

RNN with one independent variable.

Steps

which bucket
if dropped at
 $200px$

- ① ✕ Drop at random. at point (p)
- ② ✕ $(p - 300px) \rightarrow abs.$
- ③ ✕ sort. result from least to greatest
- ④ ✕ look at top K records & see which is most common
↓ sort this.

drop pos.

0	300	bucket
150	150	#1
200	100	#2

abs (pos - 300)

300	300
150	150
200	100

bucket

#1
#2
#3

after sort:

325 ; 25
275 ; 25
375 ; 75

If $K=3$
#4 }
#3 }
#4 }
#3 }
K=3
4 → common

Lodash //

$(300 - pt)$

[
[10, 0.5, 16, 1],
[200, 0.5, 16, 4],
[350, 0.5, 16, 4],
[600, 0.5, 16, 5].
]
→

[
[290, 1]
[100, 4] → sort
[50, 4]
[300, 5]
]
to get closest ✓

$\text{slice}(0, K)$

[
[4:1] ← sort by
[4:2] to pairs
[1:1]]

most common number {

$K=3$
50 4
100 4
290 1

Prediction was Bad.

- ① Adjust parameters. (K value after \sqrt{N})
- ② Add more features to explain analysis (like bounciness)
- ③ Change prediction point.
- ④ Maybe drop position doesn't matter.

① Finding ideal K

Record a bunch of data points

split (at Random)

test

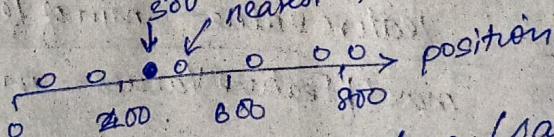
→ training

for each test, run KNN using training data

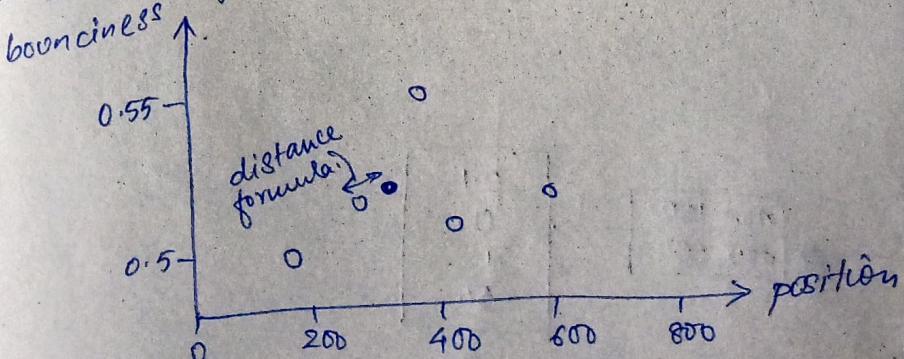
Does the result of KNN equal the 'test'?

② Adding more features.

Initially with only 1 feature it was like. ~~one-D~~



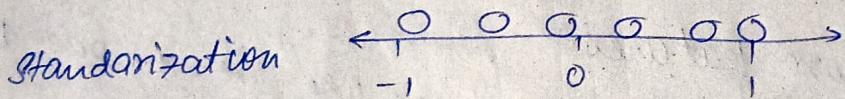
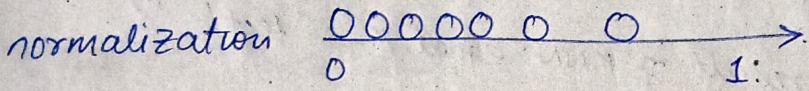
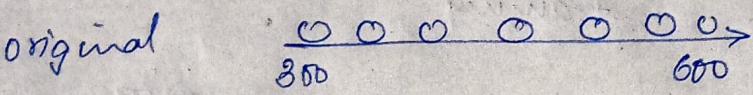
Now considering more than 1 feature. (say 2)



Now in this method of using distance formula
 position range 0 - 700 → $(300-400)^2 = 10000$
 bounciness range 0.5 - 0.55 $(0.55-0.5) = (0.05)^2 = 0.0025$

doesn't cause
any significant
change

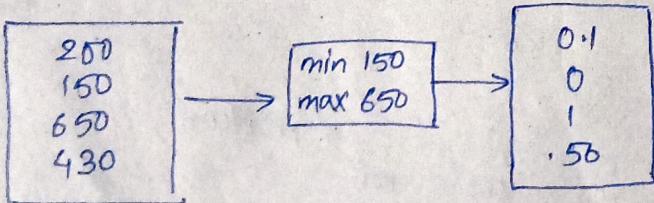
∴ we need to normalise our data
 all numbers range from zero to 1.



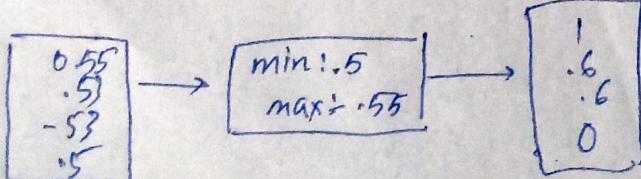
$$\text{Normalised dataset} = \frac{\text{feature value} - \text{min of feature Value}}{\text{max of feature Values} - \text{min of feature Value}}$$

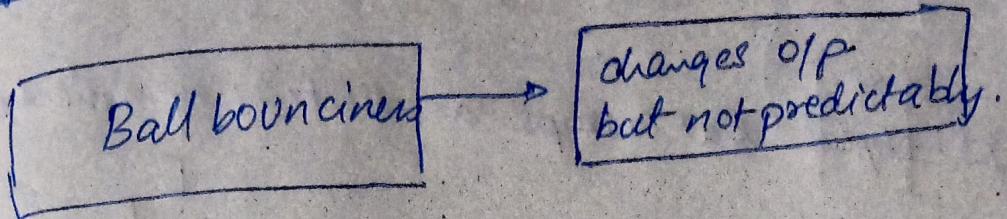
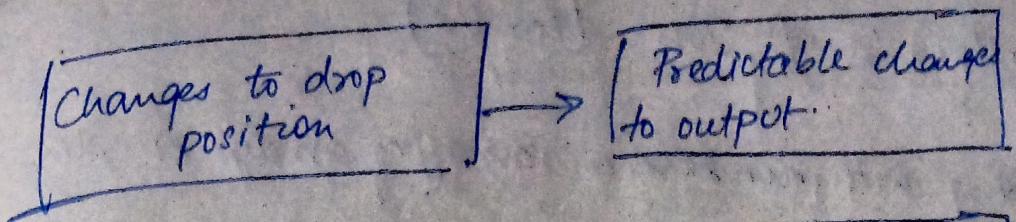
Example.

Drop position



Bounciness



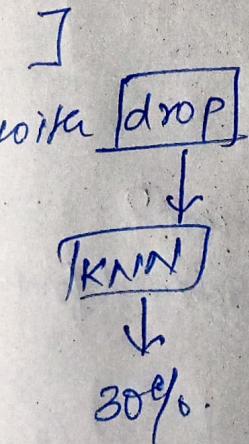


Feature Selection

Deciding which features to include in analysis

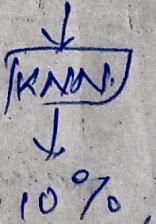
How to select?

[[300, 4]
 [350, 5]
 [416, 9]



[[5, 4]
 [52, 5]
 [53, 4]

[bounciness]



Tensorflow's # Job (For V8)

* make working with numbers in array of arrays simpler
terms.

tensors is object that wraps a collection of numbers.

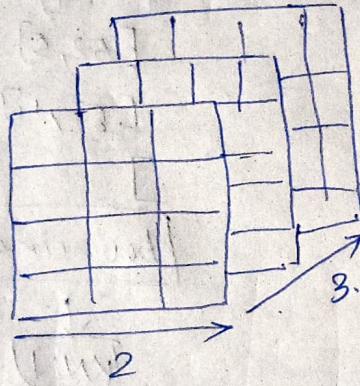
e.g. [[300, 0.5, 16, 4], [10, 2, 6], [5, 10, 15]]
[[350, 0.53, 16, 3], ,]]

2D

1D

3D

3D



shape - How many records in each dimension?
basically length considered from
out to inwards.

[5, 10, 17]

= [3]

[[5, 10, 17],
[18, 4, 2].length,
] .length]

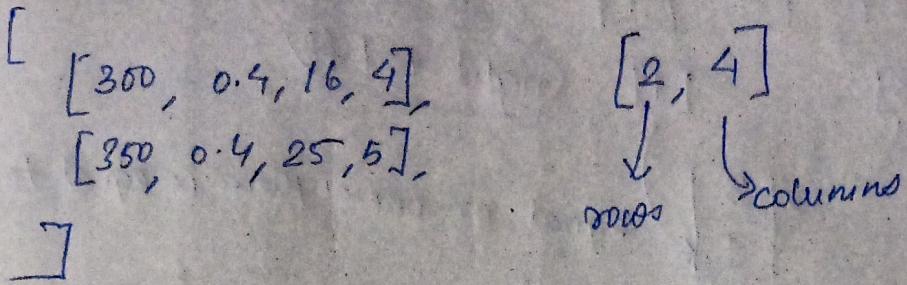
= [2, 3]

[[5, 10, 17].length,
] .length]

, length

= [1, 1, 3]

2D being most common that we will work with.



tensorflow examples:

```
const data = tf.tensor([1, 2, 3]);  
const otherData = tf.tensor([4, 5, 6]);  
data.add(otherData); // [5, 7, 9]
```

identical indices operate.
creates a new data & not alter the original.

$$\begin{array}{r} 128 \\ 456 \\ + \end{array} \quad \frac{}{6, 7, 9}$$

data.sub(bla) for subtraction

• mul()
• div()

"shape needs to match"

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix}_{[3]} + \begin{bmatrix} 4 \end{bmatrix}_{[1]} = \begin{bmatrix} 5 & 6 & 7 \end{bmatrix}_{[3]}$$

This is called broadcasting.

condition :-

from right to left the shapes are equal or ~~size~~
one has '1'

~~$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix}_{[2, 3]} + \begin{bmatrix} 1 \end{bmatrix}_{[2, 1]} = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}_{[2, 3]}$$~~

was '1'

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 4 \end{bmatrix} = \begin{bmatrix} 5 & 6 & 7 & 1 \end{bmatrix}$$

$\begin{bmatrix} 3 \end{bmatrix}$ $\begin{bmatrix} 1 \end{bmatrix}$ right to left

$\begin{bmatrix} 3 \\ 1 \end{bmatrix}$ equal or one of them is equal to 1

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|} \hline & \\ \hline & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array}$$

$\begin{bmatrix} 2, 3 \end{bmatrix}$ $\begin{bmatrix} 2, 1 \end{bmatrix}$

$$\begin{array}{|c|c|} \hline 2 & 8 \\ \hline 2 & 1 \\ \hline \end{array}$$

one, broadcasting possible
same

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|} \hline 1 \\ \hline + \\ \hline 1 \\ \hline \end{array} \checkmark$$

$\begin{bmatrix} 2, 3, 2 \end{bmatrix}$ $\begin{bmatrix} 3, 1 \end{bmatrix}$

will work as long
as others are equal
or has 1

$$\begin{array}{|c|c|c|} \hline 2 & 3 & 2 \\ \hline 3 & 1 & \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} + \begin{array}{|c|} \hline 1 \\ \hline + \\ \hline 1 \\ \hline \end{array} = \text{won't work}$$

$$\begin{array}{|c|c|c|} \hline 2 & 3 & 2 \\ \hline 2 & 1 & \\ \hline \end{array}$$

x ✓

initialising a tensor

const data = ~~to~~ tf.tensor([
[1, 2, 3]
])
data.shape // [1, 3]

to point the data in tensor
~~clg(data.print)~~ X no need to clg.
data.print() \Rightarrow

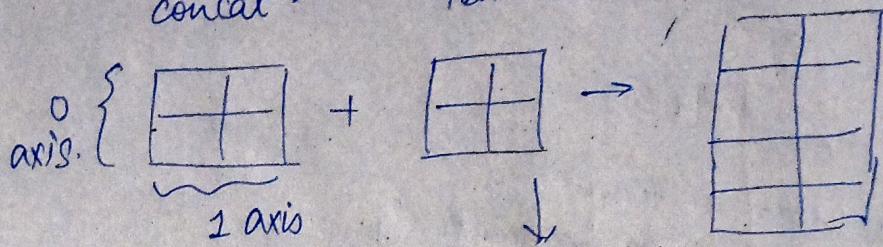
data[0] X data.get(0) $\xrightarrow{\text{1D}}$
 | ↓ index
 | data.get(0, 0) $\xrightarrow{\text{2D}}$
 | ↓ ↓ c

to get a specific column only.

data.slice (startindex, size) $\xrightarrow{\text{not zero-indexed}}$
 | [0, 1] $\xrightarrow{\text{c}}$
 | [6, 1] $\xrightarrow{\text{c}}$ width of slice
 | no rows. like no. of rows.

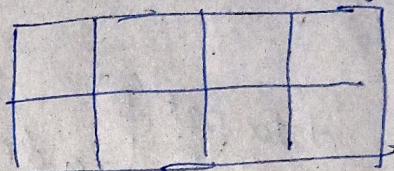
data.shape \rightarrow [rows, col]
or data.slice ([0, 1], $\xrightarrow{\text{total}}$)
or data.slice ([data.shape[0], 1])
~~([0, 1], [-1, 1])~~
-1 gives all

joining tensors.
concat → (identical dimensions)
tensorA.concat(tensorB)



tensorA.concat(tensorB, 1) → default = 0

Solution b.



(data.sum()); sums all elements
data.sum(1); sum of each row → [4]
(0); columns.

gives 1D data

[4, 3] X error.

so, data.sum(1).concat(data2)

second arg of sum, to keep dimensions same

data.sum(1, true).shape → [4, 1]

Alternate way

data.sum(1).expandDims(1) → [1, 4]

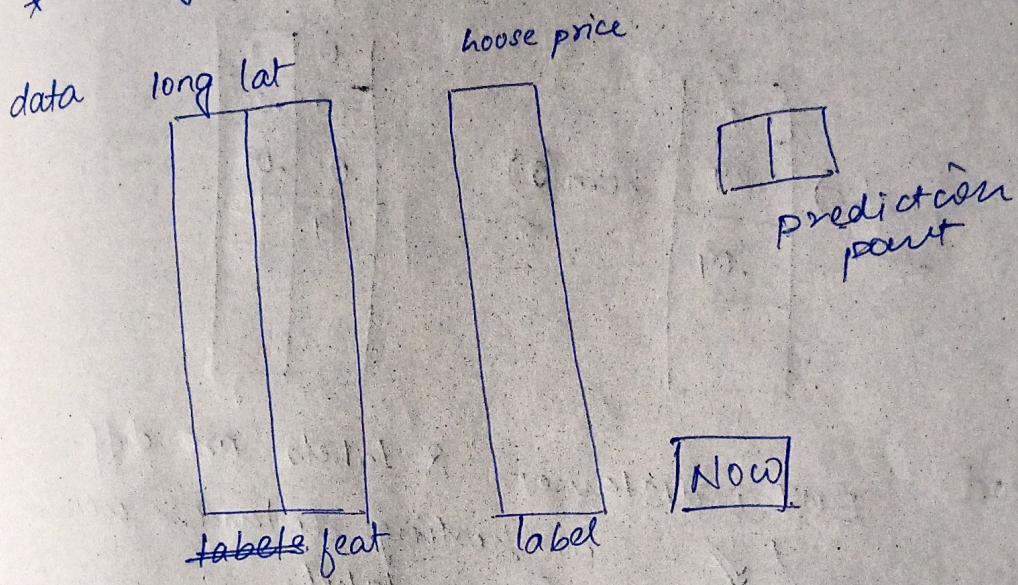
(1) → [4, 1]

Plinko → Classification

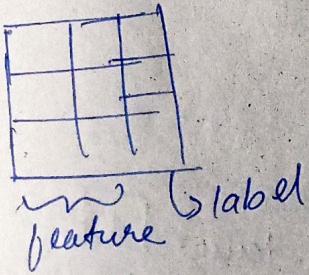
House price → Regression.

KNN algorithm

- * distance b/w features and prediction point
- * sort from lowest point to greatest
- * take the top K records
- * average the label value of those top K records.



64



change in data str

$$\text{dist} = \sqrt{(\text{lat} - \text{lat})^2 + (\text{long} - \text{long})^2}$$

lat	long
-121	47
-121.2	46.5

Sub $\begin{bmatrix} -121 & 47 \\ -121.2 & 46.5 \end{bmatrix}$

0	1
.2	.5
0	3

$\downarrow \text{pow}(2)$

1
.53
.03

1
.29
.09

$\xleftarrow{\text{sum}(a)}$

0	1
.04	.25
0	-.09
1.9	

$\downarrow \text{sort}$

BUT ① distances & labels are diff.
so sorting distances won't
change label.

② tensor can't be sorted.

\downarrow

$[10, 2]$

$\downarrow \downarrow$

1	200
.53	210
.03	250

$\xrightarrow{10}$

unstack

[

tensor

$\begin{bmatrix} 1 & 200 \end{bmatrix}$

tensor

$\begin{bmatrix} .53 & 210 \end{bmatrix}$

]

$\text{sort}((a, b) \Rightarrow a \cdot \text{argsort}([b]))$ returns array
oo use sort//

top K-records.

array \rightarrow slice $(0, K)$

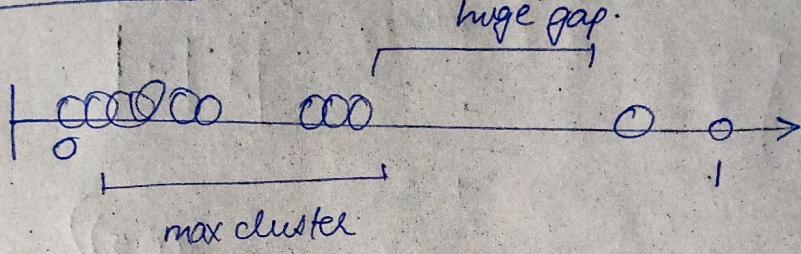
reduce $((acc, ten) \Rightarrow acc + ten \cdot \text{get}(1, 0)) / K //$

~~avg~~
~~stat~~
price

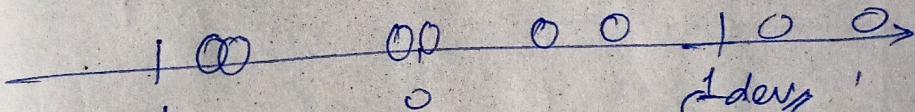
avg of prices

g error $\frac{\%}{}$

Normalization



Standardization



as if accounts for extreme values better
than normalization

Stand

Standardization

$$\frac{\text{value} - \text{Avg}}{\text{Std Deviation}} = \sqrt{\text{variance}}$$

tf. moments (tensor to calc.) → { mean :
variance :

↓
does total (all elements)

tf. moments (data, 0 or 1)
row column

$$0 \left\{ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \right. \quad \begin{array}{c} \overline{+} \\ + \end{array} \quad | + | + |$$

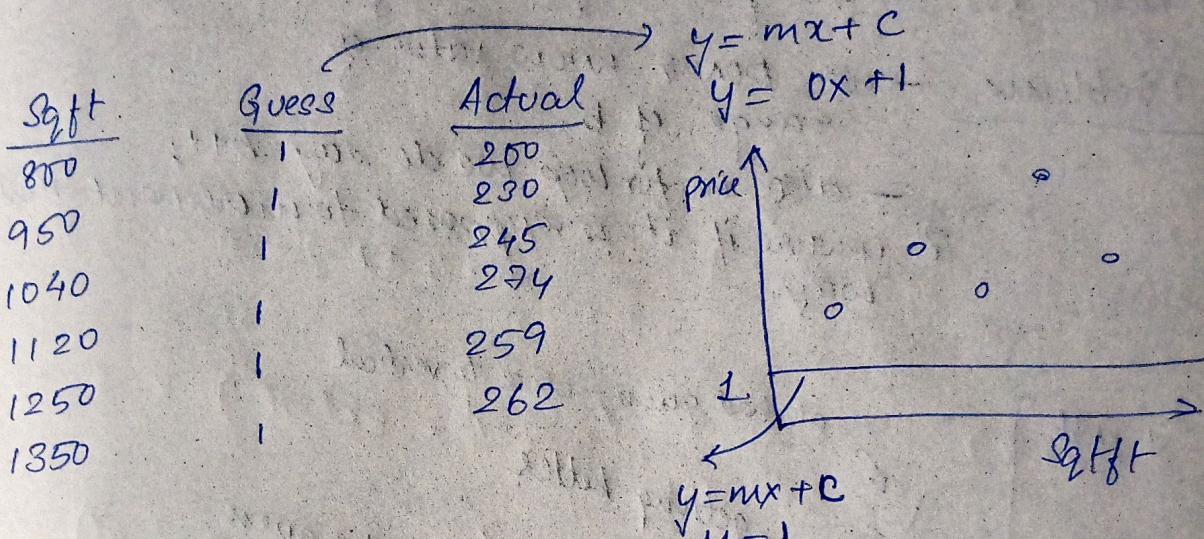
data. sub (mean). div (variance · $1000(0.5)$)

Linear Regression

Create a relation with independent variable & dependent variable

$$\text{mean squared error} = \frac{1}{n} \sum_{i=1}^n (\text{guess}_i - \text{actual}_i)^2$$

↓ tells you about how bad, your equation is



$$y = 0x + 1$$

mean sq error = $\frac{(1-200)^2 + (1-230)^2 + \dots + (1-262)^2}{6}$

MSE = $6013\frac{2}{3}$

New guess = $y = 0x + 200$.

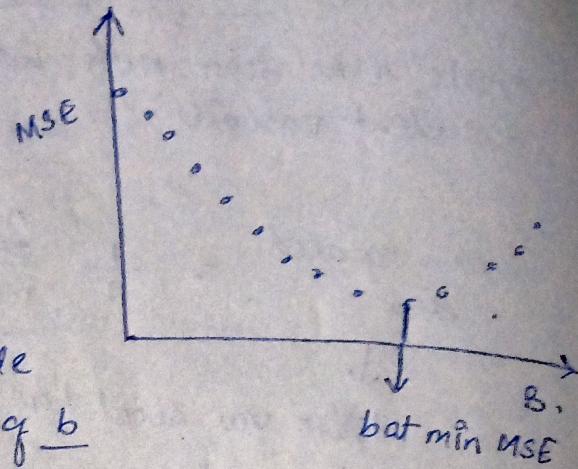
MSE = 4287.

$$y = mx + \frac{b}{j}$$

we need to guess this better

lets focus on guessing B now.

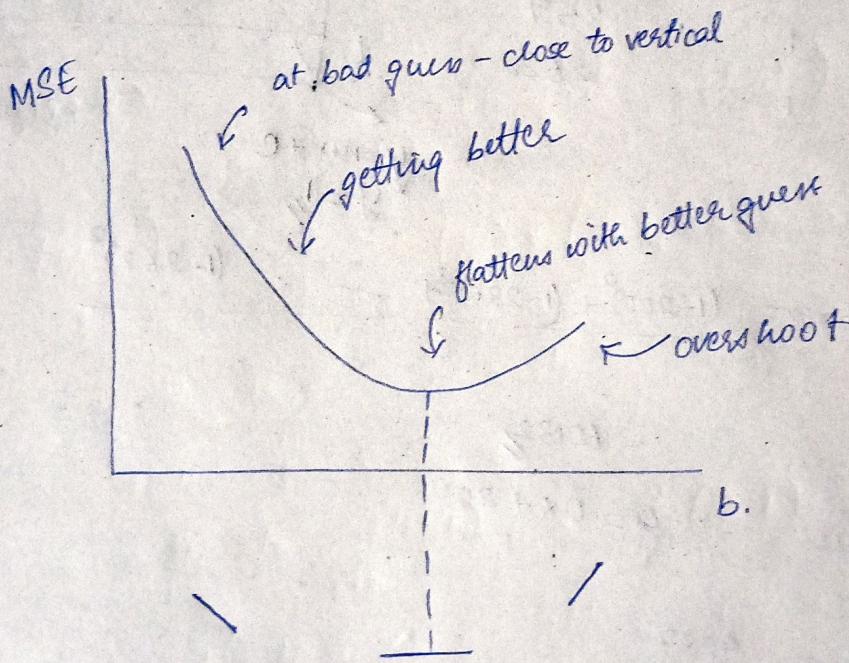
b .	MSE
1.	60132
200	4287
1	
1.	



now using for loop to provide
MSE for range of value of b

Problems :- we DON'T KNOW value of
range of b .

- also in for loop we do say $b+1$;
- 2' unsure if it is supposed to increment that way.



slopes to left
 $'b'$ is less than
optimum

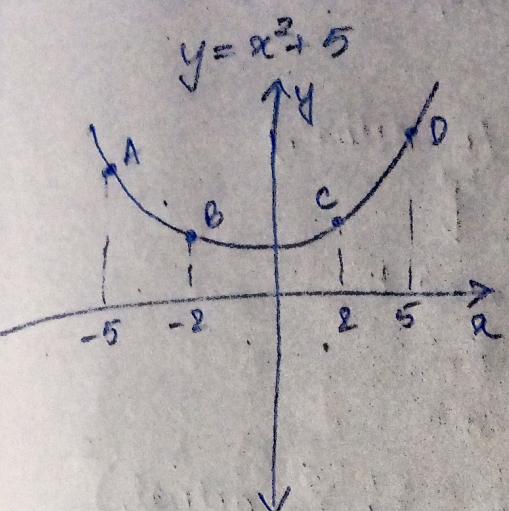
slopes to right
 $'b'$ is more than
optimum

more slope - farther
less slope - closer

therefore slope/rate of change is valuable

Derivatives

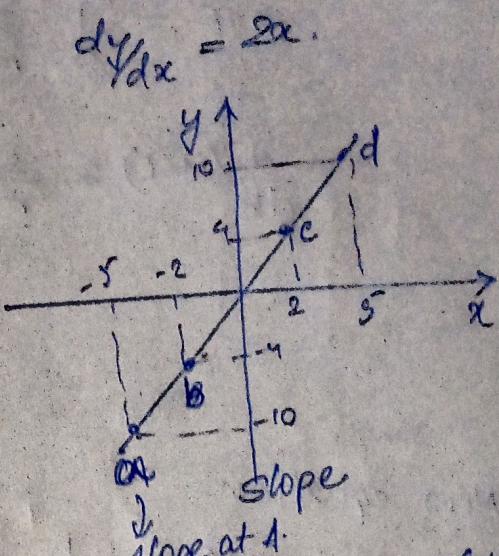
The derivative of an equation gives a new equation that tells us the slope at any location.



$$y - \text{MSE}$$

$$x - b$$

$$\frac{d(\text{MSE})}{db} = \frac{2}{n} \sum (b - \text{actual})$$



$$\text{y-axis-slope} - \frac{d(\text{MSE})}{db}$$

$$\text{x-axis} - a_1, b_1$$

Gradient Descent

- Steps—
- Pick 'b' any value.
 - calculate slope of MSE vs b.
 - if slope is very very small — work done!
 - multiply slope with an arbitrary value called learning rate.
 - sub that from b.

a - learning rate

-ve slope (m)

$b - (a \times m)$

$b + a \times m$
 $b + a/m$
right

$b - a/m$

~~b~~ at

towards left

Now to include m as well.

b alone =

$$MSE = \frac{1}{n} \sum_{i=1}^n (\text{guess}_i - \text{actual}_i)^2$$

$$\frac{d}{db} (MSE) = \frac{2}{n} \sum_{i=1}^n (b - \text{actual}_i)$$

\downarrow
 $mx + b$.
 \downarrow
0.

New equation

$$MSE = \frac{1}{n} \sum_{i=1}^n ((mx_i + b) - \text{actual}_i)^2$$

$$\frac{d}{db} (MSE) = \frac{2}{n} \sum_{i=1}^n (mx_i + b - \text{actual}) \times 1$$

$$\frac{d}{dm} (MSE) = \frac{2}{n} \sum_{i=1}^n (x_i) \times (mx_i + b - \text{actual})$$

Gradient Descent

- Pick b & m .
- calculate $\frac{d(MSE)}{db}$ & $\frac{d(MSE)}{dm}$
- if they are small — we are done
- multiply both with learning rate (α)
- sub results $b - \alpha \frac{d(MSE)}{db}$ & $m - \alpha \frac{d(MSE)}{dm}$

goo.gl/f4KdZj

here say when we apply, there is unchanging value of b , but works fine for me meaning :- scaling is required.

Project . cars.csv.

$$\text{miles per gallon} = m * (\text{car horsepower}) + b.$$

structure

class LinearRegression -
gradientDescent() → run one iteration of GD
train() → and update m & b
test() → run GD till we get good value of b & m
predict() → evaluate accuracy using testcases
make prediction using our m' & b'

1/1
1/6
1/9
1/21

$$\frac{d(\text{MSE})}{db} = \frac{2}{n} \sum_{i=1}^n ((mx_i + b) - \text{actual}_i)$$

$$\frac{d}{dm} (\text{MSE}) = \frac{2}{n} \sum_{i=1}^n \cancel{x_i} \left(\underbrace{(mx_i + b)}_{\text{current guess}} - \text{actual}_i \right)$$

Tensorflow implementation

R.T.O

Matrix Multiplication

shape $\begin{bmatrix} 4 & 2 \end{bmatrix}$ $\begin{bmatrix} 2 & 3 \end{bmatrix}$
must be same

Now how this related to ML

$$\begin{bmatrix} c & d \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \text{jet} * a + - d * b.$$

Now implementation:
the 2 is insignificant
as it can be put in learning rate
as well.
we just need slope

$$\text{Slope of MSE w.r.t both } b \text{ & } m. \quad \frac{\text{features}^T (\text{features} * \text{weights}) - \text{labels}}{n}$$

General eq.

weights: $M \otimes b$ in a tensor

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 1 & 6 \end{bmatrix} \quad \begin{bmatrix} i & 2 \\ 3 & 4 \\ 1 & 6 \end{bmatrix}$$

horsepower

x_1

x_2

x_3

weights

$$\begin{bmatrix} m \\ b \end{bmatrix}$$

Matrix multiply.

$$\begin{bmatrix} 2, 1 \end{bmatrix}$$

some append on
left hand side

a_1 1

a_2 1

a_3 1

\vdots 1

$$\begin{bmatrix} 1 & x \end{bmatrix}$$

equal

$$\begin{bmatrix} m \\ b \end{bmatrix}$$

$$\begin{bmatrix} 6, 2 \end{bmatrix}$$

$$\begin{bmatrix} 2, 1 \end{bmatrix}$$

current Guesses.

$$\begin{bmatrix} a_1 m + b \\ a_2 m + b \\ a_3 m + b \end{bmatrix}$$

$$\begin{bmatrix} 6, 1 \end{bmatrix}$$

$$\begin{array}{l}
 \begin{matrix} x_1m+b \\ x_2m+b \\ x_3m+b \end{matrix} \quad \text{sub actual} \rightarrow \begin{matrix} \text{differences} \\ \text{diff} - 2 \\ [6, 1] \end{matrix} \\
 - \quad (mx_i + b) - \text{actual} \\
 \text{done}
 \end{array}$$

horsepower

x_1

x_2

\vdots

[6, 2]

x

differences!

[6, 1]

$$= \frac{\sum (-)}{\sum x_i}$$

$$\begin{bmatrix}
 x_1 & x_2 & x_3 & \dots & \\
 1 & 1 & 1 & \dots & \\
 \end{bmatrix} \quad [2, 6]$$

transpose

actually

[6, 1]

$$\begin{bmatrix}
 d-1 & 1 \\
 d-2 & \\
 d-3 & \\
 \vdots &
 \end{bmatrix}$$

$$\begin{bmatrix}
 x_1d-1 + x_2d-2 & \\
 d-1 + d-2 & \\
 \end{bmatrix}$$

$$\begin{aligned}
 & - \sum x_i(mx_i + b - \text{actual}) \\
 & - \sum (mx_i + b - \text{actual})
 \end{aligned}$$

gradients or slopes

slopes * learning rate \rightarrow
weights - α

Gauge accuracy using 'coefficient of determination'
R².

IDE
V

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} \rightarrow \begin{matrix} \text{sum of squares residual} \\ \text{sum of squares total} \end{matrix}$$

Higher - the better of test.

$$SS_{tot} = \sum_{i=1}^n (\text{actual} - \text{average})^2$$

$$SS_{res} = \sum_{i=1}^n (\text{actual} - \text{predicted})^2$$

test()

| a_1
| a_2
| a_3
| a_4
|

b
m

$b + m a_i$
our prediction

prediction →

$$SS_{res} = \sum (\text{label} - \text{predicted})^2$$

Bad accuracy

Standardization

* apply same mean & var. in test as well

$$\frac{\text{value} - \text{mean}}{\text{std dev}} \rightarrow \sqrt{\text{variance}}$$

Multivariate-

$$MPG = b + (m_1 * \text{weight}) + (m_2 * \text{displ}) + (m_3 * \text{horsepower})$$

$$y = b + (m_1 * x_1) + (m_2 * x_2) + (m_3 * x_3) - .$$

features' weights disp' horsee

x_1	x_2	x_3
a_1	a_2	a_3
α_1	α_2	
m_1	m_2	

matmul

b
m_1
m_2
m_3

4	1
---	---

6	4
---	---

gives the prediction or guess.

$[6, 1]$

feat^T

x_1	x_1	x_1	x_1	x_1	x_1
x_2	x_2	x_2	x_2	x_2	x_2
x_3					

4, 6

4, 1

$$\overbrace{b + (m_1 * x_1) - .}^{6, 1}$$

Optimising learning rate

methods:-

Adam
RMS Prop
Adagrad
Momentum.

$$\left\{ \begin{array}{l} 155 \\ 145 \\ \hline 300 \end{array} \right.$$

Custom Learning Rate Optimization

Steps:-

- with every iteration of GD, calculate the exact value of MSE and store it.
- 1 iteration later look at MSE_{old} & MSE_{new} .
- if $MSE_{new} > MSE_{old}$ learning rate / 2
- if $MSE_{new} < MSE_{old}$ ↑ LR by 5%
i.e. 1.05 LR

MSE
vectorized mean sq error = $\frac{\text{sum}((\text{features} * \text{weights}) - \text{labels})^2}{n}$

Batch Gradient Descent

- guess starting value of B & M
- calculate slope of MSE using a portion of observations in feature set & current M/B values
- multiply slope by learning rate
- update B & M

Stochastic Gradient Descent

- guess a starting value of B & M
- calculate slope of MSE using 1 obs. in feature set & current M/B values.
- multiply the slope by the learning rate
- update B & M

Gradient Descent
uses entire feature set

Batch Gradient Descent
couple observation

Stochastic G D
ONE observation

b4

train() → gradientDescent()
↓ ↓
features labels

~~train()~~
~~GD~~

features labels
↓ ↓
train() batch of features
 batch of labels → gradientDescent()
carve them to little batches

BGD — size in each batch
— total no. of batches.

here - $GD \rightarrow$ each iteration = entire data set
~~BGD = each iteration~~ is gone through
MSE is taken based on full
then weights are adjusted
HUGE computation
slow.

$BGD \rightarrow$ calculates MSE in batches & weight
adjustments.

it's like 100 pushups in 5 workout
sessions than 100 at one go.

Now predicting.

take an array of arrays

[

[horse, weight, dispalc] } multiprediction
[" " "]

]

then matmul with weights gives the guesses.

Logistic Regression:
Predicts discrete value (classification)

Binary Logistic Regression - Binary classification

Pass / Not Pass

Spam / Not Spam.

Customer accepts / Customer declines

Apple / Android

Q. Based on a person's age, do they prefer
READ BOOKS or WATCH MOVIES.
(RB) (WM)

Sample

Age	Activity
5	watch movies
15	watch movies
25	read books
35	read books
45	read books

Encoded activity

0

0

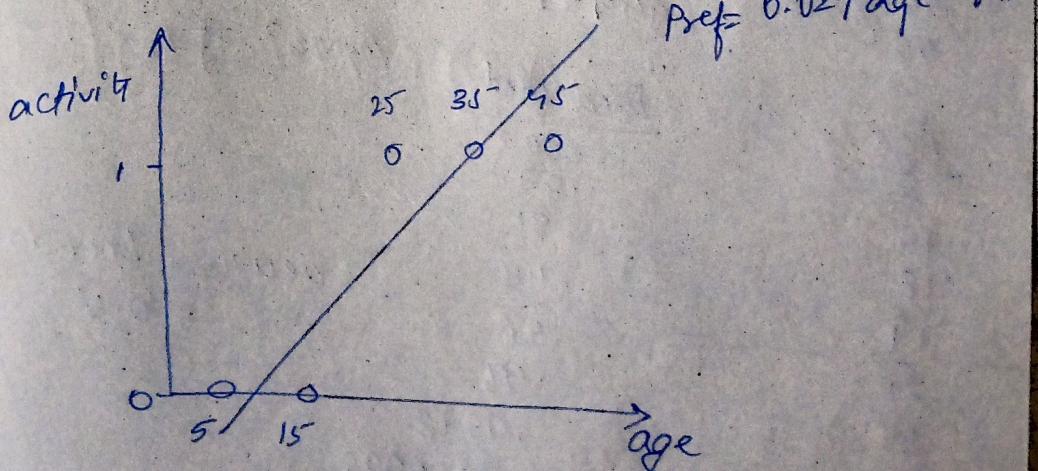
1

1

1

convert this column to
numerical data

Now we formulate an equation



If linear regression is used.

<u>Age</u>	<u>actual</u>	<u>Pred.</u>
5	0	-0.0011 X
15	0	.2969
25	1	.5945
35	1	.8923
45	1	1.19 } X
50	1(?)	1.48 }

results b/10 0 & 1 are ok
outside 0.81 X

The eq. we will use:

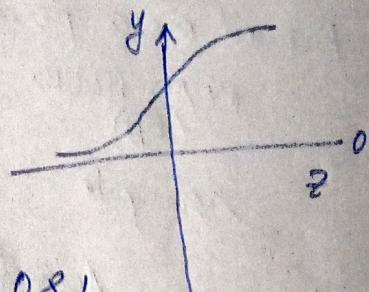
Probability
of preferring books.

$$\frac{1}{1 + e^{-(mx + b)}}$$

$$e = 2.718$$

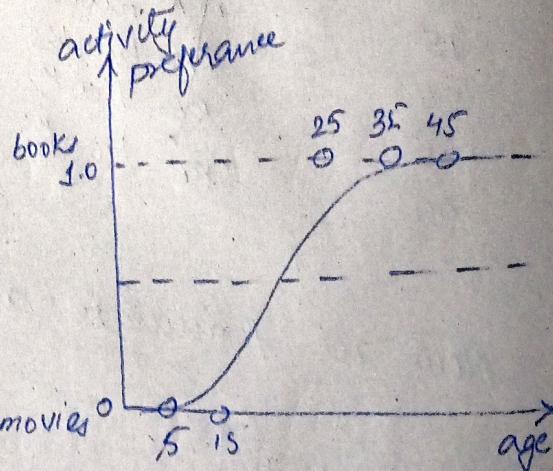
sigmoid eq. =

$$\frac{1}{1 + e^{-z}}$$



output is b/10 0 & 1

<u>Age</u>	<u>predicted</u>
5	.0002
15	.07
25	.96
35	.9998
45	.999999
55	.99999999



<u>Age</u>	<u>Predicted</u>
5	.0002
15	.07
18	.81
19	.44
20	.58
25	.96
35	.9998
45	.999999
55	.99999999

decision boundary
= 0.5

Logistic Regression Gradient Descent

Steps:

- * Encode label values as either '0' or '1'
- * guess M, B
- * cal. slope of MSE using all obs.
& ~~current~~ ^{current} M/B .
- * multiply * learning rate
- * update B and M

Linear
Regression

Given a vehicle's weight, horsepower & engine displacement

PASS or NOT PASS.
smog emissions check

Linear - MSE (metric of how bad we guessed)

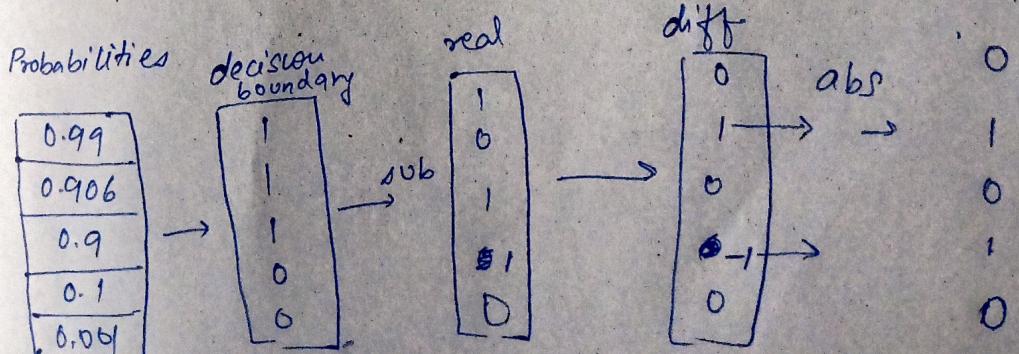
For logistic - Cross Entropy

$$-\frac{1}{n} \sum_{i=0}^n \text{actual} \cdot \log(\text{guess}) + (1 - \text{actual}) \cdot \log(1 - \text{guess})$$

guess = sigmoid($Mx + b$)

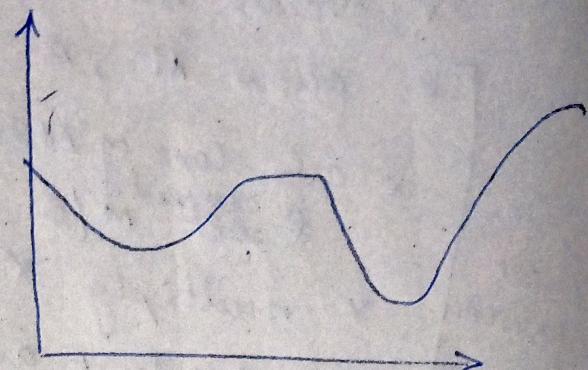
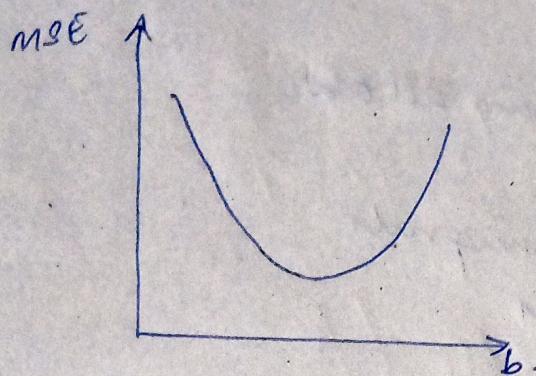
for logistic regres. $R^2 = 1 - \frac{\text{SSres}}{\text{SStot}}$

makes no sense



no. of errors = $\frac{2 \times \text{no. of abs.}}{\text{total tes.}}$

Why not MSE used with sigmoid? although it still gives good accuracy.



not sigmoid
Convex function ✓
One min. value

X
has local minima

and hence, cross entropy

cost History → gives history of ~~C~~ Entropy

Cross Entropy.

$$(-\frac{1}{n}) \sum_{i=0}^n \text{actual} \cdot \log(\text{guess}) + (1-\text{actual}) \cdot \log(1-\text{guess})$$

Vectorized

$$-\frac{1}{m} \cdot (\text{actual} \cdot \log(\text{guess}) + (1-\text{actual})^T \log(1-\text{guess}))$$

Multinomial Logistic Regression.

- * based on hrs driven - Luxury, Sedan, Truck, Compact what car
- * salary bars phone - Apple, Android, Feature, windows
- Person's age :- Books, Dance, Movies.

Method-1

<u>age</u>	<u>Pref. Activity</u>	<u>encoded Movies</u>	<u>encoded Books</u>	<u>encoded Dance</u>
5	movies	1	0	0
15	movies	1	0	0
25	books	0	1	0
35	books	0	1	0
45	books	0	1	1
55	dance	0	0	1
65	dance	0	0	1

Now,

[20 yrs old]

Method-1

diff model for each	Books model	Dance model
0.84	1	0.47

intense for multiple say 20 diff so options

method-2

combine weights

M	B.	D.
b	b	b
m	m	m

combine tensors.

Movies Books Dance

1	0	0
0	1	0
0	1	0
0	1	0
0	0	1
0	0	1

Slope of Cross Entropy

w.r.t M and B

$$= \frac{\text{Features}^T * \text{sigmoid}(\text{Features} * \text{weights}) - \text{labels}}{n}$$

Predictions

weights

1	α_4
1	α_2
1	α_1
1	α_4
1	α_1

PX2

b	b	b
m	m	m
M	B	D

2X3

~~bxnx~~

M B D.

$m_2+b_1, m_2+b_2, m_2+b_3$

PX3

Sigmoid (pred) . sub (labels) \rightarrow [PX3]

b^T

Given: horsepower, weight, displacement

fuel efficiency \rightarrow high, medium or low.

mpg \rightarrow somewhat same as fuel efficiency.

\hookrightarrow this is given as float values. so first convert it to

0 - 15
low.

15 - 30
med.

30+

high

low mid high

[1 1 0]
conditional

what if we get output

marginal probability distribution

\rightarrow considers 1 possible output case in isolation.

conditional probability distribution

\rightarrow considers all possible output cases together

we need marginal probability

what we have

movies

0.34

books

.47

dancing

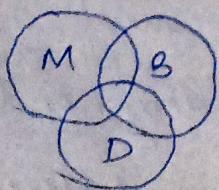
.67

$$0.34 + 0.47 + 0.67 = 1.48 \leftarrow \text{greater than 1}$$

it's like ~~A ∩ B~~ A ∪ B
cases would have been needed to be reduced.

this indicates that

say movies probability is not just ONLY movies
but some intersection of movies & something



$$M + B + D > 1$$

$$M + B + D = 1$$

(like dice)

so it would be nice to have $\sum \text{prob} = 1$

like

ONLY -

M

B

D.

.28

.32

.39

Sigmoid \rightarrow Marginal Probability Distribution

Softmax \rightarrow Conditional Probability Dist.

↳ considers all possible output cases together

Softmax Equation

Probability of being 'Y' label rather than 'O' label

$$\frac{e^{mx+b}}{\sum_{k=0}^K (e^{m\alpha+b_k})}$$

Individually :

<u>movies</u>	$mx+b = 15$
e^{15}	$e^{15} + e^{23} + e^{22}$

<u>books</u>	$mx+b = 23$	<u>dancing</u>
e^{23}	$e^{15} + e^{23} + e^{22}$	e^{22}

Together

.0002

.73

.26

so we need to use softmax on each row.

	M	S	D
softmax	$m \times 3$	$m \times b$	$m \times b$
	0.819	0.422	0.258

but here how do we take which one?
decision boundary = 0.5 (say)
then all = 0.

we need max so we use

$\text{argMax}()$

→ returns the INDEX with
highest value.

test()

L	M	H	(pred)	label	Plan
			col. of max	col. of max	≠
1	0	0	0	0	0
1	0	0	0	0	0
1	1	0	1	not Equal	0 → 0
0	1	0	1	1	0
0	1	0	2	0	0
0	0	1	2	2	1
0	0	1	2	0	0

argMax

sum.

we get
error

Handwritten Digit Recognition.

Given the pixel intensity values in an image identify whether the character is a hand-written 0, 1, 2, 3, 4, 5, 6, 7, 8, 9

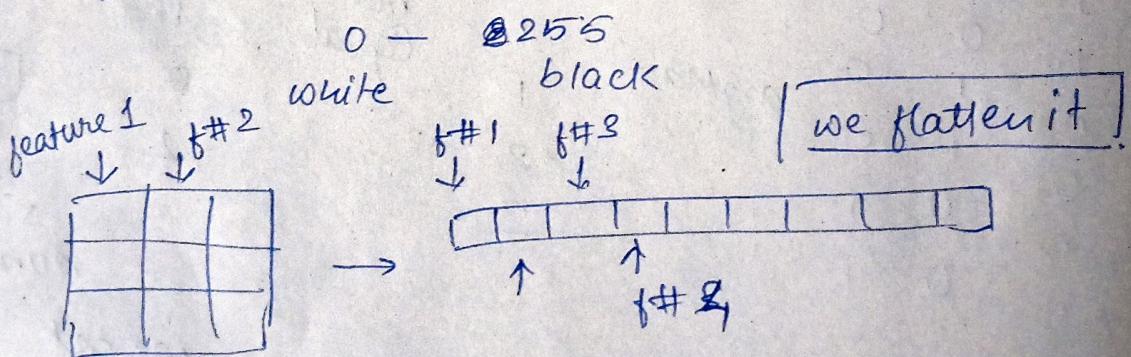
Each image

28 PX

28 pixel

784 pixels.

The input is gray scaled



here we have multiple features

784 features $[P_{x-1} \quad P_{x-2} \quad P_{x-3} \dots]$

b4 \rightarrow [horsepower, displa, weight]

next for each label

we divide them into columns.

en	→	0	1	2	3	4	5	6	7	8	9
7											
4											
2											
9											

Now, the problem

for scaling, we are using standardization.
where formula $\frac{\text{value} - \text{avg}}{\sqrt{\text{variance}}}$ we get variance as zero
in some cases.

$$\frac{\text{value} - \text{avg}}{\sqrt{\text{variance}}} \rightarrow \text{NAN}$$

eq features

$$[[0, 20, 40], [0, 20.5, 80], [0, 20.3, 90]]$$

$$\text{variance} [0, 0.0422, 466.667]$$

$$\text{feat. sub(means)} \cdot \text{div}(\text{vari. pow}(0.5)) = E$$
$$[1, -1.2977, -1.3887]$$

$$\frac{\text{subtraction}}{0} = 1$$

↓
broaden

$$[1, -1.12 - - - - -]$$

Nodejs = NAN

so what we do is

$$\text{variance} \cdot \text{cast}(\text{bool}) \rightarrow [0, 1, 1]$$

$$\cdot \text{logicalNot}() \rightarrow [1, 0, 0]$$

$$\cdot \text{cast}('float32') \rightarrow [1, 0, 0] \rightarrow \text{addable}$$

so basically turns '0' → 1 in variance

another problem:

So cross entropy formula has $\log()$
 $\log(0)$ & $\log(-\alpha) \rightarrow \text{NAN}$.

Trick to solve it

$\log(0)$. add $(-e^{-7})$ $\rightarrow 10^{-7} \neq 0$
so solved

and such a small value won't
make a diff in others.