maximum Likelihood

Derive MLE of event rate parameter λ

Let x_1, x_2, \dots, x_n be iid poisson random variables with prob mass function,

pmf
of
poisson
distn

$$p(x_i, \lambda) = \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

Now, the joint prob of all variables x_i , called likelihood function, is given by

Likelihood

$$L(\lambda) = \prod_{i=1}^n p(x_i, \lambda)$$

$$= \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

The log likelihood of pmf is,

Log
Likelihood

$$\ln L(\lambda) = \ln \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$= \sum_{i=1}^n \ln \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$= \sum_{i=1}^n [-\lambda + x_i \ln \lambda - \ln(x_i!)]$$

$$\ln L(\lambda) = -n\lambda + \ln \lambda \sum_{i=1}^n x_i - \sum_{i=1}^n \ln(x_i!)$$

To get the maximum likelihood estimate of the parameter λ , we maximize the Log Likelihood ($\ln L(\lambda)$) wrt λ .

$$0 = \frac{\partial}{\partial \lambda} \ln L(\lambda)$$

$$= \frac{\partial}{\partial \lambda} [-n\lambda + \lambda n \bar{x} - \bar{x} \ln \lambda]$$

$$0 = -n + \frac{\bar{x} n}{\lambda} - 0$$

$$n = \frac{\bar{x} n}{\lambda}$$

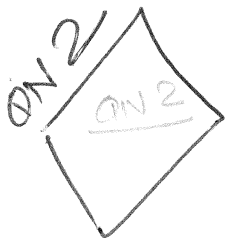
MLE of λ

$$\lambda = \frac{\sum_{i=1}^n x_i}{n}$$

number of samples \rightarrow

Ans

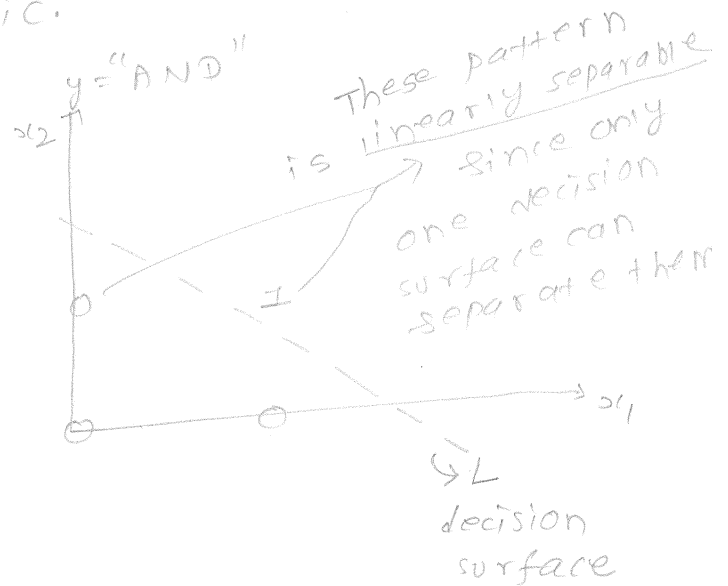
Here, the maximum likelihood estimate of the poisson distribution parameter λ is just the mean (or expectation) of the distribution.



XOR problem in Logistic Regression

To describe XOR problem in LR, I shall start with 'AND' logic.

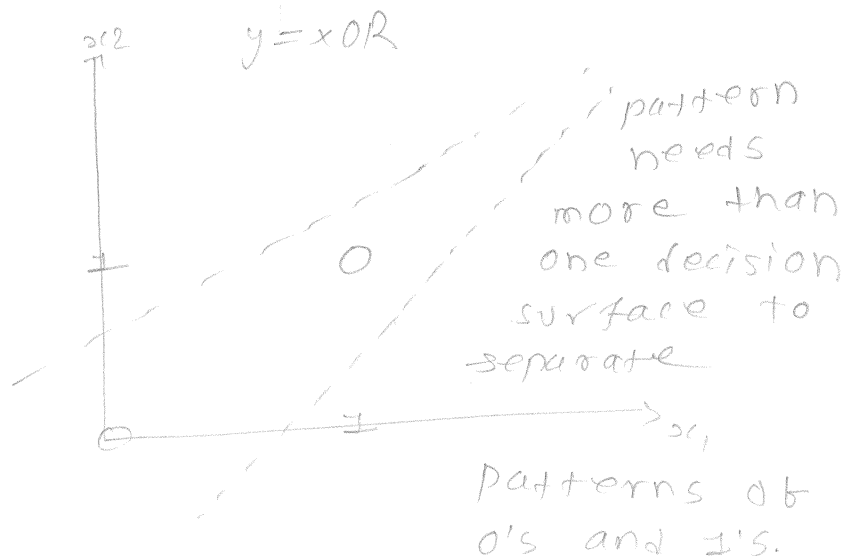
x_1	x_2	$y = x_1 \& x_2$
0	0	0
0	1	0
1	0	0
1	1	1



Now, look at XOR,

Truth table

x_1	x_2	$y = \text{XOR}$
0	0	0
0	1	1
1	0	1
1	1	0



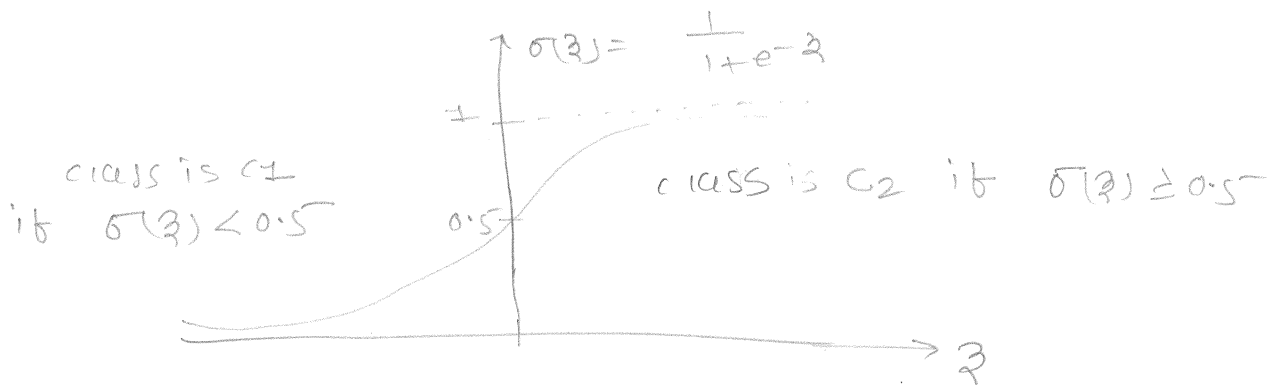
\therefore XOR logic is not linearly separable.

P.T.O.

Now we know that XOR logic is not linearly separable. We also show that Logistic Regression is linear and it can only classify binary classifications that are linearly separable.

The classifier used in LR is sigmoid

$$\text{function } \sigma(z) = \frac{1}{1+e^{-z}}$$



In Logistic Regression,

$$\text{pattern is } \pm \quad \text{if} \quad \frac{1}{1+e^{-w^T x}} \geq \frac{1}{2}$$

$$0 \quad \text{if} \quad \frac{1}{1+e^{-w^T x}} < \frac{1}{2}$$

let's set the ^{equal} bias to separator value:

$$\frac{1}{1+e^{-w^T x}} = \frac{1}{2}$$

$$2 = 1+e^{-w^T x}$$

$$\frac{1}{1} = e^{-w^T x} = \frac{1}{e^{w^T x}}$$

$$e^{w^T x} = 1$$

$$w^T x = \ln 1 = 0$$

$$\Rightarrow \boxed{\sum_i w_i x_i = 0}$$

→ this means LR is linear classifier.

conclusion: \Rightarrow LR is linear classifier

2) XOR problem is non-linear and linearly inseparable

\therefore LR can NOT perfectly classify dataset that follows XOR logic.

ANS

Q3

Gradient of Logistic Regression

for linear regression, hypothesis $h = w^T x$

for logistic regression, hypothesis $h = \sigma = \frac{1}{1 + e^{-w^T x}}$

$$h = \sigma = \frac{1}{1 + e^{-w^T x}}$$

(1) I write $w^T x$ as wx since they are just matrix dot product of two matrices

$$\text{likelihood function } P(t|w) = \prod_{n=1}^N h_n^{t_n} (1 - h_n)^{1 - t_n}$$

$$\text{-ve log likelihood, } E \text{ or } J = -\ln P(t|x)$$

$$E = -\ln \prod_{n=1}^N h_n^{t_n} (1 - h_n)^{1 - t_n}$$

$$E = \sum_n E_n = \sum_{n=1}^N (-t_n \ln h_n - (1 - h_n) \ln (1 - h_n))$$

for any n th sample the cost can be written as,

$$\text{LOSS } E = -t \ln h - (1 - h) \ln (1 - h) \quad (2)$$

$$\text{here } h = \sigma = \sigma(w^T x) = \sigma(wx)$$

Before proceeding further, I would derive derivative of sigmoid function.

$$\sigma(z) = \frac{1}{1+e^{-z}} \Rightarrow \frac{d\sigma}{dz} = \frac{-1 \cdot e^{-z} \cdot (-1)}{(1+e^{-z})^2}$$

$$\frac{d\sigma}{dz} = \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}+1-1}{1+e^{-z}}$$

$$= \frac{1}{1+e^{-z}} \cdot \left(\frac{e^{-z}+1}{1+e^{-z}} - \frac{1}{1+e^{-z}} \right)$$

$$= \frac{1}{1+e^{-z}} \cdot \left(1 - \frac{1}{1+e^{-z}} \right)$$

$$\boxed{\frac{d\sigma(z)}{dz} = \sigma(1-\sigma)}$$

similarly $\boxed{\frac{d\sigma(az)}{dz} = a \sigma(1-\sigma)}$ (here $a \neq a(z)$)

∴
derivative
of
sigmoid
function

$$\boxed{\frac{d\sigma(az)}{dz} = a \sigma(1-\sigma)}$$

Now, going back to the problem, let's calculate the gradient of loss function.

$$\frac{\partial E}{\partial w} = \frac{\partial}{\partial w} [-t \ln h - (1-t) \ln(1-h)]$$

$$= -\frac{t}{h} \frac{\partial h}{\partial w} - \frac{(1-t)(-1)}{1-h} \frac{\partial h}{\partial w}$$

$$= -\frac{t}{h} \frac{\partial \sigma(wx)}{\partial w} + \frac{1-t}{1-h} \frac{\partial \sigma(wx)}{\partial w}$$

$$= -\frac{t}{h} \times h(1-h) + \frac{1-t}{1-h} \times h(1-h)$$

Note
 $h = \sigma(wx)$

$$\therefore \frac{d\sigma(a)}{da} = \sigma(1-\sigma) \text{ and } \sigma \text{ is } h.$$

$$= -tx - \cancel{txh} + xh - \cancel{txh}$$

$$= xh - tx$$

$$\boxed{\frac{\partial E}{\partial w} = (h-t)x}$$

(here this is for a single sample, but for total loss we add up all the losses)

$$\Rightarrow \nabla_w E = \sum_{n=1}^N (h_n - t_n) x_n$$

Q.E.D.

Q24
ON 4LOGISTIC REGRESSION IN SKLEARN

part a - the cost function for L2 penalized logistic regression is given by,

$$E = \underbrace{\frac{1}{2} w^T w}_{\text{L2 regularizer}} + C \underbrace{\sum_{n=1}^N \ln(1 + e^{-t_n w^T x_n})}_{\text{cost for LR}}$$

①

(from sklearn)

we have to show, the -ve log likelihood,

$$-\ln p = \sum_{n=1}^N \ln(1 + e^{-t_n w^T x_n}) \quad \text{--- ②}$$

Now, the posterior probabilities for class 0 and 1 are,

posterior prob for LR

$$p(0|x) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

$$p(1|x) = 1 - \sigma(w^T x) = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$$

③

the likelihood function is,

likelihood for LR

$$P(\mathcal{D}|w) = \prod_{n=1}^N h_n^{t_n} (1 - h_n)^{(1-t_n)}$$

④

The -ve log likelihood is,

$$-\ln p(t|w) = -\ln \prod_{n=1}^N h_n^{t_n} \cdot (1-h_n)^{1-t_n} = -\sum_{n=1}^N \ln(h_n^{t_n} \cdot (1-h_n)^{1-t_n})$$

$$= -\sum_{n=1}^N [\ln(h_n^{t_n}) + \ln((1-h_n)^{1-t_n})]$$

cost for LR

$$E = -\ln p = -\sum_{n=1}^N [\underbrace{t_n \ln h_n}_{\text{class 0}} + \underbrace{(1-t_n) \ln(1-h_n)}_{\text{class 1}}] \quad (5)$$

where $t_n \in \{0, 1\}$

$$h_n = \frac{1}{1 + e^{-w^T x_n}}$$

$$1-h_n = 1 - \frac{1}{1 + e^{-w^T x_n}}$$

$$= \frac{1 + e^{-w^T x_n} - 1}{1 + e^{-w^T x_n}}$$

$$= \frac{e^{-w^T x_n}}{1 + e^{-w^T x_n}}$$

$$= \frac{1}{e^{w^T x_n} + 1}$$

$$1-h_n = \frac{1}{1 + e^{w^T x_n}}$$

$$E = -\ln p$$

$$= - \sum_{n=1}^N \left[t_n \ln \left(\frac{1}{1+e^{-w^T x_n}} \right) + (1-t_n) \ln \left(\frac{1}{1+e^{w^T x_n}} \right) \right]$$

$t_n \in \{0, 1\}$

$$= - \sum_{n=1}^N \ln \frac{1}{\begin{cases} 1+e^{-w^T x_n} & \text{for } t_n=1 \\ 1+e^{w^T x_n} & \text{for } t_n=0 \end{cases}}$$

$$= - \sum_{n=1}^N \ln \frac{1}{1+e^{-t_n w^T x_n}} \quad \text{where new label } t_n' \in \{-1, 1\}$$

$$E = + \sum_{n=1}^N \ln (1+e^{-t_n w^T x_n})$$

$Q \in \mathbb{D}$

putting back dummy variable t_n' as t_n and now $t_n \in \{-1, 1\}$

Qn 4b Find 'C' parameter of LR cost fn from sklearn.

Now, the L2 regularized cost function for LR

$$E = E_w + E_p$$

$$= \frac{1}{2} w^T w + (-\ln p)$$

$$E = \frac{1}{2} w^T w - \sum_{n=1}^N \left[t_n \ln h_n + (1-t_n) \ln (1-h_n) \right]$$

$$t_n \in \{0, 1\}$$

from
Bishop

$$E = \frac{1}{2} w^T w + \sum_{n=1}^N \ln (1 + e^{t_n w^T x_n}) \quad \text{--- (1)}$$

$t_n \in \{-1, 1\}$

To get maximum Likelihood Estimate of parameter w ,

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} w^T w + \sum_{n=1}^N \ln (1 + e^{t_n w^T x_n}) \right]$$

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} w^T w + \frac{1}{C} \sum_{n=1}^N \ln (1 + e^{t_n w^T x_n}) \right]$$

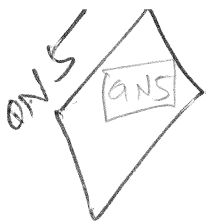
$$\text{but, } \hat{w} = \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} w^T w + C \sum_{n=1}^N \ln (1 + e^{t_n w^T x_n}) \right]$$

(from sklearn)

$$C = \frac{1}{\lambda}$$

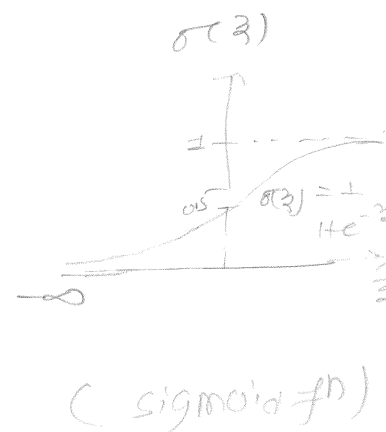
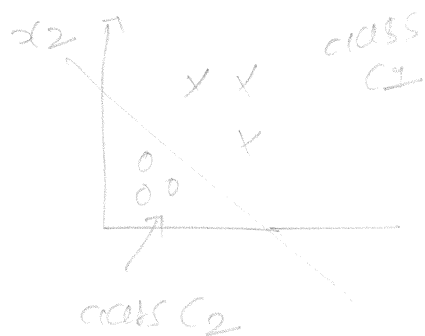
Ans

Extra note
To add bias term
we use $(w^T + c)$
term in
learn.



softmax Regression as special case of Logistic Regression

FOR LR



(binary classification)

the posterior prob of class 1 is
posterior prob

$$P(c_1 | x; w) = P(c_1 | x) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

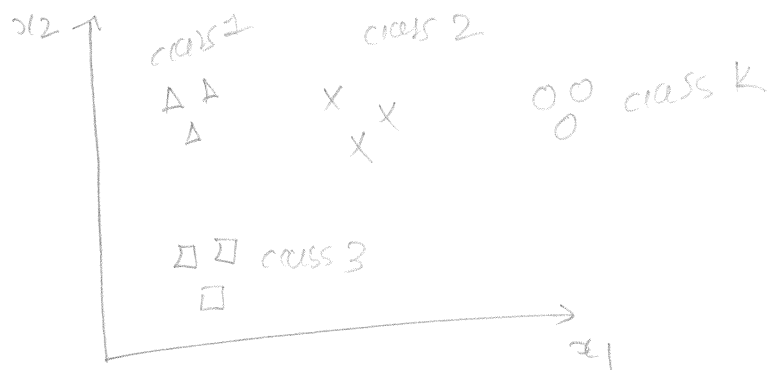
the posterior prob of class 2 is,

$$P(c_2 | x) = 1 - \sigma(w^T x) = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$$

hypothesis the hypothesis for LR,

$$L(w) = \begin{bmatrix} \frac{1}{1 + e^{-w^T x}} \\ \frac{e^{-w^T x}}{1 + e^{-w^T x}} \end{bmatrix} \begin{matrix} \rightarrow \text{for class 1} \\ \rightarrow \text{for class 2} \end{matrix}$$

Again, FOR SR (Softmax Regression or Multiclass Logistic Regression)



total number of samples = N

In softmax Regression, the posterior probability of any class C_k is given by,

input data X	target
$\begin{bmatrix} x_1 & \text{some values} \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$	$\begin{bmatrix} \text{triangle } t_1 \\ \text{cross } t_2 \\ \text{rectangle} \\ \vdots \\ \text{circle } t_N \end{bmatrix}$

$$P(C_k | x, w) = \frac{e^{w_k^T x}}{\sum_{k=1}^K e^{w_k^T x}}$$

posterior prob.

[Bishop, 4.3.4
eq 4.1104
p. 209 / Pdf 227]

note $a_k = w_k^T x$ is called activation fn for the parameter vector w_k

hyp:

$$h = \frac{1}{\sum_{k=1}^K e^{w_k^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \\ \vdots \\ e^{w_K^T x} \end{bmatrix}$$

plus sign
 $+ w_k^T x$

sigmoid is -
softmax is +

Show $h(w) = h(w - \psi)$

[SR has overparameterization property]

$$h(w - \psi) = \frac{e^{(w - \psi)^T x}}{\sum_{k=1}^K e^{(w - \psi)^T x}}$$

ψ is any fixed vector

$$= \frac{e^{w^T x} \cdot e^{-\psi^T x}}{\sum_k e^{w^T x} \cdot e^{-\psi^T x}}$$

$$= \frac{e^{w^T x} \cdot e^{-\psi^T x}}{\underbrace{e^{-\psi^T x}}_{\text{does not depend on } k} \sum_k e^{w^T x}}$$

$$= \frac{e^{w^T x}}{\sum_k e^{w^T x}}$$

$h(w - \psi) = h(w)$

\therefore If we change any parameter vector $w_k \rightarrow w_k - \psi$ we get same hypothesis for softmax Regression.

for two-class SR,

$$h = \frac{1}{\sum_k e^{w_k^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \end{bmatrix}$$

$$= \frac{1}{e^{w_1^T x} + e^{w_2^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \end{bmatrix}$$

use overparametrization property of SR

$$h(w - \psi) = h(w) \quad \text{where } \psi = w_2$$

$$= \frac{1}{e^{(w_1 - w_2)^T x} + 1} \begin{bmatrix} e^{(w_1 - w_2)^T x} \\ 1 \end{bmatrix}$$

write $\vec{w}_2 - \vec{w}_1 = \vec{w}$

$$= \frac{1}{1 + e^{-w^T x}} \begin{bmatrix} e^{-w^T x} \\ 1 \end{bmatrix}$$

$$(h)_{SR} = \begin{bmatrix} \frac{e^{-w^T x}}{1 + e^{-w^T x}} \\ \frac{1}{1 + e^{-w^T x}} \end{bmatrix} = (h)_{LR} \quad \# \text{ proved}$$

↑ this is same as h for LR

Qn6 Qn6

Gradient of softmax regression

example of softmax function

$$x = \begin{bmatrix} x_1=1 \\ x_2=2 \\ x_3=3 \\ x_4=4 \\ x_5=1 \\ x_6=2 \\ x_N=3 \end{bmatrix}$$

$$e^x = \begin{bmatrix} e^{x_1} = 2.72 \\ e^{x_2} = 7.39 \\ e^{x_3} = 20.09 \\ e^{x_4} = 54.6 \\ \vdots \\ e^{x_N} = 20.09 \end{bmatrix}$$

sum = 114.98

$$\text{softmax}(x) = \begin{bmatrix} e^{x_1}/\text{sum} = 0.024 \\ e^{x_2}/\text{sum} = 0.054 \\ e^{x_4}/\text{sum} = 0.475 \\ e^{x_N}/\text{sum} = 0.175 \end{bmatrix}$$

label 4 has maximum value.

- softmax regression is multiclass regression (classification) scheme where each data sample belongs to one of the classes.

the posterior probability of class C_K , i.e. the data x belongs to class K is,

posterior prob of class C_K

$$p(C_K | x, w) = \frac{e^{w_K^T x}}{\sum_{k=1}^K e^{w_k^T x}} \quad \text{--- (1)}$$

(Bishop 4.3.4 page 209/227 eq 4.104)

Likelihood

likelihood is, $\ell(w) = \prod_{n=1}^N p(t_n | x_n, w) \quad \text{--- (2)}$

-ve log likelihood is, $-\ln \ell(w) = -\ln \prod_n p(t_n | x_n, w)$

-ve log likelihood

$$-\ln \ell(w) = -\sum_{n=1}^N \ln \left(\frac{e^{w_K^T x}}{\sum_{k=1}^K e^{w_k^T x}} \right) \quad \text{--- (3)}$$

now, the error function (or cost function or objective) for softmax regression is given by

$$E(w) = \frac{1}{N} \sum_{k=1}^N -\ln p(c_k)$$

Loss
for
SR

$$E(w) = -\frac{1}{N} \sum_{n=1}^N \ln \left(\frac{e^{w_k^T x}}{\sum_{k=1}^K e^{w_k^T x}} \right) \quad (4)$$

this is the loss for the data D , so we write,

Loss
from
data

$$E_D(w) = -\frac{1}{N} \sum_{n=1}^N \ln \left(\frac{e^{w_k^T x}}{\sum_{k=1}^K e^{w_k^T x}} \right) \quad (5)$$

now, we add the L_2 regularizer term for each of the weight vectors w_k ,

Loss
from
 L_2 regularizer

$$E_w(w) = \frac{\alpha}{2} \sum_{k=1}^K w_k^T w_k \quad (6)$$

$$E = E_D(w) + E_w(w) \quad (7)$$

Total
loss
for
SR

$$E(w) = -\frac{1}{N} \sum_{n=1}^N \ln \left(\frac{e^{w_k^T x}}{\sum_{k=1}^K e^{w_k^T x}} \right) + \frac{\alpha}{2} \sum_{k=1}^K w_k^T w_k$$

(Regularized Softmax Regression)

now calculate the gradient (derivative) of cost,
 $\frac{\partial}{\partial w} E_W(w) = \frac{\partial}{\partial w_j} \frac{1}{2} \sum_{k=1}^K w_k^T w_k$

$$= \frac{1}{2} \frac{\partial}{\partial w_j} [w_1^T w_1 + w_2^T w_2 + \dots + \underbrace{w_j^T w_j}_{\text{only this term survives}} + \dots + w_K^T w_K]$$

$$= \frac{1}{2} \frac{\partial}{\partial w_j} w_j^T w_j \quad (\text{for } j = k\text{th term})$$

$$= \frac{1}{2} \cdot \frac{\partial}{\partial w_j} \|w_j\|^2$$

$$= \frac{1}{2} \cdot 2 w_j$$

$$\frac{\partial E_W}{\partial w_j} = w_j$$

$$\boxed{\frac{\partial E_W}{\partial w_k} = w_k}$$

δa

again, calculate derivative of E_D ,

$$\frac{\partial}{\partial w_j} E_D = -\frac{1}{N} \frac{\partial}{\partial w_j} \sum_{n=1}^N y_n \left(\frac{e^{w_k^T x_n}}{\sum_{k=1}^K e^{w_k^T x_n}} \right)$$

$$= -\frac{1}{N} \sum_n \frac{\partial}{\partial w_j} [w_k^T x_n - y_n \sum_k e^{w_k^T x_n}]$$

$$= -\frac{1}{N} \sum_{n=1}^N \left[\delta_{jk} x_n - \left(\frac{1 \cdot e^{w_k^T x_n}}{\sum_k e^{w_k^T x_n}} \cdot x_n \right) \right]$$

$\rightarrow = p(c_k)$

$$\boxed{\frac{\partial E_D}{\partial w_j} = -\frac{1}{N} \sum_{n=1}^N (\delta_{jk} - p(c_k)) x_n}$$

$$\frac{\partial E_w}{\partial w_k} = \alpha w_k = \alpha w_k^T \quad (\text{write } w \text{ as } w^T)$$

$$\frac{\partial E_D}{\partial w_k} = -\frac{1}{N} \sum_{n=1}^N (\delta_{kn} - p(c_k | x_n)) x_n$$

$$= -\frac{1}{N} \sum_{n=1}^N (\delta_{kn}(x_n) - p(c_k | x_n)) x_n^T$$

$$\text{So, } \nabla_w E = \nabla_w E_D + \nabla_w E_w$$

gradient
of

$$\nabla_w E = -\frac{1}{N} \sum_{n=1}^N (\delta_{kn}(x_n) - p(c_k | x_n)) x_n^T + \alpha w_k^T$$

regularized
softmax
Regression

Q.E.D.

Note

example of data for multiclass Regression

IFAR-10 data

training examples 5 files
test eg 1 file

classes
airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{10K} \end{bmatrix} \begin{matrix} 1024 & 1024 & 1024 \\ R & S & B \\ \text{umpy array of uint8} \end{matrix}$$

$10,000 \times 3072$
 $N \times 1$
 $=(10K, 1)$

$$t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix} \quad t_i \in (c_1, c_2, \dots, c_{20})$$

$10,000, 1$

$$\vec{w} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1,10K} \\ w_{21} & w_{22} & \dots & w_{2,10K} \\ \vdots & \vdots & \ddots & \vdots \\ w_{10K,1} & \dots & \dots & w_{10K,10K} \end{bmatrix} \quad 10K, 10K$$

$$w^T X = (10K, 10K) (10K, 1) = (10K, 1) = t$$

each $w_{ij} \in w_k$ where $w_k = \{w_1, w_2, \dots, w_{10K} = w_k\}$