

أيضاً كل mini-batch يوزع على mini-batch size / الفيديو

loss function (الخسارة) كلها هي المجموع الكلي لـ distribution

أيضاً كل mini-batch يوزع على mini-batch size / الفيديو

• mini-batch size / الفيديو يوزع على mini-batch size / الفيديو

• mini-batch size / الفيديو يوزع على mini-batch size / الفيديو

• mini-batch size / الفيديو يوزع على mini-batch size / الفيديو

• mini-batch size / الفيديو يوزع على mini-batch size / الفيديو

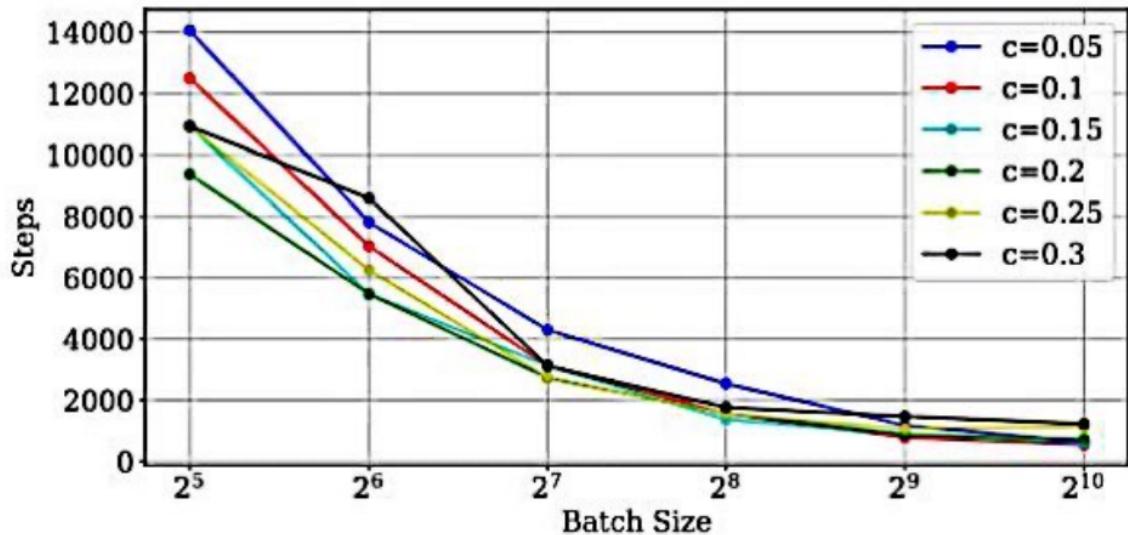


Figure 1: Number of steps for Algorithm 1 versus batch size needed to train ResNet-34 on CIFAR-10

Scanned with CamScanner

$$\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$$

↳ Batch Normalization will be explained later.

$$\bar{S}_B = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_B)^2 \quad \leftarrow \text{m.b. variance}$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \leftarrow \text{normalize}$$

$$g_i = \gamma \hat{g}_i + \beta = BN_{\gamma, \beta}(g_i) \sim \text{scale and shift}$$

size of mini-batch $\leftarrow m$ since learnable θ has $C \times V \times B$

وتحت BN Water: جمعيات اجتماعية محلية

use of a VgB learnable classifier. BN layer fixed

...invis (share) will go to;

50-148 *Sigmodon leucotis* mini bathurstensis baird

→ $\text{P}(\text{Randomness}) = \frac{1}{2^{\text{number of bits}}}$ (Randomness) \neq

• ~~multiple~~ Regularization

Numerical instability during training with gradient (stability issues with optimization) (E)

(1) Overfitting (جفون عدو و جفون بارجع بغير فرق بين تابع هادئ و تابع هادئ)

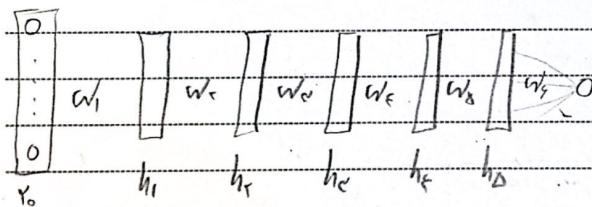
Gradient vanishing: saturation (أسباب) extreme

Gradient exploding: extreme gradients (أسباب)

Shift in performance: minimization (أسباب)

Initial values (أسباب) and initialization (أسباب)

Random HE, Xavier with weight initialization



input
layer

number of learnable parameters

عدد المعلمات

$$N_{w_1} + N_{w_2} + \dots + N_{w_4} + N_{b_1} + \dots + N_{b_4}$$

$$= 10 \times 10 + 5 \times (10 \times 10) + 10 \times 1 + 5 \times 10 + 1$$

$$= 991$$

$\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$ (output of softmax) \rightarrow $\hat{y}_i = \text{softmax}(w_i u + b_i)$

$$\hat{y} = \text{softmax}(w_s u + b_s) = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\rightarrow \hat{y}_1 + \hat{y}_2 = 1 \rightarrow \hat{y}_1 \geq \hat{y}_2 \equiv \hat{y}_1 \geq 1 - \hat{y}_2 \rightarrow \hat{y}_1 \geq 0.5 = 1$$

$$\hat{y} = \sigma(w_s u + b_s) = \frac{1}{1 + e^{-(w_s u + b_s)}} \quad : \text{Sigmoid}$$

$\rightarrow \begin{cases} \text{if } \hat{y} \leq 0.5 \rightarrow \text{label 0} \\ \text{if } \hat{y} > 0.5 \rightarrow \text{label 1} \end{cases}$

$$\hat{y}_i = \frac{e^{(w_i u + b_i)}}{1 + e^{(w_i u + b_i)} + e^{(w_s u + b_s)}} \quad : \text{Probability}$$

$$= \frac{1}{1 + e^{((w_{sr} - w_i)u + (b_{sr} - b_i))}}$$

predicted label \rightarrow 1 if $\hat{y}_i \geq 0.5$ else 0

$$w_d = w_s - w_i \quad : \text{bias adjustment}$$

$$b_d = b_s - b_i$$

W_{jk} $\leftarrow \text{class}_k \text{ probability}$

$n \times n$ $n \times 1$

w_{jk} \downarrow

w_{jk}

\rightarrow $\text{class}_k \text{ probability}$

delta_{jk}

$$w_j = \begin{bmatrix} w_{j1} \\ w_{j2} \end{bmatrix}$$

\rightarrow $\text{class}_j \text{ weight}$

$(\text{delta}_{jk}) \rightarrow \text{class}_j \text{ weight}$

Training \rightarrow fit train \rightarrow subset \rightarrow subset \rightarrow subset \rightarrow subset \rightarrow subset

\rightarrow delta error \rightarrow subset \rightarrow subset \rightarrow subset \rightarrow subset

\rightarrow subset \rightarrow subset \rightarrow subset \rightarrow subset \rightarrow subset

\rightarrow subset \rightarrow subset \rightarrow subset \rightarrow subset \rightarrow subset

$$y_{\text{com}}^k(u) = \sum_{n=1}^N w_{nk} u_n$$

combination weight depends on class k.

ب) الخطوة الخامسة: رابط (x,y) بـ (x',y')

لهم تعلم بـ bounding box الخطوة الرابعة

لهم transform الخطوة الخامسة bounding box الخطوة الرابعة

$y = stdScale \times$

لهم rotate الخطوة الخامسة bounding box الخطوة الرابعة

لهم translate الخطوة الخامسة bounding box الخطوة الرابعة

Transformation horizontal (horizontal shearing) الخطوة الخامسة

decomposition horizontal vertical translational rotation

vertical transformation vertical vertical distance

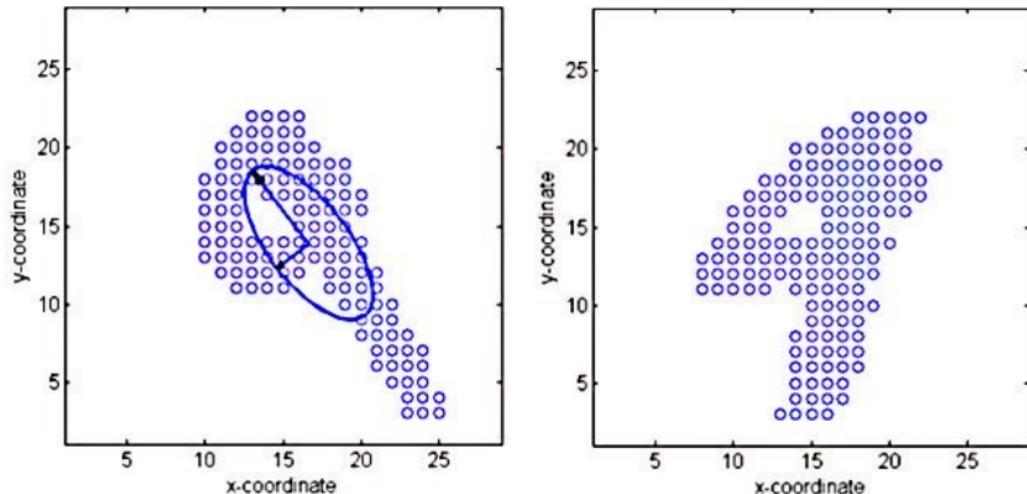


Figure 1. (Left panel) x-, y- coordinates (circles) of the original image together with the eigenvectors scaled by the corresponding eigenvalues. (Right panel) x-, y-coordinates (circles) of the deslanted image after horizontal shearing.

Scanned with CamScanner

Subject :

Date : ..

Subject: الجغرافيا

البيانات التي تم تدريب الموديل على称作 training data (s)

introduction to the new science of prefaces (with a note)

grayscale ratios (between two images) are uncorrelated), (Grayscale)

خواص ملکیت این دسته از اسناد ممکن است اینجا در متن اسناد مذکور شوند.

مختارات (ابدأ) (ابدأ) (ابدأ) (ابدأ) (ابدأ) (ابدأ) (ابدأ) (ابدأ)

سؤال ٣ Sigmoid activation function of hidden layer

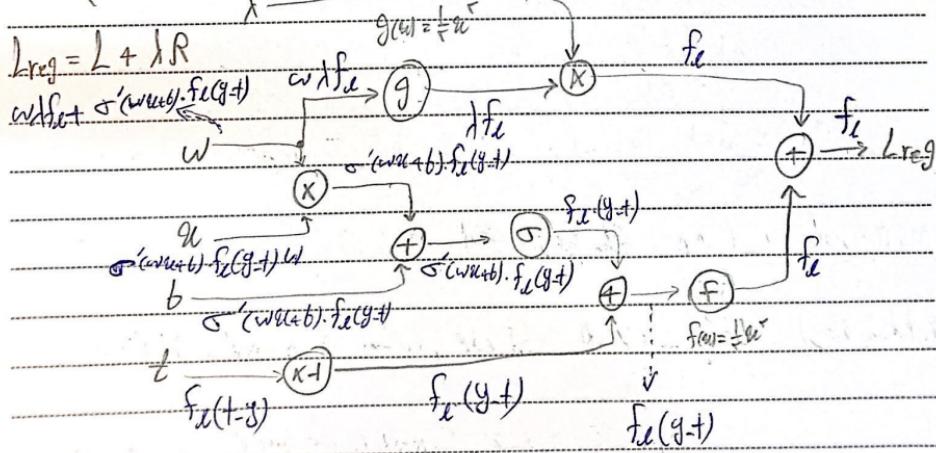
$$Z = wX + b$$

$$y = \sigma(Z)$$

$$L(a_i; w, b) = \frac{1}{T} (\sigma(a_i; w, b) - t)^2 + \frac{\lambda}{T} w^2 \quad (7)$$

$$L = \frac{1}{T} (y - t)^2$$

$$R = \frac{1}{T} w^2$$



$$\frac{\partial L_{reg}}{\partial w} = \frac{1}{\gamma} \vec{w} \quad (f_w=1) \Rightarrow \text{جبری مجموعه ای از داده ها}$$

$$\frac{\partial L_{reg}}{\partial w} = \omega_1 + \sigma'(w_0 + b)(y_1)$$

$$\frac{\partial L_{reg}}{\partial w} = \sigma'(w_0 + b)(y_1) w$$

$$\frac{\partial L_{reg}}{\partial b} = \sigma'(w_0 + b)(y_1)$$

$$\frac{\partial L_{reg}}{\partial t} = -y_1$$

$$\sigma'(w) = \frac{d}{dw} \left(\frac{1}{1+e^{-w}} \right) = \sigma(w)(1 - \sigma(w))$$

$(w^T x)$ \rightarrow $w_0 + \sum w_i x_i$ linear combination of x_i 's \rightarrow $w_0 + \sum w_i x_i + b$ \rightarrow $w^T x + b$

$\sigma(w)$ \rightarrow Activation function \rightarrow Sigmoid \rightarrow $\frac{1}{1+e^{-w}}$

ابعد آن خاصیت رفت و پنهان \rightarrow $\sigma(w)$ خود خود را درست نماید \rightarrow $\sigma(w) \approx 0.5$

$\text{Gib} \rightarrow 12 \text{ Reg. with regularization}$ \rightarrow stop gradient

\rightarrow $\nabla w = -\nabla w$ \rightarrow ∇w

\rightarrow ∇w \rightarrow ∇w \rightarrow ∇w \rightarrow ∇w \rightarrow ∇w \rightarrow ∇w

initialization of weights / initialize weights

initialization of weights / initialize weights: symmetry problem

feature scaling / normalization of input features

gradient descent / gradient descent

sigmoid function / value of sigmoid function

random initialization of weights

optional regularization / parameter space

(overfitting) / early stopping / convergence

$$\left. \begin{array}{l} \lambda = 0,1 \\ \text{learning rate } \alpha = 0,1 \\ w_0 = 1 \\ b_0 = 1 \\ \text{iteration } t \end{array} \right\} : \text{MLP}, \text{LR}$$

$$t=0 \quad \text{from G.D., we know: } w_{i+1} = w_i - \alpha \frac{\partial L_{\text{reg}}}{\partial w}$$

$$\frac{\partial L_{\text{reg}}}{\partial w} = \lambda w + \eta(y-t) \sigma'(w_0 w_1 b) = 0,1 + \eta(0 - (0,1) - 0) \sigma'(0,1)$$

$$= 0,1 + \eta(\sigma(w_1) \sigma'(w_1)) [1 - \sigma(w_1)]$$

• MICRO

Subject :

Date :

$$\frac{\partial L_{reg}}{\partial b} = (\gamma \cdot t) \sigma'(w_i^T a_i b) = \sigma(u_i + 1) \sigma(u_i + 1) [1 - \sigma(u_i + 1)]$$

$$w_i = w_0 - \alpha \frac{\partial L_{reg}}{\partial w_i} = 1 - \alpha / \times \frac{\partial L_{reg}}{\partial w_i} = 1 - \alpha / [0, 1 + \alpha \sigma(u_i + 1) (1 - \sigma(u_i + 1))]$$

$$b_i = b_0 - \alpha \frac{\partial L_{reg}}{\partial b} = 1 - \alpha / \times \frac{\partial L_{reg}}{\partial b} = 1 - \alpha / [\sigma(u_i + 1) (1 - \sigma(u_i + 1))]$$

Scanned with CamScanner

Subject:

Date:

$$\frac{\partial L_{reg}}{\partial b} = (\gamma \cdot t) \sigma'(w_i^T a_i + b) = \sigma(u_{i+1}) \sigma(u_{i+1}) [1 - \sigma(u_{i+1})]$$

$$\rightarrow \{ w_i = w_0 - \alpha \frac{\partial L_{reg}}{\partial w_i} = 1 - \sigma_i / \times \frac{\partial L_{reg}}{\partial w_i} = 1 - \sigma_i [1 + \alpha \sigma(u_{i+1}) (1 - \sigma(u_{i+1}))] \}$$

$$\{ b_i = b_0 - \alpha \frac{\partial L_{reg}}{\partial b} = 1 - \sigma_i / \times \frac{\partial L_{reg}}{\partial b} = 1 - \sigma_i [\sigma(u_{i+1}) (1 - \sigma(u_{i+1}))] \}$$

Scanned with CamScanner

الآن نصل إلى النهاية بدلالة f_+ على \mathcal{A}_+ حيث $f_+ = \bigvee_{\theta \in \mathcal{A}_+} f(\theta_{++})$

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) g_t \rightarrow \text{Gaussian distribution of momentum bias}$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad \text{④}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta^t} \quad \text{⑤} \quad \rightarrow \left\{ \text{backward shift by one momentum step} \right.$$

$$\theta_t = \theta_{t-1} - \frac{\alpha \hat{m}_t}{\sqrt{\hat{V}_t} + \epsilon} \quad (4)$$

وكان ذلك في-initialized يخاف على الماء

(bias correction) ~~will ignore~~ consider

خط (٥) الأدبي إنجمي بكر.

learning rate α (biggest) step of optimization = 1/n = ④

position \in time and all other momenta $\in \mathbb{R}^n$, it's not

آنراسته م در مورد θ_0 است که بود از تغییرات زیاد برای (θ_0) عالیکاره است.

exponential moving average (EMA) initialize job m_0 and α

at time step t we do $m_t = \alpha g_t + (1-\alpha)m_{t-1}$

$$\nu_0 = 0$$

: first progress

$$\nu_t = (1-B_r) \sum_{i=1}^t \beta_r^{t-i} g_i$$

$$\rightarrow E[\nu_t] = E[(1-B_r) \sum_{i=1}^t \beta_r^{t-i} \cdot g_i]$$

$$= E[g_t] \cdot (1-B_r) \sum_{i=1}^{t-1} \beta_r^{t-i} + h$$

$$= E[g_t] \cdot (1-B_r^t) + h$$

- moving job & running average each $(1-B_r^t)$ part

$g = \text{moving average } (1-B_r^t)$, new data point $\rightarrow 1 - \beta_r \cdot 1/r$

$(\text{avg})(1-B_r^t) + \text{new data point} \rightarrow 1 - \beta_r \cdot 1/r + 1/r$

$$\hat{\nu}_t \leftarrow \frac{\nu_t}{1-B_r^t} \quad \text{ie. step 6) bias correction}$$

: momentum (faster convergence)

$$E[m_t] = E[(1-B_r) \sum_{i=1}^t \beta_r^{t-i} \cdot g_i] = E[g_t] \cdot (1-B_r) \sum_{i=1}^{t-1} \beta_r^{t-i} + h$$

$$= E[g_t] \cdot (1-B_r^t) + h$$

(QH)sw

(r)

$$\frac{\partial w^T H w}{\partial w} = (H + H^T) w$$

$$\rightarrow H = Q \Lambda Q^T \rightarrow H^T = Q \Lambda^T Q^T = Q \Lambda Q^T = H$$

\downarrow
= unitary

$$\rightarrow \frac{\partial w^T H w}{\partial w} = P H w$$

$$\text{Gradient descent} \Rightarrow w_t = w_{t-1} - \epsilon P H w_{t-1} = (I - \epsilon P H) w_{t-1}$$

$$= (I - \epsilon Q \Lambda Q^T) w_{t-1} = Q(I - \epsilon Q \Lambda) Q^T w_{t-1}$$

$$w_t = Q(I - \epsilon Q \Lambda) Q^T w_{t-1}$$

(c)

$$\rightarrow w_1 = Q(I - \epsilon Q \Lambda) Q^T w_0$$

⋮

$$w_t = [Q(I - \epsilon Q \Lambda) Q^T]^+ w_0 \quad Q \text{ is unitary matrix.}$$

$$= [Q(I - \epsilon Q \Lambda) Q^T] \underbrace{[Q(I - \epsilon Q \Lambda) Q^T]}_I \dots [Q(I - \epsilon Q \Lambda) Q^T] w_0$$

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

~~~~~

$$\hookrightarrow \begin{bmatrix} (I - \epsilon Q \Lambda)^+ & & & \\ & \ddots & & \\ & & (I - \epsilon Q \Lambda)^+ & \end{bmatrix}$$

~~new book by John Gutfreund, (J-YA) + is the title of a new book by John Gutfreund.~~

$$\text{Für } \epsilon > 0: \quad \| -\nabla f_i \| \leq 1 \rightarrow -1 \leq \nabla f_i \leq 1$$

$$\gamma_i \rightarrow -\tau \leq -\tau \epsilon \lambda_i \leq 0 \rightarrow \begin{cases} \lambda_i > 0 \\ \epsilon \leq \frac{1}{\tau} \end{cases} \rightarrow \text{weak PSD, HCR}^{\text{alg}}$$

$$w_t \leftarrow w_{t+1} - (\nabla f(w_{t+1}))^{-1} \nabla f(w_{t+1})$$

$\downarrow$   
 $(\text{J}) G^{\text{opt}}(w)$

$$\bar{J}_2 = \left[ \frac{\partial (\nabla_w (w^T H w))}{\partial w} \right]^T = \left[ \frac{\partial}{\partial w} (Y H w) \right]^T = Y H^T = Y H$$

$$\omega_{t+1} \leftarrow \omega_t - J^{-1} \nabla (w^T H w) = \omega_t - \frac{1}{\lambda} H^{-1} (\nabla H w) = \omega_t - \lambda \omega_t = 0$$

$(W^T H w)$  represents the total cost of the solution.

وَلِمَنْ وَلِمَنْ وَلِمَنْ وَلِمَنْ وَلِمَنْ وَلِمَنْ وَلِمَنْ وَلِمَنْ

$O(n^2)$  computation cost

میں اپنے بھائی کو (n) نمبر کا سوتھا ایک ایسا نام لے دیا۔

(cont.) - *the single digit numbers*

• (O(n))  $\leq$  عدد الاعداد المماثلة  $\leq$  عدد العناصر في المدخلات =  $n^2$

$$\bar{J}_t = \frac{1}{T} (y_t - \sum_{k=1}^n (\omega_k + \delta_k) u_{t+k}) \sim \mathcal{N}(0, \sigma^2) \quad (\text{Eq 1})$$

$$E[\bar{J}_t] = \frac{1}{T} E[(y_t - z)^T] = \frac{1}{T} E[y_t^T + z^T - y_t^T z] =$$

$$\frac{1}{T} (y_t^T + E[z^T] - T y_t^T E[z])$$

$$E[z] = E\left[\sum_{k=1}^n (\omega_k + \delta_k) u_k\right] = \sum_{k=1}^n \omega_k u_k$$

$$E[z^T] = E\left[\sum_{a=1}^n (\omega_a + \delta_a) u_a \sum_{b=1}^n (\omega_b + \delta_b) u_b\right]$$

$$= E\left[\sum_{a=1}^n (\omega_a + \delta_a) u_a^T + T \sum_{a=1}^n \sum_{b=1}^n (\omega_a + \delta_a)(\omega_b + \delta_b) u_a u_b^T\right]$$

$$= \sum_{a=1}^n E[(\omega_a + \delta_a)^T] u_a^T + T \sum_{a \neq b}^n \omega_a \omega_b u_a u_b^T$$

$$= \sum_{a=1}^n (1 + \alpha) \omega_a^T u_a^T + T \sum_{a \neq b}^n \omega_a \omega_b u_a u_b^T$$

$$\Rightarrow E[\bar{J}_t] = \frac{1}{T} \left[ y_t^T + \sum_{a=1}^n (1 + \alpha) \omega_a^T u_a^T + T \sum_{a \neq b}^n \omega_a \omega_b u_a u_b^T - T y_t^T \sum_k \omega_k u_k^T \right]$$

$$= \frac{1}{T} \left( y_t^T - \sum_{i=1}^n \omega_i u_i^T \right)^T + \frac{1}{T} \alpha \sum_{i=1}^n \omega_i^T u_i^T$$

$$\Rightarrow E\left[\frac{\partial \bar{J}_t}{\partial u_k}\right] = \frac{\partial E[\bar{J}_t]}{\partial u_k} = -\left(y_t^T - \sum_{i=1}^n \omega_i u_i^T\right) u_k + \alpha \omega_k u_k^T$$

Subject :

Date :

$$E[\hat{J}_i] = \frac{1}{2} (y_i - \sum_{j=1}^n w_j x_{ij})^2 + \frac{\lambda}{2} \sum_{j=1}^n w_j^2$$

لـ  $J_i$  regularization loss function

Loss func.

جـ  $\lambda$  regularization parameter

Generative error

loss function

loss

Scanned with CamScanner