wil 1-1) and Tule Jesent of Signary ide lina with the lesent of Till and the line (Till) and the contract of t

رسره است :

رسده است : (سره است : الله عال الله ع

 $w^* = \begin{bmatrix} w_1^* \\ \vdots \\ w_n^* \end{bmatrix} \longrightarrow w = \begin{bmatrix} w_1^* \\ \vdots \\ w_n^* \end{bmatrix}$

ار (ساب لنم ،

 $J(\widetilde{\omega}) = J(\omega^*) + \frac{1}{2}(\widetilde{\omega} - \omega^*)^T H(\widetilde{\omega} - \omega^*)$ سـ الله برداری است له نقط درای ۱۲ آن فاصنر دیرایر پی است . برنسی اله ۲ به مورت زیر

 $H = [h_1, \dots, h_n] \longrightarrow H(\widetilde{w} - w^*) = -w_K^* h_K \Rightarrow (\widetilde{w} - w^*)^T (-w_K^* h_K)$

 $= w_{K}^{*2} H_{KK} \Rightarrow J(\overline{w}) = J(w^{*}) + H_{KK} w_{K}^{*2}$

يسى الر نما را هذف للم م طبع هزيم داندازه W*Hkk زماد ى نسود. در نسم جايد درايداى را انتفاب كنام كه مقدار Wk على كمترون والله. توج نشود كه فرض وكينم اللويس ب نقط بمن على minimum لماء رسده است ر مَعَي ﴿ ٢٨ است بني هزن بعداز حدث * ٨٨ انزاس ي مابد .

 $J(\vec{w}) = J(w^*) + w_k^*$ ب) الد الم هاني بانسر، ا= All براى هم ، ها درنسم خراصم الست، رَسْعِ بَالِهِ درانِهِ اي راخزت كُلْسُم له كُمترين الله وا داور، بران ابن كار آلد بخواصم ملل ١٠ دران خزت كُلْم، x درایه ای که کدهکترین انزان را داند را خزن ی کنیم.

$$X = \begin{bmatrix} -x_i^T - x_i \\ x_i \in \mathbb{R}^d \end{bmatrix}, \quad x_i \in \mathbb{R}^d$$

 $y_i = b^{T_{\alpha_i}} + \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(0, \delta^2)$

7) we hall ke jung sour

$$C = \|y - \hat{y}\|_{2}^{2} = \|Xb - y\|_{2}^{2} \Rightarrow \hat{b} = \underset{b}{\text{argmin}} \|Xb - y\|_{2}^{2}$$

e stig Law X mit pseudo inverse : Lim in the diese

$$\hat{b} = X y$$

 $\hat{b} = (X \bar{X}) X \bar{Y}$

طل جون X رتبکال نشرنی است (لم < N) رتب ا

دلل عادس بند بون XX ابن است اله تفان سطی Xx برابر است و فقال سطی X مم دارای زمل له است. جن XTX ما ترس له xله است و نفای ساری آن دارای له یاج است یعی رتب کامل است ر عارس بدر ی باند.

$$A^{T} = X (X^{T} X) X^{T}$$

 $A^2 = X(XX) X X (XX) X = X(XX) X$

رسا له X(xTx) X = A السر

. wind A = A caid, with a

ر م متعارف العب ، حون A متعارف العب ، حون A متعارف العب ، محتوان تجزیه و تعاویر دیره انجام دارد، $A = \mathcal{Q} \Lambda \mathcal{Q}^{\mathsf{T}} \longrightarrow A^2 = \mathcal{Q} \Lambda^2 \mathcal{Q}^{\mathsf{T}}$

زسم عامر درو م م توان مر عادر ويو A متندر الر آن ها الم نشان دهم ؛

$$\lambda_i = \lambda_i^2 \longrightarrow \lambda_i (\lambda_{i-1}) = 0 \longrightarrow \begin{cases} \lambda_i = 0 \\ \lambda_i = 1 \end{cases}$$

 $X = X (XX)^{-1}X^{-1}X$ والمد برنال کلید دم ده $X = X (XX)^{-1}X^{-1}$

ون (XTX) رَنْهِ کال الله علی لا الله ای موجر دنست المه و لاز XTX) بانسد هدون ون X رته کال الله ای موجر دنست که و و لا کانسد درنسی میا ۵= x X الله . مرحود نست که ۵= و کالله کانسد درنسی میا ۵= x X الله .

· > 1s

$$IE\left[\frac{1}{N}\|(I-A)y\|_{2}^{2}\right] = \frac{1}{N}IE\left[\|(J-A)y\|_{2}^{2}\right]$$

= \frac{1}{N} \fra

 $\Rightarrow \frac{1}{N} \mathbb{E} \left[\sum_{i=1}^{N-d} (\overline{q_i}, \overline{y})^2 \right] = \frac{1}{N} \mathbb{E} \left[\sum_{i=1}^{N-d} (\overline{y_i}, \overline{y})^2 \right]$

$$\Rightarrow \frac{1}{N} \left[\sum_{i=1}^{N-d} (7, y)^{2} \right] =$$

$$\frac{1}{N} \mathbb{E} \left[\sum_{i=1}^{N-d} (y^{T} q_{i})^{2} \right]$$

: جردار دیم و شاکر عمار ویره ل ر I_A است یعی ا

$$(I - A) q_{i} = q_{i} \longrightarrow y^{T}q_{i} = q_{i}^{T}y = q_{i}^{T}(I - A)y = q_{i}^{T}(y - Ay) = q_{i}^{T}\varepsilon$$
where $\varepsilon = \begin{bmatrix} \varepsilon_{i} \\ \vdots \\ \varepsilon_{n} \end{bmatrix} \Rightarrow = \sum_{i=1}^{r} [(q_{i}^{T}\varepsilon)^{2}] = E[(f_{i}^{T}\varepsilon)^{2}] = E[(f_{i}^{T}\varepsilon)^{2}] = \sum_{j=1}^{r} q_{ij}^{2} = \varepsilon^{2}$

$$= \varepsilon^{2} \int_{0}^{1} q_{ij}^{2} = \varepsilon^{2}$$

 $\frac{1}{N} |E| || (I-A)y||_{2}^{2} = \frac{1}{N} (N-d) z^{2} = \frac{N-d}{N} z^{2}$

ع) م تعجم البطال كه فالا برست آسر برصوت رابط له- N مرجو داست . الر N مل فردك ماشور ما مشل لم ی شود. رسی ما به example علی تعداد علی از تعل می این از تعداد training

· verfitting

$$\frac{\partial J}{\partial w} = 25, x; \left(\sum_{K=1}^{n} S_{K} w_{K} x_{K} - J_{cl} \right) = -25, x; J_{cl} + 25, x; \sum_{K=1}^{n} S_{K} w_{K} x_{K} \right)^{2}, \quad S_{E} \sim N(1, J^{2})$$

$$\frac{\partial J}{\partial w} = 25, x; \left(\sum_{K=1}^{n} S_{K} w_{K} x_{K} - J_{cl} \right) = -25, x; J_{cl} + 25, x; \sum_{K=1}^{n} S_{K} w_{K} x_{K} \right)$$

$$= -25, x; J_{cl} + 25, z^{2}, w; + 25, x; \sum_{K=1}^{n} S_{K} w_{K} x_{K}$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; z^{2}, w; IE \left[S_{cl} \right] = IE \left[S_{cl} \right] + 2x; w; IE \left[S_{cl} \right] = 1 + 2x; w; x_{K} + 2x; w;$$

$$IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; \sum_{K=1}^{n} w_{K} x_{K} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; \sum_{K=1}^{n} w_{K} x_{K} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; \sum_{K=1}^{n} w_{K} x_{K} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; \sum_{K=1}^{n} w_{K} x_{K} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; \sum_{K=1}^{n} w_{K} x_{K} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; \sum_{K=1}^{n} w_{K} x_{K} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w; Z^{2} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl} + 2x; w;$$

$$\Rightarrow IE \left[\frac{2J}{\partial w} \right] = -2x; J_{cl}$$

 $\nabla J = 2H\omega \longrightarrow \omega^{t} = \omega^{t-1} - \varepsilon \nabla J(\omega^{t-1}) \Rightarrow \omega^{t} = \omega^{t-1} \varepsilon ZH\omega^{t-1}$ $\Rightarrow \omega^{t} = \omega^{t-1} - 2\varepsilon H\omega^{t-1}$

 $w^{t} = (I - 2\varepsilon H) w^{t-1} \longrightarrow w^{l} = (I - 2\varepsilon H) w^{0} \longrightarrow w^{t} = (I - 2\varepsilon H)^{t} w^{0}$ $\sim w^{t} = (I - 2\varepsilon H)^{t} w^{0} \longrightarrow w^{t} = (I - 2\varepsilon H)^{t} w^{0} \longrightarrow w^{t} = (I - 2\varepsilon H)^{t} w^{0}$ $\sim w^{t} = (I - 2\varepsilon H)^{t} w^{0} \longrightarrow w^{t} \longrightarrow w^{t}$

 $J(w) = J(w^{t-1}) + PJ(w^{t-1})^{T}(w-w^{t-1}) + \frac{1}{2}(w-w^{t-1})^{T}\tilde{H}(w-w^{t-1})$ $\nabla J(w) = \nabla J(w^{t-1}) + \tilde{H}(w-w^{t-1}) = 0 \rightarrow \tilde{H}w - \tilde{H}w^{t-1} = -PJ(w^{t-1})$ $\Rightarrow \tilde{H}w - \tilde{H}w^{t-1} = -\nabla J(w^{t-1}) \rightarrow w = w^{t-1} - \tilde{H}\nabla J(w^{t-1})$ $\nabla J(w^{t-1}) = 2\tilde{H}w^{t-1}, \quad \tilde{H} = 2\tilde{H}$

 $w^{t} = w^{t-1} = 2H$ $w^{t} = w^{t-1} = 0$ $w^{t} = w^{t-1} = 0$

آلد از نقله دلخوه شرع کنم:

و الله از نقله دلخوه شرع کنم:

مشل ایسا و کند. جاتوج به اینله قامع هزن که قام می الله مسل ایسا و کند. جاتوج به اینله قامع هزن که قام به مسل ایسا و کند. جاتوج به اینله قامع هزن که قام به مسل مرتب و قامع هان خود قامع الله ترتب مسل مرتب و قامع هان خود قامع الله ترتب آلله رسم آلله رس

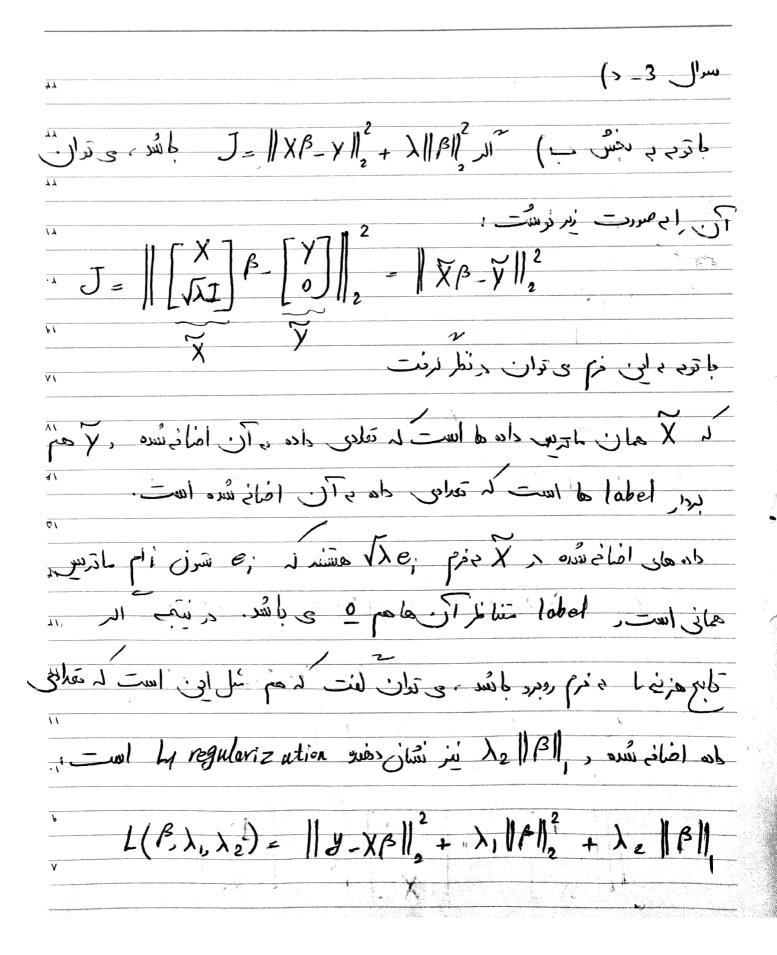
ه) اللورت نوتن نیاز ند کامی ماترس Hessian , معلوس کردن آن داد که معلوس کرفتی از ماترس از مردی آن داد که معلوس کرفتی از ماترس از مردی مردی کامید، درنسج هزن سامیلی را بر مورتی که ۱۱ بررب باشد، بسیار جالای برد.

 $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \| Y - \hat{Y} \|_{2}^{2} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ $\lim_{\beta \to \infty} \frac{1}{\beta} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ $\lim_{\beta \to \infty} \frac{1}{\beta} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ $\lim_{\beta \to \infty} \frac{1}{\beta} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ $\lim_{\beta \to \infty} \frac{1}{\beta} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ $\lim_{\beta \to \infty} \frac{1}{\beta} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ $\lim_{\beta \to \infty} \frac{1}{\beta} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ $\lim_{\beta \to \infty} \frac{1}{\beta} = \underset{\beta}{\operatorname{argmin}} \| X \beta - Y \|_{2}^{2}$ است) خوصم داست، $\hat{\beta} = (XX)XY$

 $J = \| X \beta - Y \|_{2}^{2} + \lambda \| \beta \|_{2}^{2}$

 $J = \left\| \begin{bmatrix} X \\ \sqrt{\lambda} I \end{bmatrix} \beta - \begin{bmatrix} Y \\ 0 \end{bmatrix} \right\|,$

ان عارت رای مان به صورت زیر این نوشت ، با العنفاد از حال ربعات ، جاب بهته برای می برادر خاهر بود با ا $\hat{\beta} = \left(\begin{bmatrix} x^T & \sqrt{1} \end{bmatrix} \begin{bmatrix} x \\ \sqrt{1} \end{bmatrix} \right) \begin{bmatrix} x^T & \sqrt{1} \end{bmatrix} \begin{bmatrix} y \\ 0 \end{bmatrix} = \left(x^T x + \lambda \right) \begin{bmatrix} x^T y \\ x^T y \end{bmatrix}$



$$X_{1}, X_{2} \in \mathbb{R}^{R} \quad h_{1} = \tanh\left(WX_{1} + b\right) \quad h_{2} = \tanh\left(WX_{2} + b\right)$$

$$J = \left\|h_{1} - h_{2}\right\|_{2}^{2} + \left\|W\right\|_{F}^{2}$$

$$\frac{\partial J}{\partial b} = \frac{\partial \left\|h_{1} - h_{2}\right\|_{2}^{2}}{\partial b} = \frac{\partial \left\|h_{1} - h_{2}\right\|_{2}^{2}}{\partial \left(h_{1} - h_{2}\right)} \cdot \frac{\partial \left(h_{1} - h_{2}\right)}{\partial b} = 2\left(h_{1} - h_{2}\right) \cdot \frac{\partial \left(h_{1} - h_{2}\right)}{\partial b}$$

$$\frac{\partial \left(h_{1} - h_{2}\right)}{\partial b} = \frac{\partial h_{1}}{\partial b} - \frac{\partial h_{2}}{\partial b} \Rightarrow \frac{\partial h_{1}}{\partial b} = \frac{\partial \tanh\left(WX_{1} + b\right)}{\partial b} = \frac{\partial \tanh\left(WX_{1} + b\right)}{\partial b} \cdot \frac{\partial W_{1} + b}{\partial b}$$

$$\frac{\partial \lim_{N \to \infty} \left(\tanh\left(WX_{1} + b\right)\right)}{\partial h} = \frac{\partial \lim_{N \to \infty} \left(-\frac{N}{N}\right)}{\partial h} = \frac{\partial \lim_{N \to \infty} \left(-\frac{N}{$$