

Machine Learning

- supervised learning
 - regression
 - classification

- unsupervised learning
 - clustering
 - non clustering

Linear Regression

$$J = \frac{1}{2m} \sum_i^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

- 2 - variables $h_{\theta}(x) = \theta_0 + \theta_1 x$

- Gradient Descent

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_i^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_i^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

- Multivariable Linear Regression

$$h_{\theta}(x) = \theta^T x$$

parameter for variable

$$\begin{bmatrix} & & & \\ \vdots & \vdots & \vdots & \vdots \\ & & & \end{bmatrix} \rightarrow x_0 = 1$$

$$\begin{array}{c|ccccc|c} & x_0 & x_1 & x_2 & x_3 & \cdots & y \\ \hline f & | & | & | & | & \cdots & | \\ \hline & n & & & & & \end{array}$$

- Gradient Descent

repeat f

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_i^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

- Feature Scaling +

$$\theta_0 \quad \text{slow} \quad x_{ij} := \frac{x_i}{\text{range}(x_j)} \quad \theta_1$$

- Mean normalization

$$x_{ij} := \frac{x_j - \bar{x}_j}{\text{range}(x_j)}$$

- Polynomial Regression \leftarrow Multivariable Linear Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 \rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$\begin{cases} x_1 = x \\ x_2 = x^2 \\ x_3 = x^3 \end{cases} \rightarrow \text{Feature Scaling}$$

Normal Equation ?

n small : normal equation
 n large : gradient descent
 it

No Feature Scaling

$$\theta = (X^T X)^{-1} X^T y \rightarrow \min_{\theta} J(\theta)$$

design matrix

$$X = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots \\ 1 & x_1^2 & x_2^2 & \dots \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^n & x_2^n & \dots \end{bmatrix} \quad y = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad m \times (n+1) \quad m-1 \text{ vector}$$

\uparrow

m example $[x^{(i)}, y^{(i)}]$

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \rightarrow X = \begin{bmatrix} x_0^{(1)} & x_0^{(2)} & \dots & x_0^{(m)} \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(m)} \end{bmatrix} \quad m \times n$$

$X^T X$ non-invertible

- redundant features [linearly dependent]

- too many features [$m \leq n$]

- delete

- regularization

Vectorization

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}$$

↓

$$[\theta] := [\theta] - \frac{\alpha}{m} \left[\begin{array}{c} (c_1) \cdot x_0^{(1)} + (c_2) \cdot x_0^{(2)} + \dots + (c_m) \cdot x_0^{(m)} \\ (c_1) \cdot x_1^{(1)} + (c_2) \cdot x_1^{(2)} + \dots + (c_m) \cdot x_1^{(m)} \\ (c_1) \cdot x_2^{(1)} + (c_2) \cdot x_2^{(2)} + \dots + (c_m) \cdot x_2^{(m)} \end{array} \right]$$

$$:= [\theta] - \frac{\alpha}{m} \left[\begin{array}{c} x_0^{(1)} \\ x_0^{(2)} \\ \vdots \\ x_0^{(m)} \end{array} \right] \cdot \left[\begin{array}{c} (c_1) \\ (c_2) \\ \vdots \\ (c_m) \end{array} \right]^T$$

$$+ \left[\begin{array}{c} x_1^{(1)} \\ x_1^{(2)} \\ \vdots \\ x_1^{(m)} \end{array} \right] \cdot \left[\begin{array}{c} (c_1) \\ (c_2) \\ \vdots \\ (c_m) \end{array} \right]^T$$

$$+ \left[\begin{array}{c} x_2^{(1)} \\ x_2^{(2)} \\ \vdots \\ x_2^{(m)} \end{array} \right] \cdot \left[\begin{array}{c} (c_1) \\ (c_2) \\ \vdots \\ (c_m) \end{array} \right]^T$$

$$h_\theta(x^{(i)}) = \theta^T x^{(i)} \leftarrow \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$$

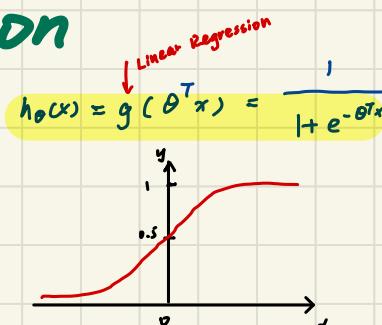
$$\begin{bmatrix} (c_1) \\ (c_2) \\ \vdots \\ (c_m) \end{bmatrix} = \begin{bmatrix} \theta^T x^{(1)} - y^{(1)} \\ \theta^T x^{(2)} - y^{(2)} \\ \vdots \\ \theta^T x^{(m)} - y^{(m)} \end{bmatrix} = X \theta - y$$

m x 1 m x 1 m x 1

$$\theta := \theta - \frac{\alpha}{m} X^T (X \theta - y)$$

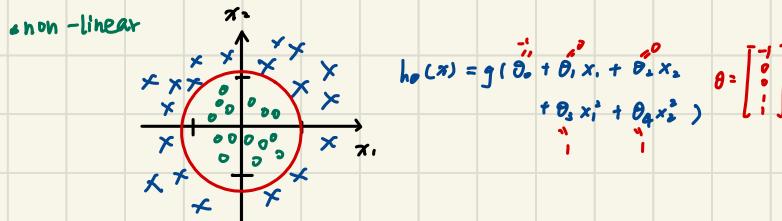
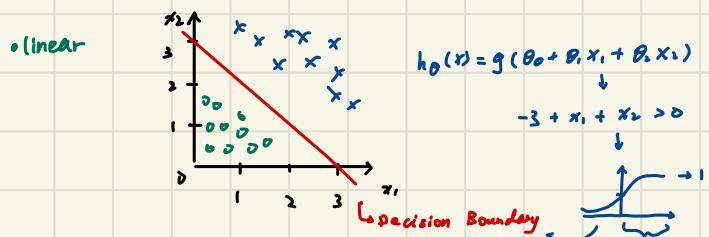
Classification

Sigmoid Function
Logistic Regression



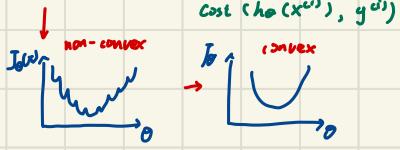
$$h_\theta(x) = P(y=1 | x, \theta)$$

Decision Boundary

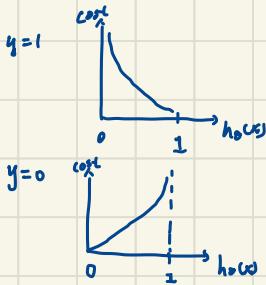


Cost function

$$\text{if } J_{\theta}(x) = \frac{1}{m} \sum_i^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & y=1 \\ -\log(1-h_{\theta}(x)) & y=0 \end{cases}$$



$$\text{cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

↓

$$J(\theta) = -\frac{1}{m} \left[\sum_i^m \text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) \right]$$

Gradient Descent

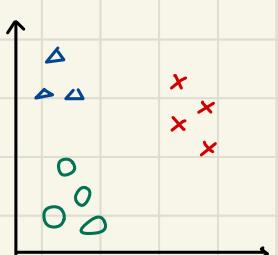
repeat {

$$\begin{aligned} \theta_j := \theta_j - \alpha \frac{\partial J}{\partial \theta_j} &= \theta_j - \alpha \left(-\frac{1}{m} \sum_i^m \frac{y^{(i)}}{h_{\theta}(x^{(i)})} \cdot \frac{\partial h_{\theta}(x)}{\partial z} \cdot \frac{\partial z}{\partial \theta_j} + \frac{(1-y^{(i)})}{1-h_{\theta}(x^{(i)})} \cdot (-1) \cdot \frac{\partial h_{\theta}(x)}{\partial z} \cdot \frac{\partial z}{\partial \theta_j} \right) \\ &= \theta_j - \alpha \left(-\frac{1}{m} \sum_i^m \frac{y^{(i)}}{h_{\theta}(x^{(i)})} \cdot h_{\theta}(x^{(i)})(1-h_{\theta}(x^{(i)})) \cdot x_j - \frac{(1-y^{(i)})}{1-h_{\theta}(x^{(i)})} \cdot h_{\theta}(x^{(i)}) \cdot (1-h_{\theta}(x^{(i)})) \cdot x_j \right) \\ &= \theta_j - \alpha \left(-\frac{1}{m} \sum_i^m [y^{(i)}h_{\theta}(x^{(i)}) - h_{\theta}(x^{(i)}) + y^{(i)}h_{\theta}(x^{(i)})] x_j \right) \\ &= \theta_j - \alpha \sum_i^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j \end{aligned}$$

Vectorization

$$\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y})$$

Multiclass Classification



$$\begin{aligned} h_{\theta}^{(1)}(x) &= P(y=1|x; \theta) \\ h_{\theta}^{(2)}(x) &= P(y=2|x; \theta) \\ h_{\theta}^{(3)}(x) &= P(y=3|x; \theta) \end{aligned}$$

$$\max h_{\theta}^{(i)}(x)$$

Overfitting → Regularization

- perfectly fit example
- × generalize to new data, × predict well

Solution:

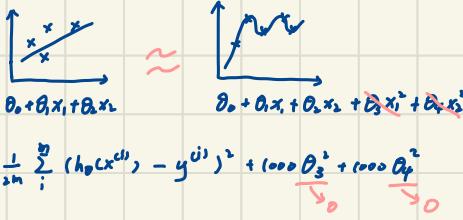
1. Reduce number of features (i)

• manually

• model selection algorithm

2. Regularization

- keep all features, reduce values of θ → $\min_{\theta} \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2$
- works well with lots of features, each contributes a little to predict y



Regularized Linear Regression

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

\downarrow Regularization parameter

Gradient Descent ($x\theta_0$)

$$\theta_j := \theta_j \left(1 - \frac{\partial J}{\partial \theta_j} \right) - \frac{\alpha}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

\downarrow vectorization

$$\theta := \left(1 - \frac{\alpha}{m} \right) \theta - \frac{\alpha}{m} X^T (X\theta - y)$$

Normal equation

$$\theta = (X^T X + \lambda \begin{bmatrix} 0 & 1 & \dots & 1 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix})^{-1} X^T g$$

Regularized Logistic Regression

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log (1-h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Gradient Descent ($x\theta_0$)

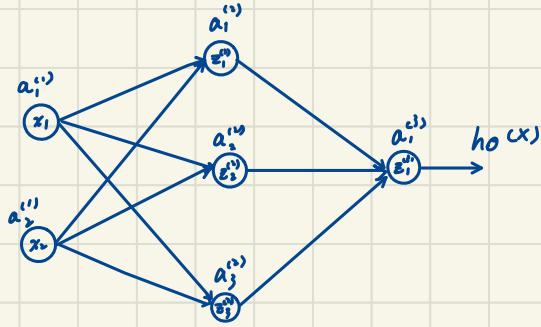
$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j$$

$$\theta_j := \theta_j \left(1 - \frac{\partial J}{\partial \theta_j} \right) - \frac{\alpha}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

\downarrow vectorization

$$\theta := \left(1 - \frac{\alpha}{m} \right) \theta - \frac{\alpha}{m} X^T (g(X\theta) - y)$$

Neutral Network



Layer 1

$a^{(1)}$

$$a^{(1)} = [a_0^{(1)} \ a_1^{(1)} \ a_2^{(1)}]$$

Layer 2

$a^{(2)}$

$$a^{(2)} = [a_0^{(2)} \dots a_b^{(2)}]$$

Layer 3

$a^{(3)}$

$$a^{(3)} = [a_0^{(3)} \dots a_s^{(3)}]$$

$$a_1^{(3)} = g(\theta_{1,0}^{(3)} a_0^{(2)} + \theta_{1,1}^{(3)} a_1^{(2)} + \theta_{1,2}^{(3)} a_2^{(2)})$$

$$a_2^{(3)} = g(\theta_{2,0}^{(3)} a_0^{(2)} + \theta_{2,1}^{(3)} a_1^{(2)} + \theta_{2,2}^{(3)} a_2^{(2)})$$

$$\hookrightarrow \theta_j^{(3)} \Rightarrow s_{j+1} \times (s_j + 1)$$

$$x^0 = a^0$$

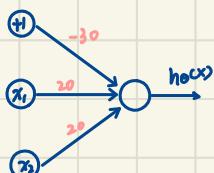
$$\hookrightarrow \theta^0 \cdot a^0 = z^0$$

$$\hookrightarrow g(z^0) = a^0$$

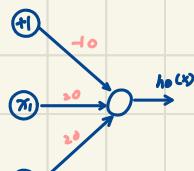
$$\hookrightarrow \theta^{(1)} \cdot a^{(0)} = z^{(1)}$$

$$\hookrightarrow g(z^{(1)}) = h_0(x)$$

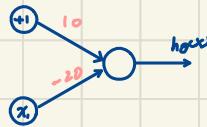
AND



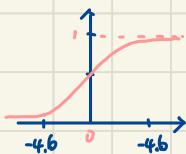
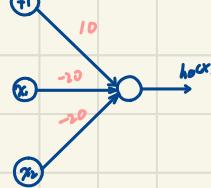
OR



NOT



(NOT x1) AND (NOT x2)



Cost Function

$$J_\theta = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)})) + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\theta_{j,i}^{(l)})^2$$

m - batch size

k - number of classes in multi-classification

L - number of layers

s_l - number of units in Layer l

Backpropagation

$\Delta_{ji}^{(l)}$ use for all $\{x^{(m)}, y^{(m)}\}$

for $m : M$

$$\rightarrow 1) x^{(m)} \Rightarrow a^{(0)}$$

2) for $l : L$

$$z^{(l+1)} = \theta^{(l)} \cdot a^{(l)}$$

$$a^{(l+1)} = g(z^{(l+1)})$$

$$\rightarrow y^{(m)}$$

$$3) \delta^{(l)} = a^{(l)} - y^{(m)}$$

for $l \rightarrow 2$

$$\delta^{(l)} = [\theta^{(l), T} \quad \delta^{(l+1)}] \cdot \underbrace{g'(z^{(l)})}_{a^{(l)} \cdot (1-a^{(l)})}$$

$$\rightarrow 4) \Delta^{(l)} = \Delta^{(0)} + a^{(l)} \cdot \delta^{(l+1)}$$

$$D_{ji}^{(l)} = \frac{1}{n} (\Delta_{ji}^{(l)} + \lambda \theta_{ji}^{(l)}) \quad \text{if } j \neq 0$$

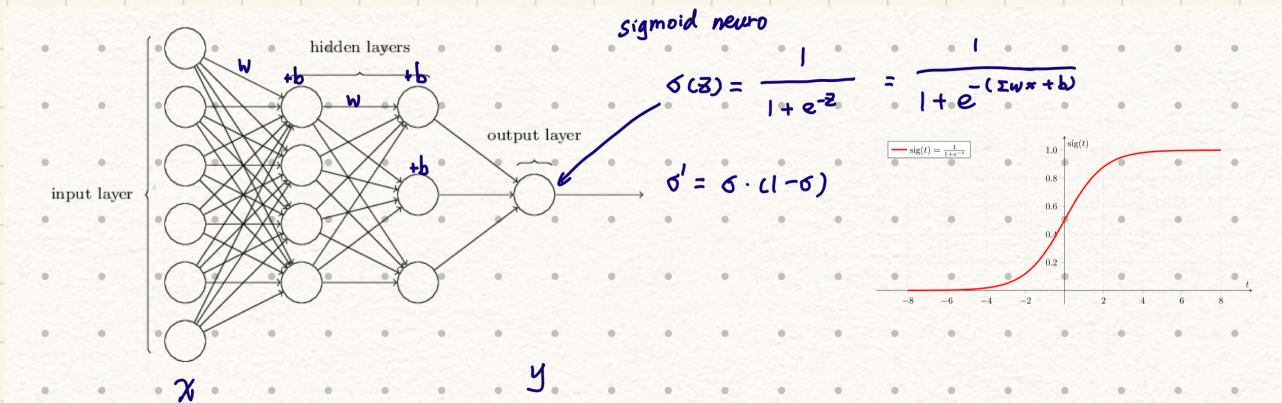
$$D_{ji}^{(0)} = \frac{1}{n} \Delta_{ji}^{(0)}$$

$$j=0$$

$$\rightarrow \theta_{0i}$$



NNN structure



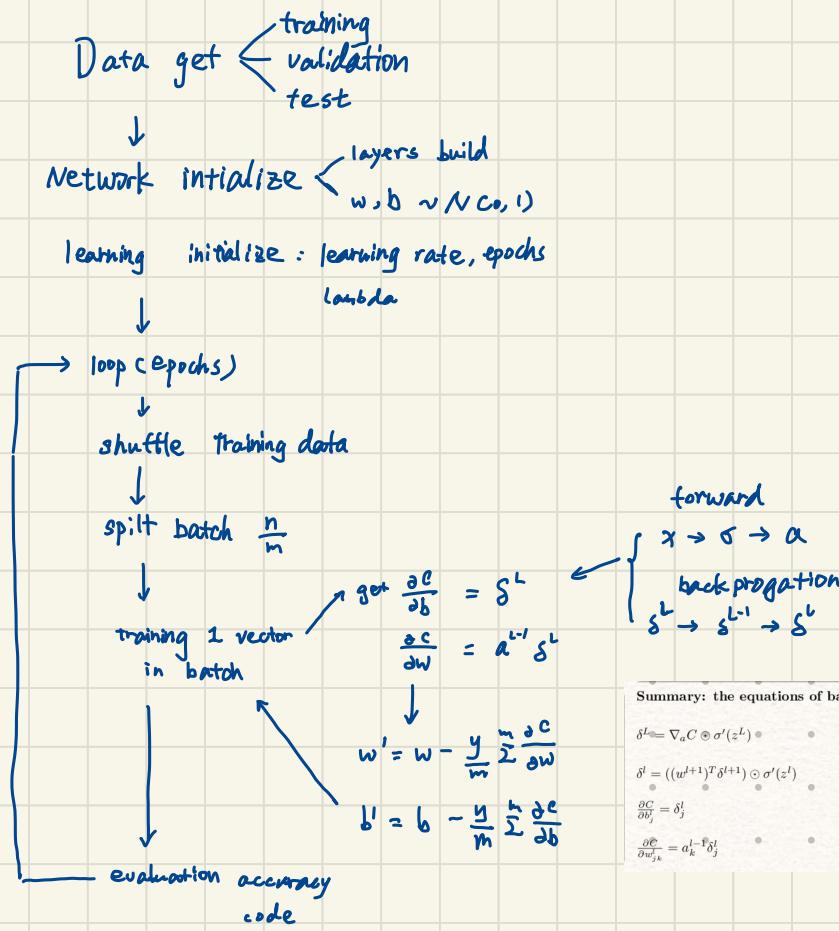
$$y_i \rightarrow x_i$$

$$x_j^J = \sum_k w_{jk}^J \sigma(x_k^{J-1}) + b_j^J$$

$$X^j = W^j \sigma(X^{j-1}) + B^j$$

$$\star C = \frac{1}{2n} \sum [y(x) - a]^2$$

Logic



Summary: the equations of backpropagation

$$\delta^L = \nabla_a C \odot \sigma'(z^L) \odot \dots \odot \quad (\text{BP1})$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (\text{BP2})$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad (\text{BP3})$$

$$\frac{\partial \mathcal{C}}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \quad (\text{BP4})$$

• Proof

$$BP_1 = \frac{\partial C}{\partial z_1} = \frac{\partial C}{\partial z_1^c} \cdot \frac{\partial z_1^c}{\partial z_1^t} = \frac{\partial C}{\partial z_1^c} \circ'(z_1^t)$$

$$s^L = \nabla_a C \odot \sigma'(z^L)$$

$$= (a^L - y) \theta'(z^L)$$

$$\delta_j^L = \frac{\partial C}{\partial z_j^L} = \sum_k \frac{\partial C}{\partial z_k^L} \frac{\partial z_k^L}{\partial z_j^L}$$

$$Z_{\hat{X}}^{(1)} \in \sum_{N=1}^{\infty} G(N) \subset E(\hat{X})$$

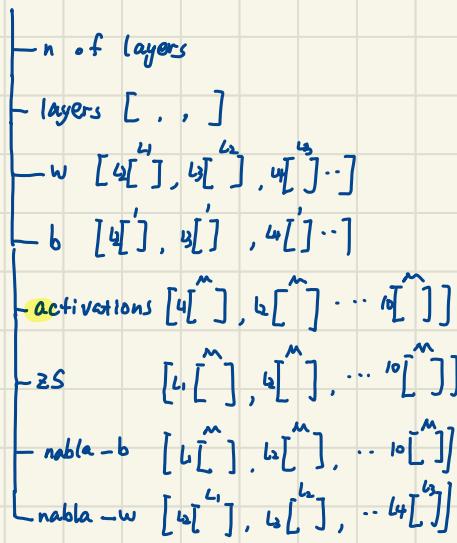
$$\frac{\partial Z_{\hat{X}}^{(1)}}{\partial \hat{x}_i} = W_{\hat{X}}^{(1)} \phi'(c_i)$$

$$\textcircled{1} \quad \begin{matrix} \text{1st layer} \\ w_1 x + b \rightarrow z \end{matrix} \rightarrow a$$

$$\textcircled{2} \quad (\text{last} \rightarrow 1) \\ S^L \rightarrow S^{L+1} \rightarrow S^1 \dots = \frac{\partial \epsilon}{\partial b_j^1}$$

Data structure

Class Network



P1

Batch vector → Batch matrix

$$\begin{aligned}
 &\text{Backprop } x, y \\
 &\downarrow \\
 &\nabla_w Cx, \nabla_b Cx \\
 &\downarrow \\
 &\sum \nabla_w Cx, \sum \nabla_b Cx
 \end{aligned}$$

① 1 → last layer
 $\mathbf{u}^1 = \mathbf{u}^0 \mathbf{W}^1 + \mathbf{b}^1 \rightarrow z^1 \rightarrow a^1$
 $\mathbf{u}^2 = (\mathbf{a}^1 \cdot \mathbf{y}) \odot \sigma'(z^1) \rightarrow z^2 \rightarrow a^2$
 \vdots
 $\mathbf{u}^L = \mathbf{u}^{L-1} \mathbf{W}^L + \mathbf{b}^L \rightarrow z^L \rightarrow a^L$

② last layer → 1
 $\delta^L = (A - Y) \odot \sigma'(z^L)$
 $\delta^{L-1} = (\mathbf{W}^L)^T \delta^L \odot \sigma'(z^{L-1}) = \frac{\partial C}{\partial b^L}$
 \vdots
 $\delta^1 = (\mathbf{W}^1)^T \delta^L \odot \sigma'(z^1)$

$\frac{\partial C}{\partial w_{ik}^L} = a_k^L \delta_i^L$
 $10 \left[\dots \right] n \left[\dots \right]^{L-1} \downarrow$
 $(a^{L-1})^T$

Neuron Function

Softmax

$$a_j^l = \frac{e^{z_j^l}}{\sum_k e^{z_k^l}}$$

$$j \neq k: \frac{\partial a_j^l}{\partial z_k^l} = -\frac{e^{z_j^l}}{(z_k e^{z_k^l})^2} \cdot e^{z_k^l} < 0$$

$$j=k: \frac{\partial a_j^l}{\partial z_j^l} = e^{z_j^l} (\sum_k e^{z_k^l})^{-1} - e^{z_j^l} (\sum_k e^{z_k^l})^{-2} \cdot e^{z_j^l} \\ = e^{z_j^l} (\sum_k e^{z_k^l})^{-1} \left(1 - \frac{e^{z_j^l}}{\sum_k e^{z_k^l}}\right) > 0$$

P3

inverting the softmax layer $z_j^l = f(a_j^l)$

$$a_j^l = e^{z_j^l} / \sum_k e^{z_k^l}$$

$$\ln a_j^l = z_j^l - \ln \sum_k e^{z_k^l}$$

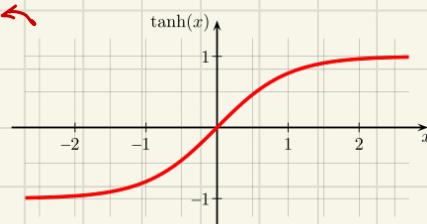
independent of j

Tanh

$$a = \tanh(wx + b)$$

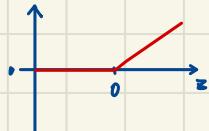
$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \rightarrow \sigma(z) = \frac{1 + \tanh(\frac{z}{2})}{2}$$

normalize



Rectified Linear

$$a = \max(0, wx + b)$$



Cost Function 2 Cross Entropy

$$C = \frac{1}{n} \sum (y - a)^2$$

↓

$$C = -\frac{1}{n} \sum [y \ln a + (1-y) \ln(1-a)]$$

① $C > 0$

$$\textcircled{2} \quad \frac{\partial C}{\partial w_j} = -\frac{1}{n} \sum \left[\frac{y}{\sigma(z)} \cdot \sigma' \cdot x_j - \frac{1-y}{1-\sigma} \cdot \sigma' \cdot x_j \right]$$

$$= -\frac{1}{n} \sum \frac{y - \sigma}{\sigma(1-\sigma)} \sigma' x_j$$

$$\downarrow \sigma' = \sigma(1-\sigma)$$

$$= \frac{1}{n} \sum (\sigma - y) x_j \rightarrow \frac{\partial C}{\partial w_{jk}} = \frac{1}{n} \sum a_k^{L-1} (a_j^L - y_j)$$

$$\textcircled{3} \quad \frac{\partial C}{\partial b} = \frac{1}{n} \sum (\sigma - y) \rightarrow \frac{\partial C}{\partial b_j^L} = \frac{1}{n} \sum (a_j^L - y_j)$$

P2

choose Cost function to eliminate x_j

$$\text{goal: } \frac{\partial C}{\partial w} = x_j(a - y)$$

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial a} \cdot \frac{\partial a}{\partial w} = \frac{\partial C}{\partial a} \cdot \sigma' \cdot x_j$$

$$\frac{\partial C}{\partial a} = \frac{(a-y)}{\sigma' x_j}$$

$$\downarrow \begin{aligned} \sigma' &= \sigma(z) (1-\sigma(z)) \\ &= a(1-a) \end{aligned}$$

$$\frac{\partial C}{\partial a} = \frac{(a-y)}{a(1-a)x_j}$$

$$C = -\frac{1}{x_j} [y \ln a + (1-y) \ln(1-a)] + \text{constant}$$

$$\frac{\partial C}{\partial b} = \frac{\partial C}{\partial a} \cdot \frac{\partial a}{\partial b} = \frac{(a-y)}{a(1-a)x_j} \cdot \sigma'$$

$$= \frac{1}{x_j} (a-y)$$

Cost function

$$C = -\ln a_y^L$$

$$\begin{aligned} \frac{\partial C}{\partial b_j^L} &= \frac{\partial C}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L} \cdot \frac{\partial z_j^L}{\partial b_j^L} = \frac{-1}{a_j^L} \cdot \frac{e^{z_j^L}}{\sum_k e^{z_k^L}} \left(1 - \frac{e^{z_j^L}}{\sum_k e^{z_k^L}} \right) = a_j^L - 1 \\ &\stackrel{j \neq L}{=} -\frac{1}{a_j^L} \cdot \frac{(e^{z_j^L})^2}{(\sum_k e^{z_k^L})^2} = a_j^L \quad \left. \begin{aligned} \frac{\partial C}{\partial b_j^L} &= a_j^L - y \end{aligned} \right\} \end{aligned}$$

$$\begin{aligned} \frac{\partial C}{\partial w_{jk}^L} &= \frac{\partial C}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L} \cdot \frac{\partial z_j^L}{\partial w_{jk}^L} = \frac{a_j^L}{a_j^L} (a_j^L - 1) \cdot x_k^{L-1} \\ &\stackrel{j \neq k}{=} a_j^L \cdot x_k^{L-1} \quad \left. \begin{aligned} \frac{\partial C}{\partial w_{jk}^L} &= (a_j^L - y) \cdot a_k^{L-1} \end{aligned} \right\} \end{aligned}$$

Backpropagation

$$\begin{aligned} \delta_j^L &\equiv \frac{\partial C}{\partial z_j^L} \\ &= \frac{\partial C}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L} = \frac{\partial C}{\partial b_j^L} = a_j^L - y_j \end{aligned}$$

Regularization

L2

$$C = C_0 + \frac{\lambda}{2n} \sum_w w^2$$

↳ original cost function

$$\frac{\partial C}{\partial w_{jk}^L} = \frac{\partial C_0}{\partial w_{jk}^L} + \frac{\lambda}{n} w_{jk}^L$$

$$\frac{\partial C}{\partial b_j^L} = \frac{\partial C_0}{\partial b_j^L}$$

learning

$$b_j^L = b_j^L - y \frac{\partial C_0}{\partial b_j^L}$$

$$w_{jk}^L = w_{jk}^L - y \frac{\partial C_0}{\partial w_{jk}^L} - \frac{\eta \lambda}{n} w_{jk}^L$$

$$= (1 - \frac{\eta \lambda}{n}) w_{jk}^L - y \frac{\partial C_0}{\partial w_{jk}^L}$$

weight decay

$$w = (1 - \frac{\eta \lambda}{n}) w - \frac{y}{m} \sum_x \frac{\partial C_x}{\partial w}$$

↳ training set

$$b = b - \frac{y}{m} \sum_x \frac{\partial C_x}{\partial b}$$

minibatch size

L1

$$C = C_0 + \frac{\lambda}{n} \sum_w |w|$$

$$\frac{\partial C}{\partial w} = \frac{\partial C_0}{\partial w} + \frac{\lambda}{n} \operatorname{sgn}(w)$$

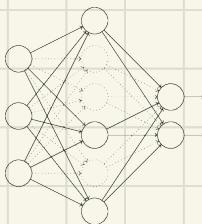
learning

$$w_{jk}^L = w_{jk}^L - \frac{y \lambda}{n} \operatorname{sgn}(w_{jk}^L) - y \frac{\partial C_0}{\partial w_{jk}^L}$$

mini-batch

$$w = w - \frac{y \lambda}{n} \operatorname{sgn}(w) - \frac{y}{m} \sum_x \frac{\partial C_x}{\partial w}$$

Dropout



① delete half of neurons

② forward $x \rightarrow a$
backpropagate $\nabla C \leftarrow a$

③ $w = w - y \nabla C$
 $b = b - y \nabla C$

④ restore dropout neurons

new batch

Weight initialization

$$z = \sum w_j x_j + b$$

$$w_j \sim N(0, 1)$$

$$b \sim N(0, 1)$$

$$\begin{array}{c} n=1000 \\ \text{input} \\ \downarrow \\ \frac{1}{2} : 1 \\ \downarrow \\ \frac{1}{2} : 0 \end{array}$$

$$\sim N(0, 500 \cdot 1^2)$$

$$\sum w_j x_j \sim N(0, (500)^2) \text{ not the same } v_j$$

$$+ b \sim N(0, (500)^2 + 1^2)$$

$$+ b \sim N(0, 501 \cdot 1^2)$$

$$z \sim N(0, (\sqrt{501})^2)$$



$$w_j \sim N(0, (\frac{1}{\sqrt{n}})^2)$$

$$b \sim N(0, 1)$$

$$\begin{array}{c} n=1000 \\ \text{half 1} \\ \text{half 0} \end{array}$$

$$\sum w_j x_j \sim N(0, 500 \cdot \frac{1}{1000})$$

$$+ b \sim N(0, \frac{1}{2} + 1)$$

$$z \sim N(0, (\sqrt{\frac{3}{2}})^2)$$

Compare L2

$$w \sim N(0, 1) \quad w = (I - \frac{\gamma \lambda}{n})w - \frac{\gamma}{m} \sum_x \frac{\partial C_x}{\partial w}$$

$$b = b - \frac{\gamma}{m} \sum_x \frac{\partial C_x}{\partial b}$$

$$\text{weight decay: } \frac{\gamma \lambda}{n} w + \frac{\gamma}{m} \sum_x \frac{\partial C_x}{\partial w}$$

$$\frac{\partial C_x}{\partial w} \text{ random, no use}$$

• λ too small, weight decay → first epoch learning

• $\gamma \lambda \ll n$, decay $\rightarrow e^{-\frac{\gamma \lambda}{m}}$

• λ too large, $w \approx \frac{1}{\sqrt{n}}$ → no decay

$$\begin{matrix} i & i \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

$$i \cdot j + j \cdot k + \dots = n_w$$

$$\begin{array}{l} x \cdot a' \cdot a^2 \\ \downarrow \quad \downarrow \quad \downarrow \\ x \sigma(w^T x + b) \end{array} \rightarrow \sigma \left[w^T \sigma \left(w^T x + b \right) + b \right]$$

$$\xrightarrow{\text{epoch } m_1 \rightarrow m_2} \frac{n}{m} \cdot \lambda$$

$$\left\{ \begin{array}{l} (1 - \frac{\gamma \lambda}{n})^{\frac{n}{m}} \\ n \gg \gamma \lambda \\ \lim_{n \rightarrow \infty} (1 + \frac{1}{n})^n = e \end{array} \right.$$

$$\begin{aligned} & \rightarrow \lim_{n \rightarrow \infty} (1 + \frac{\gamma \lambda}{n})^{\frac{n}{m}} \xleftarrow{\gamma \lambda} (1 + \frac{\gamma \lambda}{m}) \\ & = e^{(\frac{\gamma \lambda}{m})} \end{aligned}$$

$$\begin{array}{l} \frac{\partial C}{\partial w} = 0 \rightarrow \frac{1}{n} \sum w^2 + C_x = C \\ \frac{\partial C}{\partial b} = 0 \end{array}$$

$$b = \text{constant}$$

$$\begin{array}{l} C_x = \text{constant} \\ C_x = \frac{1}{n} \sum (y - a)^2 \end{array}$$

?

•

Early Out

no - improvement - in - n epochs

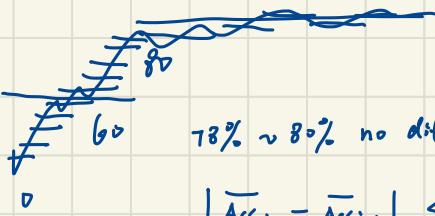
PS₁₎

epochs w b nn

↓
validation accuracy \approx μ_{nn}

simple check if $i > i-n$

2)



78% \approx 80% no differences

$$|\bar{Acc}_i - \bar{Acc}_{i-1}| < \Delta_{Acc} \text{ stop at } i$$

$$\bar{Acc} = \frac{\sum_{j=1}^{Macc} Acc_j}{Macc}$$

if $n_{\text{epoch}} > Macc$ start

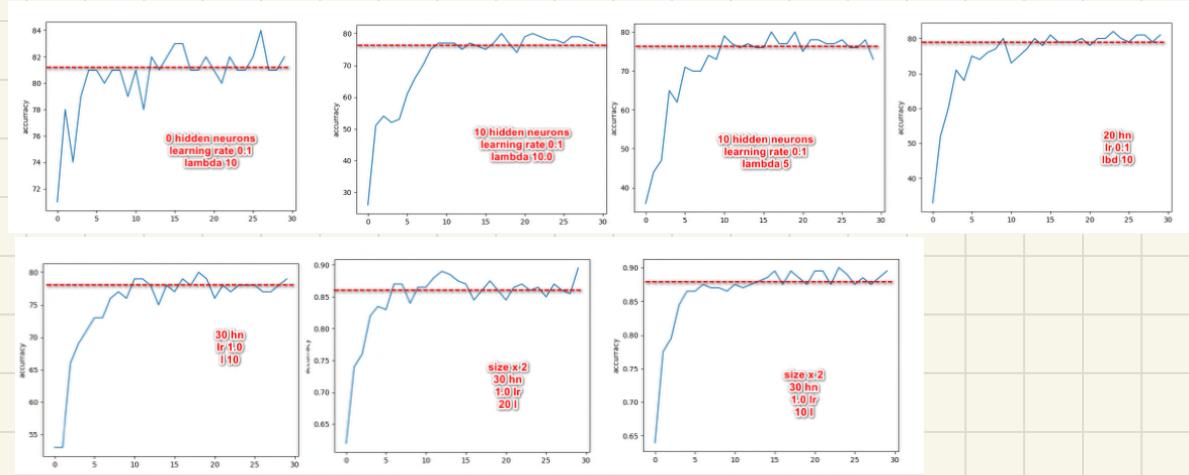
$$\text{for each epoch } j: \bar{Acc} = \frac{\sum_{j=n+1}^j Acc_j}{Macc}$$

if len C

too
specific
↓
rough!

Get Hyper parameter

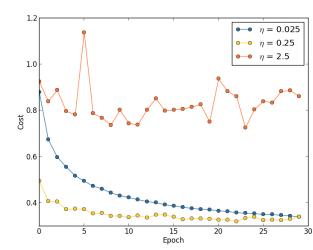
1. Broad strategy
 - ① 0 hidden neurons
small training data & validation data → fast feedback
 - ② $\lambda \downarrow$ same rate
 - ③ learning rate $\eta \downarrow$ or \uparrow to best
 - ④ $\lambda \uparrow$ or \downarrow to best
 - ⑤ + hidden neurons
 - ⑥ ③ ~ ⑤ = best



2. learning rate η : same epochs
mini-batch-size
 λ

} compare

- ① choose estimate η
- ② if Cost \downarrow
 ↓
 if Cost \uparrow , MW
 ↓
 largest η , C \downarrow
- ③ $\eta \downarrow$, not too slow



3. regularization parameters λ

- ① $\lambda=0 \rightarrow y$
- ② initial $\lambda=1$
- ③ $\lambda+=10 / \lambda-=10$

3. auto select

Random search for hyper-parameter optimization
grid search
Practical Bayesian Optimization of Machine Learning Algorithms

4. mini batch size
 - x small specific
 - x big x weights change

Hessian Technique

goal : choose $\Delta w \rightarrow \min(C(w + \Delta w))$

$$\begin{aligned}
 C(w + \Delta w) &= C(w) + \sum_j \frac{\partial C}{\partial w_j} \Delta w_j + \frac{1}{2} \sum_{j,k} \frac{\partial^2 C}{\partial w_j \partial w_k} \Delta w_j \Delta w_k + \dots \\
 &= C(w) + \nabla C \cdot \Delta w + \frac{1}{2} \Delta w^T H \Delta w + \dots \quad \text{converge} > \text{standard gradient descent} \\
 &\approx C(w) + \nabla C \cdot \Delta w + \frac{1}{2} \Delta w^T H \Delta w \\
 &\downarrow \text{constant} \\
 \frac{\partial}{\partial \Delta w} &= \nabla C^T + \frac{1}{2}(H + H^T) \Delta w = 0 \\
 \nabla C + \frac{1}{2} \Delta w^T (H^T + H) &= 0 \\
 &\quad \xrightarrow{\frac{\partial^2 C}{\partial j \partial k}} \Rightarrow H^T = H \\
 \nabla C + \frac{1}{2} \Delta w^T 2H &= 0 \\
 \Delta \tilde{w} &= -\frac{\nabla C}{H}
 \end{aligned}$$

Logic:

$$\begin{array}{l}
 \textcircled{1} \quad w \xrightarrow{\quad} \\
 \textcircled{2} \quad w' = w - H^{-1} \nabla C \xrightarrow{\quad} \\
 \textcircled{3} \quad w'' = w' - H'^{-1} \nabla C' \xrightarrow{\quad} \\
 \dots
 \end{array}$$

Momentum-based gradient descent

$$\begin{aligned}
 \text{velocity} \\
 v' &= \mu v - \gamma \nabla C
 \end{aligned}$$

$$w' = w + v'$$

$$\textcircled{1} \quad \mu = 0 \rightarrow w' = w - \gamma \nabla C$$

$$\textcircled{2} \quad \mu = 1 \rightarrow v \Delta w = v - \gamma \nabla C > 0$$

$$\Rightarrow \Delta w = \frac{v - \gamma \nabla C - \gamma \nabla C}{\gamma v} < \Delta w, \dots \text{slow}$$

$$\textcircled{3} \quad 0 < \mu < 1 \rightarrow \textcircled{1} \quad \Delta w = \mu v - \gamma \nabla C \quad \textcircled{2} \quad \Delta w = \mu(\mu v - \gamma \nabla C) - \gamma \nabla C$$

$$\textcircled{1} < \textcircled{2} \quad \dots \Delta w \rightarrow \textcircled{1} \quad \text{gentle slow}$$

$$\textcircled{4} \quad \mu < 0 \rightarrow v \Delta w = \mu v - \gamma \nabla C$$

$$\Rightarrow \Delta w = \mu(\mu v - \gamma \nabla C) - \gamma \nabla C$$

$$\Delta w < \textcircled{1} \quad \Delta w < \textcircled{2} \quad \text{too slow}$$

$$\textcircled{5} \quad \mu > 1 \rightarrow \textcircled{1} \quad \Delta w = \mu v - \gamma \nabla C$$

$$\Rightarrow \Delta w = \mu(\mu v - \gamma \nabla C) - \gamma \nabla C$$

$$\gg \textcircled{2} \quad \gg \textcircled{2} \quad \text{too fast}$$

Visual Learning

