

به نام خدا

امیر محمد یعقوبیان زاده ۱۳۰۰ هجری قمری
تکمین ۲ یا دگرین تکمین

$$L = f(\omega_0) + \nabla f(\omega_0)^T (\omega - \omega_0) + \frac{1}{2} (\omega - \omega_0)^T H (\omega - \omega_0) \quad (1) \text{ د آ}$$

چون با GD همگرا شده ایم پس $\nabla f(\omega_0) = 0$. هدف ما این است که یک بار وزن ها موجود در ω را تغییر دهیم به طوری که L (خطا) کمترین افزایش را داشته باشد.

$$\hat{\omega} = \omega_0 - \omega_i$$

تغییرات را به آن صورت به فرم ω_i نام
که برابر $\omega_{0(i)}$ است.

$$\omega_i = \begin{bmatrix} \vdots \\ \omega_{0(i)} \\ \vdots \end{bmatrix} \leftarrow \text{در اینجا}$$

$$\Rightarrow L_i = f(\omega_0) + 0 + \frac{1}{2} (-\omega_i)^T H (-\omega_i) = f(\omega_0) + \frac{1}{2} \omega_i^T H \omega_i$$

$$L_i = f(\omega_0) + \frac{1}{2} |\omega_{0(i)}|^2 H_{ii}$$

از بین وزن ها موجود باید وزنی را معرفی کنیم که مقدار $|\omega_{0(i)}|^2 H_{ii}$ برای آن کمینه باشد.

$$H = I \rightarrow H_{ii} = 1 \rightarrow$$

(ب) باید وزن ها را حذف کنیم که کمترین نرم L_2 داشته باشند.

$$b^* = \underset{b}{\operatorname{argmin}} \frac{1}{N} \|y - Xb\|^2 \xrightarrow{\frac{\partial}{\partial b}} \frac{\partial}{\partial b} (y^T y + b^T X^T X b - 2y^T X b) = 0 \quad \pi, \textcircled{P}$$

$$\Rightarrow 2X^T X b = 2X^T y \Rightarrow b^* = (X^T X)^{-1} X^T y$$

$$A^T = X(X^T X)^{-1} X^T = A \rightarrow A \text{ مقعر است.} \quad (\text{ب})$$

$$A^T A = A^T = A \Rightarrow Av = \lambda v \Rightarrow \begin{cases} A^T v = \lambda Av = \lambda^2 v \\ A^T v = Av = \lambda v \end{cases} \Rightarrow \lambda^2 = \lambda \Rightarrow \boxed{\lambda = 1 \text{ یا } 0}$$

با نوشتن تجزیه SVD برای $X = U \Sigma V^T$ و بدست آوردن A می بینیم که ماتریس A دارای رتبه یک ماتریس قطری با ± 1 است.

در $N-d$ درجه آزادی است.

$$E\left[\frac{1}{N} \left\| \overbrace{(X(X^T X)^{-1} X^T - I)}^A \epsilon \right\|_F^2\right] = \frac{1}{N} E\left[\sum_{i=1}^{N-d} \epsilon_i^2\right] = \frac{1}{N} \times \sum_{i=1}^{N-d} E[\epsilon_i^2]$$

ماتریس قطری با $N-d$ درجه 1 و d درجه صفر در قطری

$$= \frac{1}{N} (N-d) \sigma^2$$

(ج) هرچه مقدار d به N نزدیک تر شود مقدار $N-d$ به سمت صفر می رود و خطای همبسته شده رفته رفته کاهش می یابد.

$$\lim_{d \rightarrow N} \left(\frac{N-d}{N}\right) \sigma^2 = 0$$

$$\hat{\beta} = \arg \min_{\beta} \|Y - X\beta\|_F^2 \rightarrow \frac{\partial}{\partial \beta} (Y^T Y - 2\beta^T X^T Y + \beta^T X^T X \beta) = 0 \quad (1) \quad (2)$$

$$-2X^T Y + 2X^T X \beta = 0 \Rightarrow \beta = (X^T X)^{-1} X^T Y$$

$$\hat{\beta} = \arg \min_{\beta} \|Y - X\beta\|_F^2 + \lambda \|\beta\|_F^2 \rightarrow \frac{\partial}{\partial \beta} (\|Y - X\beta\|_F^2 + \lambda \|\beta\|_F^2) = 0 \quad (ب)$$

$$\Rightarrow -2X^T Y + 2X^T X \beta + 2\lambda \beta = 0 \Rightarrow \beta = (X^T X + \lambda I)^{-1} X^T Y$$

$$X = \Sigma^{-1} X F \quad (ج) \quad F \leftarrow \sum_{N \times N} X_{N \times K} = X_{N \times K} F_{K \times K}$$

F ماتریس غیرتکین است و بنابراین در آن وارون است.

$$\hat{\beta} = (X^T X)^{-1} X^T Y = (\underbrace{F^T X^T \Sigma^{-1} X}_{X^T})^{-1} \underbrace{F^T X^T \Sigma^{-1} Y}_{X^T} = (X^T \Sigma^{-1} X)^{-1} F^T F^T X^T \Sigma^{-1} Y$$

$$= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y$$

$$(X^T X)^{-1} X^T Y = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y \Rightarrow X^T Y = X^T X (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y \quad (\text{توجه! اگر})$$

$$\xrightarrow[\text{چون باید به هم برابر باشد}]{\text{برقرار باشد}} X^T = X^T X (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} \xrightarrow[\text{transpose}]{X \Sigma} \Sigma X = X (X^T \Sigma^{-1} X)^{-1} X^T X \equiv F$$

برای آنکه سعی بر قرار دهی باید $\Sigma X = XF$ باشد. طبق تعریف، F یک ماتریس مکتوب پذیر و غیرتکین است.

(د) اگر $X \in \mathbb{R}^{N \times d}$ باشد در نمونه ها d ویژگی باشند با افتاده بودن λ نمونه به صورت $(y=0, x=e_1, \dots, x=e_d)$

نمونه به صورت $(y=0, x=e_1, \dots, x=e_d)$ و λ نمونه به صورت $(y=0, x=e_1, \dots, x=e_d)$ به دیاست و اعمال یک رگرسیون L_1 به تابع

هدف می توان تابع خطای نشان داده شده را نوشت:

$$L = \|y - X\beta\|^2 + \lambda_1 \|\beta\|_1 = \|y_{old} - X_{old}\beta\|^2 + \lambda_1 \sum_{i=1}^d (0 - e_i \beta)^2$$

$$+ \lambda_2 \|\beta\|_1 = \|y_{old} - X_{old}\beta\|^2 + \lambda_1 \sum_{i=1}^d \beta_i^2 + \lambda_2 \|\beta\|_1$$

$$= \|y_{old} - X_{old}\beta\|^2 + \lambda_1 \|\beta\|_2^2 + \lambda_2 \|\beta\|_1$$

$$\frac{\partial \mathcal{L}}{\partial w_i} = -2(y_d - \sum_{k=1}^n \delta_k w_k x_k) \times \delta_i x_i$$

(۱) (۲)

$$\Rightarrow E\left[\frac{\partial \mathcal{L}}{\partial w_i}\right] = E\left[-2y_d x_i \delta_i + 2x_i \delta_i \sum_{k=1}^n \delta_k w_k x_k\right] = -2y_d x_i + 2x_i \sum_{\substack{k=1 \\ k \neq i}}^n w_k x_k$$

$$+ 2x_i^2 w_i (\sigma^2 + 1) = 2x_i^2 w_i \sigma^2 - 2y_d x_i + 2x_i \sum_{k=1}^n w_k x_k \quad (I)$$

(ب) ابتدا بردارین را برای حالتی که این نوع شبکه ریزایون وجود ندارد محاسبه می کنیم:

$$\frac{\partial \hat{\mathcal{L}}}{\partial w_i} = -2(y_d - \sum_{k=1}^n w_k x_k) \times x_i = -2y_d x_i + 2x_i \sum_{k=1}^n w_k x_k \quad (II)$$

با مقایسه روابط (I) و (II) متوجه می شویم که عبارت $2x_i^2 w_i \sigma^2$ در بردارین حالت dropout اضافه شده است.

آپری تابع خطای $\hat{\mathcal{L}}$ مقدار $\sigma^2 \|x\|^2$ اضافه شود هنگام $\frac{\partial}{\partial w_i}$ گرفتن مقدار $2x_i^2 w_i \sigma^2$ تولید خواهد شد. بنابراین چنین ریزایونی به تابع خطای اضافه شده است که اندکی به ریزایون L_2 تفاوت و در آن مقادیر بزرگتر و ابتدا باعث scale درایه ها وزن می شود و سپس ریزایون L_2 بزرگتر برای جدید اعمال می گردد.

$$w_t = w_{t-1} - \epsilon \nabla_w (w^T H w) \Big|_{w=w_{t-1}} = w_{t-1} - \epsilon H w_{t-1} = (I - \epsilon H) w_{t-1}$$

(۳) (۴)

$$w_t = (I - \epsilon H) w_{t-1} = (I - \epsilon H)^2 w_{t-2} = \dots = (I - \epsilon H)^t w_0$$

(۵)

$$I - \epsilon H = Q Q^T - \epsilon Q \Lambda Q^T = Q (I - \epsilon \Lambda) Q^T$$

(۶)

$$w_t = Q (I - \epsilon \Lambda)^t Q^T w_0$$

$$w \rightarrow (1 - \epsilon \lambda_i)^t$$

اگر $|1 - 2\epsilon\lambda_i| < 1$ باشد این الگوریتم همگرا خواهد شد.

$$\Rightarrow -1 < 1 - 2\epsilon\lambda_i < 1 \quad \forall \lambda_i \Rightarrow \begin{cases} \epsilon\lambda_i \geq 0 \\ \epsilon\lambda_i \leq 1 \end{cases}, \forall \lambda_i$$

$$\epsilon \leq \frac{1}{\lambda_{\max}}, \quad \lambda_i \geq 0$$

(د)

$$\omega_{t+1} = \omega_t - (\nabla H)^T \times \nabla H \omega_t = \omega_t - H^{-1} H \omega_t = 0$$

با این روش در یک گام به جواب بهینه $\omega = 0$ می‌رسیم.

(۵) تابع خطا در شبکه‌های عصبی پیچیده تر است و می‌تواند مایه سردی و یکنواختی کردن آن در یک گام بهینه‌تر است.

⑥ در این شبکه $W=0$ را یاد می‌گیرد و اصولاً مفهوم خامی را آموزش نمی‌بیند. اگر $W=0$ باشد $h_i = h_r = \tanh(b)$ و $J = 0 + 0 = 0$ خواهد شد که کمترین مقدار تابع هدف است.

(ب)

$$|h_i - h_r|_r^2 = \sum_{k=1}^m \left(\tanh \left[\sum_{i=1}^n w_{ki} x_{1i} + b_k \right] - \tanh \left[\sum_{i=1}^n w_{ki} x_{ri} + b_k \right] \right)^2$$

$$\frac{\partial |h_i - h_r|_r^2}{\partial w_{e,t}} = 2 \left(\sqrt{} \right) \times \left[(1 - \tanh^2 \left[\sum_{i=1}^n w_{ei} x_{1i} + b_e \right]) \times x_{1t} - x_{rt} \times \right]$$

$$(1 - \tanh^2 \left[\sum_{i=1}^n w_{ei} x_{ri} + b_e \right]) = 2 (h_{1e} - h_{re}) \left[(1 - h_{1e}^2) x_{1t} - (1 - h_{re}^2) x_{rt} \right]$$

$$= 2 \left[(h_1 - h_r) \circ \left[(1 - h_1 \circ h_1) x_1^T - (1 - h_r \circ h_r) x_r^T \right] \right]_{e,t}$$

\rightarrow اپراتور ضرب درایه به درایه (ضرب هادامارد)

$$\Rightarrow \left\{ \frac{\partial J}{\partial w} = 2 (h_1 - h_r) \circ \left[(1 - h_1 \circ h_1) x_1^T - (1 - h_r \circ h_r) x_r^T \right] + 2W \right.$$

$$\left. W_{t+1} = W_t - \epsilon \frac{\partial J}{\partial w} \right|_{w=W_t}$$

$$\frac{\partial |h_i - h_r|_r}{\partial b_j} = r(h_{ij} - h_{rj})[(1 - h_{ij}^r) - (1 - h_{rj}^r)]$$

$$\Rightarrow \begin{cases} \frac{\partial \mathcal{J}}{\partial b} = r(h_i - h_r) \circ (h_r \circ h_r - h_i \circ h_i) \\ b_{t+1} = b_t - \epsilon \frac{\partial \mathcal{J}}{\partial b} \Big|_{b=b_t} \end{cases}$$