ \underbrace{(y(x_n, w) - t_n)}_A + \underbrace{\left(\sum_{i=1}^D w_i \epsilon_i \right)}_B + \underbrace{y(x_n, w) - t_n \left(\sum_{i=1}^D w_i \epsilon_i \right)}_C \right\}^2$$

~~$$E_D(w) = \frac{1}{P} \sum_{n=1}^N \left\{ \underbrace{(y(x_n, w) - t_n)}_A + \underbrace{\left(\sum_{i=1}^D w_i \epsilon_i \right)}_B + \underbrace{y(x_n, w) - t_n \left(\sum_{i=1}^D w_i \epsilon_i \right)}_C \right\}^2$$~~

$$E_E[E_D(w)] = E_E \left[\frac{1}{P} \sum_{n=1}^N (A + B + C) \right] = \frac{1}{P} \sum_{n=1}^N E_E[A] + \frac{1}{P} \sum_{n=1}^N E_E[B] + \frac{1}{P} \sum_{n=1}^N E_E[C]$$

$$E_E[A] = E_E[(y(x_n, w) - t_n)^2] = (y(x_n, w) - t_n)^2 \quad \checkmark$$

$$E_E[B] = E_E \left[\left(\sum_{i=1}^D w_i \epsilon_i \right)^2 \right] = E_E \left[\sum_{i=1}^D \sum_{j=1}^D w_i w_j \epsilon_i \epsilon_j \right] = \sum_{i=1}^D \sum_{j=1}^D w_i w_j E_E[\epsilon_i \epsilon_j]$$

$$= \sigma^2 \sum_{i=1}^D \sum_{j=1}^D w_i w_j \delta_{ij} \rightarrow \left. \begin{array}{l} E_E[\epsilon_i \epsilon_j] = \sigma^2 \delta_{ij}, \epsilon_i \sim \mathcal{N}(0, \sigma^2) \text{ i.i.d.} \\ \delta_{ij} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases} \end{array} \right\} E_E[B] = \sigma^2 \sum_{i=1}^D w_i^2 \quad \checkmark$$

$$E_C[C] = E_C \left[\gamma \left(\sum_{i=1}^D w_i \epsilon_i \right) (y(x_n, w) - t_n) \right]$$

$$= \gamma (y(x_n, w) - t_n) E_C \left[\sum_{i=1}^D w_i \epsilon_i \right] = \gamma (y(x_n, w) - t_n) \sum_{i=1}^D w_i E_C[\epsilon_i]$$

• $\epsilon_i \sim \mathcal{N}(0, \sigma^2) \Rightarrow E_C[\epsilon_i] = 0 \Rightarrow E_C[w_i \epsilon_i] = 0$

$$\Rightarrow E_C[C] = 0$$

$$\Rightarrow E_C[E_D(w)] = \frac{1}{\gamma} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2 + \frac{\sigma^2}{\gamma} \sum_{i=1}^D w_i^2$$

$$\Rightarrow E[E_{D-\text{noisy}}^{(w)}] = E_D(w) + \frac{\sigma^2}{\gamma} \sum_i w_i^2 \quad \checkmark$$

نکته: فرض شد $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ ، همچنین

$$\sigma^2 \delta_{ij} = E(\epsilon_i \epsilon_j)$$

که داریم:

$$\delta_{ij} = \begin{cases} 1 & i=j \\ 0 & \text{o.w.} \end{cases}$$

$$L(w) = \sum_{i=1}^n f_i (y_i - w^T x_i)^2 \Rightarrow w^* = \underset{w}{\operatorname{argmin}} L(w)$$

(4 سوال)

derivative of $L(w)$ with respect to $w = \nabla_w L(w) = 0$

$$\nabla_w L(w) = -2 \sum_{i=1}^n f_i (y_i - w^T x_i) x_i^T = \sum f_i y_i x_i^T - \sum f_i w^T x_i x_i^T = 0$$

$$\Rightarrow \sum_{i=1}^n f_i y_i x_i^T = w^T \sum_{i=1}^n f_i x_i x_i^T \Rightarrow w^T = \sum_{i=1}^n f_i y_i x_i^T \left(\sum_{i=1}^n f_i x_i x_i^T \right)^{-1}$$

$$\Rightarrow w = \left(\sum_{i=1}^n f_i y_i x_i^T \left(\sum_{i=1}^n f_i x_i x_i^T \right)^{-1} \right)^T = \left(\sum_{i=1}^n f_i x_i x_i^T \right)^{-1} \left(\sum_{i=1}^n f_i y_i x_i^T \right)^T$$

$$= \left(\sum_{i=1}^n f_i x_i x_i^T \right)^{-1} \sum_{i=1}^n f_i x_i y_i^T$$

$$\sum_{i=1}^n f_i x_i x_i^T = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} F_1 & & & 0 \\ & F_2 & & \\ & & \ddots & \\ 0 & & & F_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = X F X^T$$

$$\sum_{i=1}^n f_i x_i y_i^T = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} F_1 & & & 0 \\ & F_2 & & \\ & & \ddots & \\ 0 & & & F_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = X F Y^T$$

$$\Rightarrow w^* = (X F X^T)^{-1} X F Y^T$$

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}, Y = \begin{bmatrix} y_1 & y_2 & \dots & y_n \end{bmatrix}, F = \begin{bmatrix} F_1 & 0 & \dots & 0 \\ 0 & F_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & F_n \end{bmatrix}$$

$$X = \begin{bmatrix} x_1^{(1)} & \dots & x_m^{(1)} \\ \vdots & & \vdots \\ x_1^{(n)} & \dots & x_m^{(n)} \end{bmatrix}, Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \Rightarrow w^* = (X^T F X)^{-1} X^T F Y$$

یا با این روش:

$$E_D(w) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2$$

(a)

$$\Rightarrow w^* = (X^T X)^{-1} X^T Y$$

از طریق نقطه میسر ز را برای ما می‌دهد، پس هر داده x به شکل زیر دیده می‌شود:

$$x = [x_1, x_2, \dots, x_m] \longrightarrow x = [0, 0, \dots, x_j, 0, \dots, 0]$$

$$X = \begin{bmatrix} 0 & 0 & \dots & x_j & 0 & \dots & 0 \end{bmatrix}$$

پس عملاً ماتریس X به این شکل دیده می‌شود:

$$\Rightarrow X^T X = x_j^T x_j \Rightarrow (X^T X)^{-1} = (x_j^T x_j)^{-1} = \frac{1}{x_j^T x_j}, \quad X^T Y = x_j^T Y$$

$$\Rightarrow w_j^* = \frac{x_j^T Y}{x_j^T x_j}$$

(b)

$$w^* = (X^T X)^{-1} X^T Y$$

$$X \rightarrow \text{orthogonal} \Rightarrow X^T X = \begin{bmatrix} \|x_1\|^2 & & 0 \\ & \ddots & \\ 0 & & \|x_m\|^2 \end{bmatrix}$$

$$\Rightarrow w^* = \begin{bmatrix} \frac{1}{\|x_1\|^2} & & 0 \\ & \ddots & \\ 0 & & \frac{1}{\|x_m\|^2} \end{bmatrix} X^T Y \Rightarrow w_j = \frac{x_j^T Y}{\|x_j\|^2} = \frac{x_j^T Y}{x_j^T x_j}$$

از طریق (a) ثابت کردیم که w به دست آمده از روش میسر ز، $\frac{x_j^T Y}{x_j^T x_j}$ می‌شود.

پس مسئله حل شد!

$$w^* = (X^T X)^{-1} X^T y$$

$$X = \begin{bmatrix} 1 & 0 & 0 & \dots & x_j^{(1)} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & x_j^{(n)} & 0 & \dots & 0 \end{bmatrix}$$

معمولی توان به X اضافه نگذاشت:

که شب مثل این می ماند که X کلاً ویژگی داشته باشد، پس عملاً برای X داریم:

$$X = \begin{bmatrix} 1 & x_j^{(1)} \\ \vdots & x_j^{(n)} \end{bmatrix} \Rightarrow X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_j^{(1)} & x_j^{(2)} & \dots & x_j^{(n)} \end{bmatrix}$$

$$\Rightarrow X^T X = \begin{bmatrix} n & nE[x_j] \\ nE[x_j] & \sum x_j^2 \end{bmatrix}, \text{Var}(x_j) = nE[x_j^2] - E[x_j]^2 = \frac{\sum x_j^2}{n} - E[x_j]^2$$

$$\Rightarrow \sum x_j^2 = n(\text{Var}(x_j) + E[x_j]^2)$$

$$X^T y = \begin{bmatrix} nE[y] \\ x_j^T y \end{bmatrix}, \text{Cov}(x_j, y) = \sum_i (x_j^{(i)} - E[x_j])(y^{(i)} - E[y]) \times \frac{1}{n}$$

$$\Rightarrow x_j^T y = n(\text{Cov}(x_j, y) + E[x_j]E[y])$$

$$\Rightarrow (X^T X)^{-1} X^T y = \frac{n}{n \times \text{Var}(x_j)} \begin{bmatrix} \text{Var}(x_j) + E[x_j]^2 & -E[x_j] \\ -E[x_j] & 1 \end{bmatrix} \begin{bmatrix} E[y] \\ \text{Cov}(x_j, y) + E[x_j]E[y] \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\text{Var}(x_j)E[y] - E[x_j]\text{Cov}(x_j, y)}{\text{Var}(x_j)} \\ \frac{\text{Cov}(x_j, y)}{\text{Var}(x_j)} \end{bmatrix} = \begin{bmatrix} w_0 \\ w_j \end{bmatrix} \Rightarrow w_j = \frac{\text{Cov}(x_j, y)}{\text{Var}(x_j)}$$

$$w_0 = E[y] - w_j E[x_j]$$