\_\_\_\_\_/

-> In modern real world supervised learning is applied more & hos vorting of applications.

# Supervised Learning part - 1

- → takes set of inputs which has cossect outputs, leaves from it & will become capable enough to predict the outputs forme w inputs. Example: house price prediction
- -> This gene xally fits under seguession which is a major past of supervised learning another major one is called classification.

# Supervised Learning part -2

- -> Classification algorithms learn basedon the dataset provided & classifies the input into classes or calegories; depending on type of dataset we can have two to any in 12 of classes.
- -> Classification predicts certain first values of output unlike segmension which can predict infinite different possibilities. -> Here we dear some sort of boundary to classify
- -> We can have multiple sets of input too in classification as segmention algorithms. Objects into diff calegories.

## Unsupervised learning Past-I

- $\rightarrow$  Here here are no subput labels, instead the algorithm terms to find the patherns in the dataset.
- -> Es 19 he algorithm decides that three 13 two groups forming allow assessing the dataset, such algorithm is called clustering algorithm , this classifies data into clusters.
- -> The dataset have do not have any labels, the job of algorithm reto find patterns as group the datasets which are sender.

### Un supervised learning port-2 (USL)

(USE.) → Another type of USL is called aromally detection, which detects when anything unusual occurs Es fraud detection algorithms detect & flag very large & supprising → - - - - - dimensionality seduction, it take a big defined & compresses with minimal data loss.

## Jupyter notebook

-> Take a look at them when possible.

## Lineas sequession model past-1

$$\frac{S_{ije} in feet^2}{(\chi^{i})} = i^{th} \frac{1}{1} \frac{1}{$$

## Lineas segression model part-2

 $\rightarrow f_{wb}(x) = wx+b$ 

w, b > parsameters or coefficients or weights

 $\begin{array}{c} & & & \\ & & & & \\ & & & \\ & & & & \\ & & & & \\ & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ &$ 

ું-ને ક**લ્ક્ર**શ્વ

$$\frac{\left(0,0\right)^{2} \frac{1}{\mathcal{L}_{m}} \sum_{i=1}^{m} \left(3^{(i)} - y^{(i)}\right)^{2}}{9^{(i)} \cdot f_{min}\left(x^{(i)}\right)}$$

Cast function intuition

 $\rightarrow$  Greal: minimize  $\mathcal{J}(\omega, b)$  is using to write it.

$$\frac{\sum_{i=1}^{n} k_{i} k_{i}}{\int_{\omega} (x) = \omega \cdot x \quad i \in b + 0}$$

$$\frac{J(\omega) - \frac{1}{2m} \sum_{i=1}^{m} \left[ \int_{\omega} (x^{(i)}) - y^{(i)} \right]^{2}}{\prod_{i=1}^{m} \int_{\omega} (\omega \cdot x^{(i)} - y^{(i)})^{2}}$$

$$\frac{J(\omega) - \frac{1}{2m} \sum_{i=1}^{m} (\omega \cdot x^{(i)} - y^{(i)})^{2}}{\sum_{i=1}^{n} \int_{\omega} (\omega \cdot x^{(i)} - y^{(i)})^{2}}$$

$$\frac{J(\omega) - \frac{1}{2m} \sum_{i=1}^{m} (\omega \cdot x^{(i)} - y^{(i)})^{2}}{\sum_{i=1}^{n} (\omega \cdot x^{(i)} - y^{(i)})^{2}}$$

Visualising cost function

 $\rightarrow$  Back to original model where  $b \neq 0$ 

-b 32D plot fox 2 vourable ie wyb Gwaph looks somewhat likea soupboud



- -> this gas on becoming very complex for further increase in features, hence we use some thing called contour plat, this has ellipses as evals in it allat some height in 3.0 plat.
- The points on a portrialor ellipse sepsesent diff values of w&b (hence they are diff eq. ") but have some value of J.
- $\rightarrow$  Imagine contours as top view of mountain, which is alread.



## Vaining Lineax seq session - Gradient descent

- $\rightarrow \text{Greadient descent unables for a function with many parameters too } \underbrace{\text{Fr}}_{\substack{\omega_1,\dots,\omega_n,b}} \overline{\mathcal{T}}\left(\omega_1,\omega_2,\dots,\omega_n,b\right)$ 
  - We short with some value of wand b
  - Generally set than to 0.0
- $\rightarrow$  We keep changing w,b until we set the at ox near a minima. <u>Nobe</u> : function may have move than 1 minimum.
- -> Function is always bow as harmonic shaped in equated every cast func. but if we train a cost func, based on reason returned networks the graph will be very complex.
  - I soduent descent chooses the path of steepest descent. until we seach at as near a local minima.

Implementing goodient descent

-> w & b are simultaneously updated

Learning sate (x)

 $\rightarrow$  If  $\alpha$  is too small, model works too slowly hence it decreases the efficiency of the model

ightarrow ig

-> what happens when we are already at a minima?

the values further don't update as  $\frac{\partial J}{\partial T} = 0$ 

 $\rightarrow$  As we approach minima gradient descentration becomes smaller steps until minima is attained (as derivative value gress on d horse steps becomes smaller)  $\rightarrow$ : Gradient descent reaches minimum even with constant  $\alpha$ .

Saining Lineax segression & Gradient descent for Lineax regression

$$J(\omega,b) = \frac{1}{2m} \sum_{i=1}^{m} (f_{\omega,b}(z^{(i)}) - y^{(i)})$$

$$\begin{split} & \omega \colon \omega - \alpha \left( \frac{\partial}{\partial \omega} \mathcal{J}(\omega, b) \right) \implies \frac{\partial}{\partial \omega} \mathcal{J}(\omega, b) \coloneqq \frac{1}{m} \sum_{i=1}^{m} \left( f_{\omega, b}(x^{(i)}) - g^{(i)} \right) x^{(i)} \\ & b \coloneqq b = \alpha \left( \frac{\partial}{\partial \omega} \mathcal{J}(\omega, b) \right) \implies \lambda = \mathcal{J}(u) = 1 \quad \sum_{i=1}^{m} \left( f_{\omega, b}(x^{(i)}) - g^{(i)} \right) x^{(i)} \\ & b \coloneqq b = \alpha \left( \frac{\partial}{\partial \omega} \mathcal{J}(\omega, b) \right) \implies \lambda = \mathcal{J}(u) = 1 \quad \sum_{i=1}^{m} \left( f_{\omega, b}(x^{(i)}) - g^{(i)} \right) x^{(i)} \\ & b \coloneqq b = \alpha \left( \frac{\partial}{\partial \omega} \mathcal{J}(\omega, b) \right) \implies \lambda = \mathcal{J}(u) = 1 \quad \sum_{i=1}^{m} \left( f_{\omega, b}(x^{(i)}) - g^{(i)} \right) x^{(i)} \\ & b \coloneqq b = \alpha \left( \frac{\partial}{\partial \omega} \mathcal{J}(\omega, b) \right) \qquad \lambda = \mathcal{J}(u) = 0 \quad \lambda = 0$$

$$b = b - \alpha \frac{\partial b}{\partial b} \qquad \Rightarrow \frac{\partial b}{\partial b} \frac{\partial (w, b)}{\partial b} = \frac{\partial b}{\partial b} \frac{\partial (w, b)}{\partial b} = \frac{1}{1} \sum_{i = 1}^{i} (f_{w, b}(x^{(i)}) - y^{(i)})$$

The type of function produced in squared environ is called convex function is that convex function is the sonly global minima is to lovel minima function. Further of the sonly global minima is a lovel minima function of the sonly global minima is a lovel minima function.

 $\rightarrow$  Gradient descent on a house price predictor with fit datapoints,

this is called batch gradient descent as each shep of the descent uses all of the training examples.

These are other kinds of segression too which need not take the entire batch for each sky

Linear regression with multiple vouvables - Multiple features

→ Multiple features & Ze, Ze, ..., Ze, ..., Ze, Zy = j<sup>th</sup> feature n: total number of fatures Z<sup>(1)</sup>, features of i<sup>th</sup> training example Lists a list of i<sup>th</sup> numbers, its a rective which includes all the fatures of i<sup>th</sup> training example

X; > value of jth feature of ith training example

Example

 $\rightarrow f_{\mathfrak{B}_{b}}(\tilde{x}) \models \omega, x, + \omega_{2}x_{2} + \omega_{3}x_{3} + ... + \omega_{n}x_{n} + b$ 

 $\overline{\omega} = \begin{bmatrix} \omega_1 & \omega_2 & \dots & \omega_n \end{bmatrix} \longrightarrow \text{ sow vector of all weights}$ 

 $\overline{\mathbf{x}} \cdot [\mathbf{x}, \mathbf{x}_{\mathbf{z}} \dots \mathbf{x}_{\mathbf{n}}]$ 

 $f_{\vec{x},b}(\vec{x}) = \vec{w}_{\vec{x}} \vec{x} + b$ 

dot product of vectores

This is called multiple linear regression.

(not multivariate seg ression) \*\*\*

Vectorization, Past - I

 $\rightarrow$  This makes cade shorter 4 more efficient, uses madern numerical linear algebra libraries, uses GPU.

$$\begin{array}{c} \overrightarrow{u^{2}} \cdot \begin{bmatrix} w_{0} & w_{1} & \cdots & w_{n-1} \end{bmatrix} & \overrightarrow{u^{2}} \cdot \begin{bmatrix} w_{0} & x_{1} & \cdots & x_{n-1} \end{bmatrix} & \overrightarrow{u^{2}} & \overrightarrow{u^{2}} & \overrightarrow{u^{2}} \\ \overrightarrow{x^{2}} \cdot \begin{bmatrix} x_{0} & x_{1} & \cdots & x_{n-1} \end{bmatrix} & \overrightarrow{y^{2}} & \overrightarrow{u^{2}} \\ \begin{pmatrix} 0 & inducing & for coding \end{pmatrix} \\ \begin{pmatrix} 0 & inducing & for coding \end{pmatrix} \\ \begin{pmatrix} 0 & inducing & for coding \end{pmatrix} \\ \hline{f = 0} & (include & Numpy) \\ \overrightarrow{f = 0} & (include & Numpy) \\ for & j & in & wange & (o, n) \\ \hline{f or } f \cdot f + w [j] \cdot x[j] \\ f \cdot f + b \end{array}$$

Using vectorisation

f. np. dot (w, 2) + b ~~ when data set & larger this sans much more farler than the loop

Vectorization, part-2

-> Why is vectorization faster?

in Loop each element is account separately multiplied, added but in neclosisation onlise arrays are account of multiplied at only then computer performs very efform addition. Introduction is very efforted of the computer performance of the computer perfo

Gisadient descent for multiple linear regression

$$\overrightarrow{u} \cdot [\omega_{1} \dots \omega_{n}]$$

$$f_{3k} \cdot \overrightarrow{u} \cdot \overrightarrow{\tau} + b$$

$$J(\overrightarrow{u}, k)$$

$$\text{We rat } \begin{cases} \\ \omega_{1} = \omega_{1} - \alpha \cdot \frac{\delta}{2} \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) \end{cases}$$

$$b = b - \alpha \cdot \frac{\delta}{2} \cdot J(\overrightarrow{u}, k)$$

$$f_{2} \cdot (f_{2} \cdot k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) \\ \frac{\delta}{2} \cdot U(\overrightarrow{u}, k) - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k)$$

$$b \cdot b - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k) - J(\overrightarrow{u}, k)$$

$$b \cdot b - \alpha \cdot \frac{1}{2} \cdot \sum_{i=1}^{n} (f_{2} \cdot k) \cdot J(\overrightarrow{u}, k)$$

-> Alternative way of finding w, b:

losmal cai

- works only for linear segression

+ Solves for us b without Herations

Disadvantages

+ ion't generalised to other learning algorithms

→ is also when number of features are larger than 10,000

more complex ML disbusives use normal equ<sup>ar</sup> in backend (generally).

## Feature scaling port -1

ightarrow This makes gradient descent much more faster.

Example

Рейску: щ ч,+ щ ч, + b 2, 0 sige = 300-2000 safeet 5 lange Kange 2, 5 жет 6 = 0-5 5 small .

the larges the size of feature, lower will be the value of its weight 6 vice versa.



5-0

The ellipses are so small the descent bounces a lot before reaching minima, hence we need some sort of a hack.

-> Make both x, & x\_ comparable is make some charges to units & vary both x, & x. & y in range of 0 to 1, then the plats begin to look like this:



Eature Scaling part -2

Mean normalisation

weld \_\_\_\_\_

Suppose

0 < x, 1 2000 0 4 x 2 3

Mean of all obs x, = 4, Mean of all obs x = 42

lethod-

-0.18 = x, 60.82 -0.16 & x2 < 0.54

$$\xrightarrow{\int_{0}^{1} \frac{1}{2} j_{1}} \xrightarrow{\chi_{1}} \alpha_{1}$$

$$\frac{300 \le x_1 \le 2000}{x_1 \le 2000}$$

ev bell shaped curve (some as 
$$\sigma$$
)  
buttom  
coop  $0 \in \mathbb{Z}_{n} \notin 5$   
 $\frac{\mu_{1}}{\sigma_{2}}$   $\mathbb{Z}_{2}$  model =  $\frac{\mathbb{Z}_{2} - \mu_{2}}{\sigma_{2}}$   $\xrightarrow{-3}$   $(1 \to 1)^{2} \to 3$   
 $-16 \in \mathbb{Z}_{n} \in 19$ 

Checking a sodient descent for convergence

After each idention I must develope , if J P even in L step it means a is chosen posely as there is some buy in the case

\* 500 . Joseans to flotten that means its no longer de creasing hence J has more as loss converged

There is no telling after how many sterrations gradient descent converges

-> Automatic convergence lest

et E be very small number like 10<sup>-3</sup>

if  $\mathcal{T}(\vec{u}, b)$  decreases by  $\leq \mathcal{E}$  in one iteration than declare convergence

Diado: finding constration for E is very hand have many perfor Learning curves as we can tell filter ica mittake or cale or not too in that method

Choosing learning sate

 $\rightarrow$  In leasning curve if the graph 1 than  $\alpha$  is not chosen perpendy. (or bug in code)

 $\rightarrow$  If Learning curve look like this, then either disvery high or extent in code possibly  $-\omega_{*}\omega_{\pm}\alpha d$ 

 $\rightarrow$  Even when  $\alpha$  is small graph t then there is bug in case

 $\rightarrow$  Tay with many  $\alpha$ , generally  $3\pi \, \theta \, perv \, \alpha \, lihe 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1.$ 

Feature engineering

->its using intuition to design new fatures, by transforming or combining original fatures.

3, 5 (x) \* w, x, + w<sub>2</sub> x<sub>2</sub>+b x, » frontage of plot x, » depth ...

nos na save of plot

f J, b ( x) + w1 2, + w2 x2 + w8 23+b

A sea helps predict price more accurately, here we create a new feature

### Polynomial seq session

-> we may need a curve which file datawt purpady as their is no proper ob line, we choose diff polynomial degree age to do this -

Quad value func & after a point (is after maxima) 171 if we stor a func which only I use which or y, what we want to

-> Feature scaling becomes seally important here : the powers are varying.

## Module - 3

Mo tivations

## -> Classification

Binary classification ; out comes have only 2 possibilities

- Why linear seguession cannot be used to classify ?

\* Adding a training shifts the devision boundary which is not desixable.

Logitalis descense		5(≅ ) _ ^
- No. 6 like our	> ≈∞0 ⇒ g(₹)≈1	
$\rightarrow$ The function used have is called example lumbion or basichic function of $\sigma(\pi)$ .	Z=0 + g(Z)=0.5	
and the start and a second solution the start of the function of the	$\frac{1}{1+e^{-x}} \rightarrow \pi (x) \pi (x) \pi (x)$	
	2 and 1963	<>



Es for turner detection we can chose 0.2 as threshold such that above 0.2 is malignant below is benign, this reduces the chose of false positions

# Cast function for logistic regression

-> E having having m training examples with n flatures with two tangets > 06.1 > (fermos delection)

Squard ever cat forc. 
$$\gg T(\vec{w}, \mathbf{k}) = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} \left( f \vec{w}, \mathbf{k}^{(2^{i})} - y^{(1)} \right)^2$$

for linear regression

peoplatus too many local minimus hance this is not a good cast function for logistic seguession.



-Overall cost function will be convex when this is implemented

$$\begin{bmatrix} \left(\vec{\omega}_{3}, \mathbf{b}\right) &= \frac{1}{m} \sum_{i=1}^{m} L\left(f_{\mathbf{a},\mathbf{b}}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}\right) \\ L\left(f_{\mathbf{a},\mathbf{b}}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}\right) &= \begin{cases} -\log\left(f_{\mathbf{a},\mathbf{b}}(\mathbf{x}^{(i)})\right) & \text{if } \mathbf{y}^{(i)} \\ -\log\left(1 - f_{\mathbf{a},\mathbf{b}}(\mathbf{x}^{(i)})\right) & \text{if } \mathbf{y}^{(i)} \end{bmatrix}$$

Simplified cost func. for logistic segression

-> The above loss func can be written as follows:

$$L\left(f_{a,b}\left(x^{(1)}\right)y^{(1)}\right) = -y^{(1)}\log\left(f_{a,b}\left(x^{(2)}\right)\right) - \left(1-y^{(2)}\right)\log\left(1-f_{a,b}\left(x^{(2)}\right)\right)$$

How did we end up with this complex less func. Is its based on something called <u>maximum likelihood estimation</u> which is part of statistics  $T(\vec{w}, b) = -\frac{1}{m} \sum_{t=1}^{m} \left[ g^{(t)} \log (f_{2,t}(x^{(t)}) + (\cdot \cdot y^{(t)}) \log (1 - f_{2,t}(x^{(t)})) \right]$ This cast func. is now convex & has only one single global minima.

$$\begin{array}{c} \omega_{3}^{*} \cdot \omega_{3}^{*} - \kappa \xrightarrow{3} \mathcal{J}(\vec{\sigma}, \mathbf{k}) \\ b = b - \kappa \xrightarrow{3} \frac{1}{\partial b_{3}} \mathcal{J}(\vec{\sigma}, \mathbf{k}) \\ \xrightarrow{d} \mathcal{J}(\vec{\sigma}, \mathbf{k}) = \frac{1}{m} \sum_{i=1}^{m} \begin{pmatrix} f_{i,\mathbf{k}}(x^{\alpha}) - g^{\alpha} \end{pmatrix} \mathbf{r}_{i}^{(\beta)} \\ \xrightarrow{d} \mathcal{J}_{i}(\vec{\sigma}, \mathbf{k}) = \frac{1}{m} \sum_{i=1}^{m} \begin{pmatrix} f_{i,\mathbf{k}}(x^{\alpha}) - g^{\alpha} \end{pmatrix} \mathbf{r}_{i}^{(\beta)} \\ \xrightarrow{d} \mathbf{b} \mathcal{J}(\vec{\sigma}, \mathbf{k}) = \frac{1}{m} \sum_{i=1}^{m} \begin{pmatrix} f_{i,\mathbf{k}}(x^{\alpha}) - g^{\alpha} \end{pmatrix} \end{array}$$

-> this also can be monitored i.e analyzed with learning curves ightarrow Vectorization, flature scaling is applicable here too.

Size

=> wordesefits the curve as has high bias

biase the program thinks that the curve is a st line is touts but any st line will not fit this model, honce the model was brased is has high bias

⇒ fits the dataset poetby well, this is called generalisation, it has restlar high bias nor high remience



\$ fit the dataset exactly, cast funce is zero but this will not do good job in prediction with hist cases, the model is overfit, it has high variance

03 even of one dataparint changes slightly the graph will completly change in order to fit that too, hence it has high vanience

-> Collect mose training examples

-> Select features to include as exclude  $\Rightarrow$  this is called feature selection  $\bigcirc$ including too many fatures may also cause over fitting disade: useful feature might be last

-> Regularisation

Reduce the size of pasameters (w;)



f(a)= 132-0.23 22+0.000014 23

segularization includes all chosen katures but prevents katures from having an overly large effect

Note: whether we change b or not makes very less difference practically

 $\rightarrow \text{ Lato say its } f(x) = \omega_1 x_1 + \omega_2 x_2^2 + \omega_3 x_3^3 + \omega_4 x_4^6 + b$ to prevent this from overfitting we must penalize us & w,

we can do this by,

 $\min_{m \geq 1} \sum_{2m}^{m} \left[ f(\mathbf{x}^{(n)}) - g^{(n)} \right]^2 + 1000 \omega_s^{-1} + 1000 \omega_h^{-1}$ 

1000 cg its large num, hence wasting, would almost be zero.

the base idea a that that if the weights one small its like having a simple madel with ferror fatures hence is los prone to availit

But in a model with many fatures we do not know which to penaltyse ( which is not have we penaltyse all fatures is all up parameters, this makes madel correctly fit the data

$$\Gamma\left(\vec{\omega}, b\right) = \frac{1}{2m} \sum_{i=1}^{m} \left[ f(x^{in}) - g^{in} \right]^{2} + \frac{\lambda}{2m} \sum_{j=1}^{n} \omega_{j}^{2}$$

mean squared ersor X => segularization parameter 2>0

 $\lambda$  is scaled by multiplying by  $rac{1}{2m}$  cog this makes it easy to change  $\lambda$ 

Some include  $\frac{\lambda}{2}$  be also, but it makes very little difference if we penalize box not

I's value decrebs the solative tradeoff between minimizing cast & minimizing assights



Desivation of pastial desire of cast func of logistic seguession

 $\mathcal{J}(\vec{\omega}, b) = -\frac{1}{m} \sum_{f=1}^{m} \left[ g^{(i)} \log(f_{\vec{\sigma}, b}(\mathbf{x}^{0})) + (1 - g^{(i)}) \log(1 - f_{\vec{\sigma}, b}(\mathbf{x}^{0})) \right]$ 

$$\begin{aligned} \int_{1}^{1} \cdot \frac{1}{m} \left[ y \log(f(x)) + (1-y) \log(1-f(x)) \right] \\ \frac{\partial}{\partial \omega} - \frac{1}{m} \left[ \frac{Q}{f(x)} \cdot \frac{Q}{h^{\alpha}} + \frac{Q}$$

$$J = -\frac{1}{m} \left[ y \log(f(x)) + (1-y) \log(1-f(x)) \right]$$

$$\frac{\partial W}{\partial t} = -\frac{1}{t} \sum_{i=1}^{\infty} \left[ \left[ y^{(1)} - f_{id_i} - f_{id_i} \right]_{i=1}^{2} + \frac{1}{t} +$$

# Kegulowied logistic sequession

 $\rightarrow$  in general when we brain the model with let of features than there is a higher with of over fitting

$$\int \left(\vec{w}_{j,k}\right)_{s} = \frac{1}{m} \sum_{l=1}^{s} \left[ g^{(s)} \log \left( f_{\vec{w}_{j,k}} \left( \vec{x}^{(s)} \right) \right) + \left( 1 - g^{(s)} \right) \left( \log \left( 1 - f_{\vec{w}_{j,k}} \left( x^{(s)} \right) \right) \right) \right] + \frac{\lambda}{2m} \sum_{l=1}^{m} \sum_{l=1}^{m} \omega_{l}^{l}$$

$$\frac{\partial}{\partial \omega_{i}} \mathcal{J}(\omega, \mathbf{k}) = \frac{1}{m} \sum_{i=1}^{m} \left( f_{\omega, \mathbf{k}}(\mathbf{x}_{i}) - g_{ij}\right) \mathbf{x}_{ij}^{(i)} + \frac{1}{m} \omega_{i}^{(i)}$$

Second: neural & deep

x<sub>o</sub><sup>(1)</sup> 1 ⇒ taken by us genesally

# Mound on home to : an analytical method to solve for O

J

 ightarrow if dataset is very large then each iteration to modify value takes O(n) time to process through the whole dataset.

shehashi gradint newa converges as it amends its values farench home, but when we have a very large dataset ikis fine as it souks very new to minimum & itsisonly is applied instead of botch gradient descent.

$$= A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

$$\nabla_{A} f(A) \approx \left( \begin{array}{c} \frac{\partial f}{\partial A_{u}} & \frac{\partial f}{\partial A_{u}} \\ \frac{\partial f}{\partial A_{s_{1}}} & \frac{\partial f}{\partial A_{s_{2}}} \end{array} \right)$$

$$\begin{cases} f(\mathbf{A}) = \mathbf{A}_{11} + \mathbf{A}_{12}^{T} \text{ any thing like this} \\ \frac{1}{2} \mathbf{y} \mathbf{A} = \frac{1}{2} \mathbf{x} (\mathbf{A}) = \sum_{i=1}^{n} \mathbf{A}_{ii} \\ \xrightarrow{\mathbf{n}_{ii}} \\ \xrightarrow{\mathbf{n}_{$$

$$\rightarrow A^{T}A = \sum_{i} A^{2} (sum of sq. of ell.)$$

Xy what if X is non invertible? it means we have steller dant features ine features which are linearly depen -dent. This can be solved by: (1) Psuedo inverse (1) Exembing which of the fatureser in extract linearly dependent.

VJ(0) ⇒ sepsesents desivative of J(0) OER"+1

$$\begin{array}{c} \displaystyle \bigvee_{\boldsymbol{\mathcal{G}}} \ \boldsymbol{\mathcal{J}}(\boldsymbol{\mathcal{G}}) := \left[ \begin{array}{c} \frac{\lambda}{2\boldsymbol{\partial}_{\boldsymbol{\mathcal{G}}}} \ \boldsymbol{\mathcal{J}}(\boldsymbol{\mathcal{G}}) \\ \frac{\lambda}{2\boldsymbol{\partial}_{\boldsymbol{\mathcal{G}}}} \ \boldsymbol{\mathcal{J}}(\boldsymbol{\mathcal{G}}) \\ \frac{\lambda}{2\boldsymbol{\partial}_{\boldsymbol{\mathcal{G}}}} \ \boldsymbol{\mathcal{J}}(\boldsymbol{\mathcal{G}}) \\ \frac{\lambda}{2\boldsymbol{\partial}_{\boldsymbol{\mathcal{G}}}} \ \boldsymbol{\mathcal{J}}(\boldsymbol{\mathcal{G}}) \\ \vdots \\ \frac{\lambda}{2\boldsymbol{\partial}_{\boldsymbol{\mathcal{G}}}} \ \boldsymbol{\mathcal{J}}(\boldsymbol{\mathcal{G}}) \end{array} \right] \end{array} \right]$$

V J(0) = 0 for reaching global minima

$$J(\theta) = \frac{1}{2} \sum_{i=1}^{n} (h_{\theta}(x^{(\theta)}) - g^{(\theta)})$$

$$X\theta = \begin{pmatrix} - & (x^{(\theta)})^{T} \\ \theta_{\theta}(x^{(\theta)}) = \frac{1}{2} \begin{pmatrix} y^{(\theta)} \\ y^{(\theta)} \\ y^{(\theta)} \end{pmatrix} = \begin{pmatrix} h_{\theta}(x^{(\theta)}) - y^{(\theta)} \\ h_{\theta}(x^{(\theta)}) - y^{(\theta)} \\ h_{\theta}(x^{(\theta)}) - y^{(\theta)} \\ \vdots \\ h_{\theta}(x^{(\theta)}) - y^{(\theta)} \\ \vdots \\ h_{\theta}(x^{(\theta)}) - y^{(\theta)} \end{pmatrix}$$

$$J(\theta) = \frac{1}{2} \begin{pmatrix} x^{(\theta)} - y \end{pmatrix}^{T} \begin{pmatrix} x^{(\theta)} - y \\ x^{(\theta)} - y \end{pmatrix}$$

$$F_{\theta} J(\theta) = \nabla_{\theta} \left[ \frac{1}{2} \begin{pmatrix} x^{(\theta)} - y \end{pmatrix}^{T} \begin{pmatrix} x^{(\theta)} - y \\ x^{(\theta)} - y \end{pmatrix} \right]$$

$$= \frac{1}{2} \nabla_{\theta} \left[ (\theta^{T} x^{T} - y^{T}) (x^{(\theta)} - y) \right]$$



## week 1 - already applied neural net.

## Week 3 ~ kchniques

week h - devision trees

## Neurons and the brain

- just like billing rad nausons artificial nauson(a) take some input is some nontra, compute sambling than author which all as an input for next nauson(a).

-> These was need for the as classic ML models couldn't handle the bombardment of data in our one whuse named network madels perform highly well even when dataset is vory

lorge

-> at activation ; term from neuro science have are higher level features & are also called activations

-> Example: Shirt fecture sale psediction



-> Newson or newsons forming a group is called a layer

- -> All else of layor have access to all of the input but different neurons choose only specific repute and ignores others. This is not done manually but we do it algorithminal -> If we have the first column we are just left out with a logishi stayeesion model which predicts whether the Tahirit would be a top-seller or roty it also
  - does feature orgineering automatically
- -> There can be multiple hidden layors too of various sizes comprising of diff no of naurons.
- $\rightarrow$  The serif hidden layers & the serif nations in each hidden layer  $\Rightarrow$  this has large impact on efficiency of model

> deciding this comes under neural net work orchi keture.

» When here are multiple hidden layors such models are also called multilayer perception

-> mxn pixel ing is a mxn gred in which intensity values from O-255 is stored

-> On flattening it into a 10 vector we get more ell as colourne vector. The model takes this as input 6 must output the identity of passon in ing









→ The progress is used show as which a biological neuron does is way nove complicated from a simple logistic or linear regression & we do not know how brain functions completely. → One karning Algorithm hypothesis

How neural networks are implemented efficiently?

## -> Vectorization



- by vectorization. the program becomes highly efficient

### Matrix multiplication

$$\rightarrow \text{Det product}$$

$$\begin{bmatrix} a \\ b \end{bmatrix} \cdot \begin{bmatrix} c \\ d \end{bmatrix} = \begin{bmatrix} a \cdot c \\ b \cdot d \end{bmatrix}$$

$$\vec{z} \cdot \vec{z} \cdot \vec{w} \cdot \vec{a} \cdot \vec{x} \cdot \vec{w}$$

## Matrix multiplication code



l ensor flow implementation

-> Hardwritten digit classification example.



 $\rightarrow$  type of less function  $\rightarrow$  explaine later.

## Training details

→ talking →epoch:	an Iri
	→ taking →epoch:



-> to train data & minimize loss > Greadient descent

### Alternatives to sigmoid activation function



Anaxoness is not a 0 well classified throughout here we use another activation function called Rectified linear unit function, its given by  $\frac{1}{(\text{ReLU})} g(z) = \max(0, z)$ 

Search about tan h and Leaky relu most of the places we prefer leaky rely over rely...for reason refer what are the things a function must qualify to become an activation function....(1) it must be non-innear...(2) It must be diffable and deriv must not be zero else it will suffer with vanishing gradient problem, which is the case with tanh, sigmoid and step function ( which is also an activation function but not at all preferred).....Leaky relu does not have this prob, basically leaky rely is max( 0.001z, z) ....0.001 can be replaced with any small number....



l→ Commonly used activation functions → taught dates

called both ways cog  $a^{(:)} \cdot g(z)$  if g(z) = z then,

a<sup>[1]</sup> = Z which is like having no activation function.

Choosing activation functions

-> Different layers can have different activation functions.

- Activation func. of output loyer:

Depending on y characteristics of model the activation funce of output layor is decided:

-> if output layer only has 2 class of output: Binasy classification => sigmonit func.

 $\rightarrow$  output can have both the G-ve values : Ex-stock price : Regression  $\gg$  Linner activation function

-> if ... only the ... Ex-house pelice: ... > ReLU

## - Activation func. of hidden layers:

 $\rightarrow$  most common choice is ReLU, seasons:

(i) ReLU to Faster

(ii) Relu becomes flat only in one part whereas signaril flatters at both ands, when the graph flatters gradient descent becomes stower, when it becomes too signerial func. The sums stopesholy scaling in slowing down the overall pargers.



tan hactivation fore. Swish activation fore.

Leaky xelu - ·



-> Here associated with each img flore could be multiple labels.

```
→ E A self-leiving an gut an ing bask(1) whither Roursa car infundsf ne?
(i) . . . . bus . . . ?
(ii) . . . . padeshan . . . ?
```

-> These are 2 approaches to solve this:

```
(1) Tanat it as 3 separati peolor; 3 matrix which investigate one of the lines quarters & these 3 combine to form the seq, matrix
(ML matrix)
```

(ii) Alteratively train one neural network with three outputs

Advanced optimization

 $\rightarrow \omega_{j}: \omega_{j} - \alpha \frac{\partial}{\partial \omega_{i}} (\vec{\omega}, b)$ 

-> if somehow we can control alphas value to be longo initially than decrease the algo will sure even factor, this can be done by "Adam Algo sither"

Adam Algorithm

-> If algo notices the stops are becoming smaller in same dire" then it I d's value

-> If - . . Function has begun to overshort then it automatically let so that the func can seach the minima efficiently.

Adam: Adaptive moment estimation

Implementation:

$$\begin{split} & \omega_1 = \omega_1 - \frac{\alpha}{2} \frac{\partial}{\partial \omega_1} \mathcal{J}(\vec{\omega}_2, \mathbf{b}) \\ & \omega_2 = \omega_2 - \frac{\alpha}{2} \frac{\partial}{\partial \omega_1} \mathcal{J}(\vec{\omega}_2, \mathbf{b}) \\ & \vdots \\ & \omega_{10} = \omega_{10} - \frac{\alpha}{2} \frac{\partial}{\partial \omega_1} \mathcal{J}(\vec{\omega}_2, \mathbf{b}) \\ & b = \mathbf{b} - \alpha_{11} \frac{\partial}{\partial b} \mathcal{J}(\vec{\omega}_2, \mathbf{b}) \\ & \mathbf{b}^T \end{pmatrix} \\ & \mathbf{b}^T \quad \text{where here, on maxing in same diver: increase of j$$
 $. . . . . escillating ______ ; decrease of j \end{split}$ 

This is done through complex steps & is beyond the scope of this course

→ it node some initial value of cl. The multiple cl. langue & smaller & pick the one which gives the best performance state

Adam algo, is more sobust : a is funed in & works much fisher than gradient descent.

### Additional layer types

-> Danse layer: Each neuron output is a func of all the activation outputs of the previous layer.



### What is a desivative?

-> python package Sympy => used for calculus

## Computation graph

7. forward prop.

## Ly like flowshart but of math

to calculate destructives back word peop is used as if we want to calculate it from left to sight we would have to calculate everything again Gayain

### Large neural network



## -s for all whatim of t2, w2, at a w et we stort from lift & more bounds vight of forward prop.

- to find derivative we start from vight & more to left i.e. we first start from  $\frac{\partial J}{\partial \sigma^2}$  > backword .

 $\rightarrow$  if instead we use forward pump to calcular derive for a model having in nodes 6 p parameters it would only taking (n-p) skeps simifficient coloness forward prop will do it in (n-th) skeps s efficient

- this is done in knowflow & called auto diff or automatic differentiation.



### Deviding what to try next?

## -> Debugging learning algo.

- Hickory agene	J> possible solutions							
Berley Service Paperson	→ in this module we mainly leaven about diagnostic i.e what to do next when our	model is not	weeking	Ŧĥe	way	we	want	ł

### Evaluating a model

-> when we have just 2 falsess them we can just hashed graph & have an take about it's bits & variance but this is not the case with makeds with many falsess. How do you tack the graph of

🔸 We generally split data take 2 parts, francing & lat. We have model an having date she's performance on how well it products with Amazing 30% het data.



## Model selection and training / cross validation / test set

duyee ←	La(Indhim wed hard) septembridgere. Therming error is an own optimatic extingt of generalization energe cas flumeted was bained on flub-dataset. When we chose a model based on degree of wave looking. That as volumer flum That became crundy optimatic extiment generalization error too tooks flui?
-> We modify the preciduur of Italia splitt	*9
	e Stegulowischion is only droce when minimizing.
rjurt in Junk First	mulel for costinate of greater state cover a sufficience while choosing degrees of curve G use best dates to product the accuracy of multiple size the markine has never seen the height it will be a function and of greater transmission encourt
$\rightarrow D_0$ not touch	i the bat set until the machine or braining, once it has barred fully then we can let its efficiency based on her est performance.

## Degenitions bas and variance Pullipse cause of CV and view - of train " - Buy is for this example - Comple of how an algo can have both high variance & high bries. - Comple of how an algo can have both high variance & high bries. - Comple of how an algo can have both high variance & high bries. - Comple of how an algo can have both high variance & high bries. - Comple of how an algo can have both high variance & high bries. - Comple of how an algo can have both high variance & high bries.

-> High bias & volviore doean't happen to linear stegerssions but can happen who we train neural network Indicator & Jov >> Jtrain

## Regularization and bias/variance



## Establishing a baseline level of performance

 $\rightarrow$  Sporth secondition example

Lets bay training error \$ 10.2%

We may think this is high but human level performance is 10.6%

So its ok for Thain but Jou must be seduced . this madel has unione prob.

## table shifting birds ... To evaluate a model we need some baseline

the baseline can be goes too for some madels but for any madel the free thing to do is to establish baseline for ennor.

earning curve



Je as datapoints t ourset as it is hand to fit them all

> Jhain will always be lower than Jov coz the model was trained on that

THE REAL

Care : fitting a st. live to a dataset Care Battore out coz no mattin how much move dato points we collect it will have admost same any, ensur

Example no we I data points we can find a betty indegrew curve to fit them well

⇒ madel may de se well en haiving date that it may surpare humans but Jev will be halfer but as we get more date the is resolved initially

-> It might be computationally expensive doing this process again in again but it does holp get an idea of how madel is working

## Deciding what to try next sevisited

ightarrow If flatures are more than it gives the model too much flaxibility to fit very complicated models.

. Delta	1000	e1+9	Testine.
	- 1		
_	10.0		
10.0		-	

Bias/vorvierce & neural net works

ightarrow we always have to tradeopt between bias and vorwine but in neural network this is not necessarily the case.

Longe neural networks are low bros machines



sprog it has high vanuionce.

Will large neural network creati a high source? No, a larger model with well chosen regularisation will do better or attenst as well as smaller model.

## But it becomes more computationally heavy



As long as training data wint too big the numal network madel will have low bias; ofter we face high vervience probe

Here here is a writ missing continued on next page

### Error in metrics for showed datasets

-> Example.

Late say a madel has to predict a rare denser from input ; it has 1% enser G if only 0.5% of population has directer from this ensert is very very high. Bisked a single comme How we cannot behad owner % that about madel when there is a store close or "skewed class" involved.

Sunonume

Const

```
Ne use commentational and a property Precision Expected to be be there situations.
```





## Faitness, bias and ethics





## Cat classification example



## In devision trees features take on discrete values

ke	face is	sound os	not sound
	вож 4	floppy or	pointy
	whicher .	one present	ox absent

## Oval/sectangles are called nodes



Dedater Tree west note ->decision nodes -> leaf nodes



Amongst all the decision trees some will do good & some bad on taxining set on CVOR text sub

· Job of algo is out of all publi dec. trues to prix one that does well on training set 6 than also generalizes well to new data



-> this is devided based on purity

Process of building a devision tree -> What feature must be used at root node? -> what " at subsequent nodes.



Why we pat for trees w lower depths ? →So that the doesn't get too big and unweeky → Big thus are more prome to over fitting.



ie basically if we fal that splitting further down is unworthy.

44



 $\rightarrow$  Entropy function denoted by  $H(P_i)$ 

P. : fraction of examples which are associt.

we have highest entropy

$$H(p_1) = -p_1 \log_2(p_1) - p_0 \log_2(p_0) = -p_1 \log_2(p_1) - (1-p_1) \log_2(p_1)$$

Entropy function is a measure of impusity of a set of data

Stands from O goes to 1.6 comes back to 1 as funce of fraction of two scamples in comple

Guini estituia

There are many other flux. which look like extropy fluction like givi citizeia which will be present in open source package. This can also be used to make decision teres but for the sake of simplicity, in this course we willjust both at entropy function, which will work will for most of the applications.



- We take log base 2 cog version will be at height 1 which makes it cannot to predict published at
- -> On taking base a graph shifts below just that setting a threshold value becomes companitively difficult

## Choosing a split : information gain

 $\rightarrow$  The way we devide what feature to explicit on at a node is to check which reduces entropy the most.

-> reduction of entropy is called information gain

How to compute information g	ain ?			
weighted any	Info gain & punity .: which fature has higher .info gain that should be .chater	winghish ang. Prachion which is anotect:	∴ Information gain :	$H(P_{i}^{\text{(me)}}) - \left(\omega^{k_{i}k_{i}} H(P_{i}^{\text{(me)}}) + \omega^{k_{i}k_{i}k_{i}} H\right)$
Petting the taggifter process of building a decursion free				

	I man and a man
Contraction (contraction)	-> frest check at seot node for which feature has higher info gain.
00	-> Repeat process at all nodes
in the	→ ide build this using stewarions if we are doing by sealth hence its called securaries splitting

. we built an overall deusion the by patting together volutions small sub-deusion thes.

- -Some libraries have good default choices for choasing manimum depth
- -> Langur thetree will be for higher more depther of this 12 like fitting higher dogue polynomials while training lange named networks, twillowedy allow doe tox to town more complex make but will also risk available.
- -> In there, we can use can which tim set to pick man depth, is by picking many values & finally chassing one which does best on CV set In practice open-source libersite have even better ways to chasse this parameter for us.

ightarrow Amethor way is to stop when inforgain by splitting further is below a threshold value

### Using one-hot encoding of categorical feature

-s How to tackle features which can take on more than 2 values

Let say there are 9 car shape ous , pointy and around

Lets split ear shape into 3 features

> Ear shapes split into 3 features.

 $\rightarrow \pm f$  a categorical fature can take on k values, create k binary fatures (O or 1 valued)

- The value which is I is the hot feature

-> when features have things like present, about an sound notwound at are called catigorial fatures if we convert them to 1.0 then it becomes non-anial fifther that any be fed to neural notworks

## Continuous valued features

- What if a feature can take on continuous values?

How do you tackleit?

 $\rightarrow$  Let weight also be a feature in the example, then,



 $\rightarrow$  we pick deft thresholds  $\epsilon_i$  compute info gain, the one which gives may. Info gain is chosen.

- the convertion would be to sail all examples It to weight on It value of the continuous feature ( take all the value that are miliparts between the sorted his of training examples the have had over beer

# Regression tree

 $\rightarrow$  Lot say you have to predict weight of the animal.

Note: there is nothing using with a Jewission tree which chances to use some features to split along left Exight

-> Once the algo straches to a log ande it predicts the value as the ang of all training treample which strached that log made





-> If you were construction a dec tree from scratch using this data set in order to predict the weight, hey dewirm is to chanse which had have to split on...

- Suppose Reve are 3 possibilities as follows, how do you choose which a best?

- ightarrow In classification we tried to studiace contropy, have we try to studiace variance of weight i.e. y variable
- → Total weighted any of nonzerve of a nodely is subtracted from nonzionce at node >> this quantity is called seduction in variance.

value of steduction in variance K good oplit

ightarrow go on splitting this way and calculating while stopping certexia is not stoched

## Using multiple decision trees

- $\rightarrow$  Disadu of using single tree: its highly sensitive to new delta
- $\rightarrow$  To tackle this we build a bells decision trees, its called a tree ensemble.

# ightarrow 5% if we charge just 1. Animing example we get an entitedy diff two which makes the makel not so subwet.

-> Now, consider we have this Oncomble of freed. If a neuclate comes it is san in all of an G flag are made to weter. Majority soile will become the final prediction

## Sampling with seplacement

- -> This technique is used to build an ensemble of tress
- → Pick, put it back, pick again ....





→ uk tabe all easy data, patril in a long & perform sampling with veplar ment (same example can veper too) → By dering this use and up with data annian to training data Eithan use make doe bue for each of these

### Random forsest algorithm



ightarrow One disador fillin method is one end up with more or less same split at work node. To solve this: Enter Random forest.

ightarrow This make algo explose how smallcharge in data spects the model, so when new data comes prediction will be more or low accusate.

-> We generally build 64-128 fews after a certain no of tess it just becomes just computationally heavy & not much better at performance

## XG1 Boost

- -> One of major algas used to barild ensem lok of trees.
- $\rightarrow$  For successive transity picks these examples which had been misclossified earlier
- $\rightarrow$  This way we basist the algo which ends up making the model work very rively



-> three had been mis closeified and now picking them is higher probab.

- Any tree k will have higher chances of picking these examples which had been mid-closerfied by prev has is from to k-1. How it is done is complementh which can be explored later....

XCr boost stands for extreme gradient boosting

-> It assigns diff weights to diff training examples hence it need not make too many sandom ,



(Com-	COMPLEX A REAL	1 month
	방문	

### 4.1 Decision Tree

Part de landersperieren en en scherten Tex agenters het seiner Me all an eine bestellen eine de an de an einer einer pe Artigenammen bereg is het ein ein einer seinersperitiete an det umgestellen auf de

Party and a state of the second second second second

- M. Artony, J.M. Toronto, a serie of parameters in and the all conservation.
   S. Straday, Experience, and an exact the series of adds and may make the series of the series and the conservation and the test.
   S. Straday, Experience, and an experience of a set of a

### 4.3 Random Forest

tion of any be former first against per sense to birth agriculture

Of other transmission have a descension and reaction data service the segment and an extension of a second service theory. These - Descentions of transmission of the second service and service and second service of the second service - Descentions of the second second service and second service and second second second service the second service - Descentions of the second second

Surveyles and by a bound from the survey descent states of the balance of surveys states a state of the balance of the

- $\sim$  for each product the the contrast the loss of the contrast of the contrast of the theorem is defined by the contrast of t

Unsupervised Learning

# Module - 1

Welcom

-> Clustering, anomally detec	film ~> Module-1
→ Recommonder sys.	~~ <del>~</del> ~ -2
-> Rein foste ment leavining	~~ - 3

## What is clustering?

-> It sees unlabelled data points & groups them & tries to find pattern in them

→ Er Grebuping similar news, grouping DNAdato, astronomical dato analysis

## K-means intuition

- -> Randomly picks points equal to re of clusters we need (say we need 2 hora)
- -> It first assigns then it starts moving the centroid.
- -> First it assigns conter at random then it goes through all training examples & assign each point to its obserst controld.
- $\rightarrow$  Then it moves the curtivial to any location of all points assigned to that controld.
- $\rightarrow$  Now appin it checks all training examples  $\xi$  seassigns them then again  ${}^{\mathcal{J}}$
- -> when after an iteration nothing charges then it means that the algorithm has converged.

### K-Means algo sithin

- -> All µ; are vertices having same dimensions as training examples.
- $\rightarrow$  L2 norm way of working distance  $\Rightarrow$   $||z^{(i)}-\mu_k||$ ; in algo we mining squared dist min\_k|| $x^{(i)}-\mu_k||^2$
- what if a cluster has zero training examples assigned to it?
- -> M-1 seliminate that cluster more common method
  - $\underline{M-2}$  > set in this lige that cluster

### Optimization objective

- c<sup>(1)</sup>, index of cluster (1,2,...,K) to which which example a<sup>(1)</sup> is currently assigned
- py . cluster can trovid k
- et (1), during contract of churley to which example x (1) has been assigned

Cast furction : Distortion function

$$\overline{\mathcal{T}}\left(\mathcal{L}_{2}^{(0)},\dots,\mathcal{L}_{n}^{(n)},\mu_{1},\dots,\mu_{k}\right) = \frac{1}{m} \sum_{\substack{\{x\} \\ x \neq y}}^{m} \left\| x^{(0)} - \mu_{2} \omega \right\|^{2}$$

sq. diet blu all points a this Juster contraid

. to minimize the machine auto matrically starts to assign points to clasest controlds

And to more cluster contraits .....

it assigns contrad to mean value of coexclination of examples assigned to that contraid.

## Initialising K-mean

-> 20 of clubber(K) must be loss than training examples(m) i.e. K.Km ; etc. thue would be less than I training example par clubbe

## -> Randomly pick K training Geomptes & awign to my the to them.



-> so we sondonly initialize & build multiple models & than compute cast func for than the one with lowest one will be chosen.

optima



## Choosing the number of clusters

### -> Elbow method

-> Jus K a plotted K a plotted ..... the se of clusters till which I deriverses sapidly of after which it shows down is chosen as our preglamed value of K ; that porticular point is

-> K is after ambiguous cay many-a-times cay. Jus K waves are many-a-times smooth & allow cannot be spotted, chassing <u>pe</u> of clusters just to studiuse J is not god regiter when below indiancy, only with longest rules of K.

TET.

> twinde off blue how good img boks us how much you can compress

→ bude of bbu - of t-shurt size G coat of manufacturing.

Here most of the times its us who doub how many clusters are stop wired according to our need.

Anomally defection

Finding unusual events

Anomally detection example : Airoraft engine features

x, + heat gene stated

X2 = vibration intensity

Dataset: { x(1) x(2), ..., x(m) } New angine: x test

Density estimation : fixet thing we do with dataset is build a madel which calculate peobability of points being seen in dataset



Usage :

faul decktation opsil Balances, how often log in 2 no of unknows wished ? transactions ? typing speed? If flagged security han looks of it checkly may giv capital or triggen alones at ....



Graussian (Normal) Distribution

or bell-shaped dis tribution

X -> xandom variable

peobab. Je karmined by 12.5 m²

Say we get a 100 m & plot histogram it will look like this

 $p(x) \cdot \frac{1}{\sigma - \sqrt{2\pi}} e^{-\frac{(x+1)^2}{2\sigma^2}}$   $P(x) \cdot \frac{1}{\sigma - \sqrt{2\pi}} e^{-\frac{(x+1)^2}{2\sigma^2}} e^{-\frac{1}{\sigma - \sqrt{2\pi}}} e^{-\frac{(x+1)^2}{2\sigma^2}}$   $P(x) \cdot \frac{1}{\sigma - \sqrt{2\pi}} e^{-\frac{(x+1)^2}{2\sigma^2}} e^{-\frac{(x+1)^2}{2\sigma^2}}} e^{-\frac{(x+1)^2}{2\sigma^2}} e^{-\frac{(x+1$ 

Anomally detection algorithm





theograph's height represents p(r1)·p(r2)

### Developing & evaluating an anomally detection system

In practical situation its ok even if there are some anamabus arample in training set its a



In general its better if CVG bet set have an amplies in them.

When we do collect date of for we know that a small number of examples are anamotives horse we use them with labels to train, check & evaluat. For example,



. having some anomalists helps the market, why not sun a supervised algo. on this data.

### Anomally detection vs supervised learning

-> : future anomalies maybe entirely diff from priorenes its safer to use anomally detec (if our date is of this type) cos anything deviating away too much will be flagged.



### Choosing what features to use

-> In unsupervised learning choosing what features to use a what not to becomes very impostant.

-> Tsy to make features gauge ian.



smaller values of c is more effective

We transform the Reatures to make it gaussian

Directly log & might give exces of 2 was zero .: we write log (x+0001) to avoid the conoxs.

-> Same transformation must be applied to CV & test sets too.

### Exxos analysis



- then we have to look at that example & try to figure out why was it an anomally & try to come up with a new feature which can recognize these kinds of anomallies



Module-2

Kecommender systems

Making recommendations

surf o.t.	COST LANS DO	TWO OTHER AND IN CASE OF	A CONTRACTOR
. We have useds as well as some no of items	**	And her tree	And Address of
∩ <sub>u</sub> ⊨ of users	Section .	ALC: N	a and a
non · · · items (have movies)	States a sum	9.92	
$s(i,j) \in 1$ if user j has soled the movie	100		
8(1,1) +1 8(item 10, user 10) 4(1	em ne , Usex 🗠	e)	
8(3,1)+0			108. 1
g(i,j)=dating giv	en by User	j to movie i	(defined only if

-> Look at movies user hourit rated yet, predict their rating for it, if that high then steenmend it to them



-> None sefus becaue that matter water have sated the same many collaboratively, giving us some of what this more may be like that allows taging what are appropriate farbase for that more, This in turn allows us to produce how other was that haven't get rated that same more may decide to sate it in fixer.

## Binary labels : forus, likes and clicks



### Mean normalisation

-> Algo becomes more efficient if we perform mean normalization



Let's say a new user Eve comes to hus not valed any manife yet, now if we sure formable it will predict that all ratings by Eve is zery which is not seasonable



- Due to mean normalisation algo becomes a little faster too.

-> We can also normalize colourns tricked of results can be done with new movie awaries but generably its not done in his way cay we don't know how much the wear may like it.





## -> Its done in 2 skps setsieval & sonking



## Isincipal Component Analysis

## Reducing the number of features

-> PCA hulps in visualisation of data, say we have 50 as 100 as 1000 falway, with PCA we seduce this to 2 as 3 falway so that we can plat it.



## PCA algorithm

-> Feature should be normalized first to have zero mean, if readed give this Rature scaling must also be applied.

->PCA sceness all and G storts to make its own area, on a good and the data must be spiled widely i.e. it must capture allow manimum variation of data.



## Module - 3

## What is seen forcement barning?

- $\rightarrow$  Reward function : when model does something distribute its rewarded 6 when it does something undertrable its penalised.
- -> Applications of seconfoscement lowering: controlling volvots, factory optimization, financial (stock) trading, playing games (including video games)

## Mors sover example



1 66 ove terminal state

## The setusion in seinforcement laruning

R. > seward of ; the state

( » discount factor (generally a number close to 1)

## $\mathsf{R}_{e}\mathsf{furn} = \mathsf{R}_1 + \mathsf{f} \mathsf{R}_2 + \mathsf{f}^2 \mathsf{R}_3 + \dots$

and a set of the set o	

-> Discount factor also pushes away the -ve rewords is for there is pushes away lesser the -ve value will be added.

### Making decisions: policies in sein forcement learning

- Our job in skinforcement banning is to come up with a policy (T) which takes state (2) as input & map it to an action (a) which is to be performed.

## Review of key concepts

-> Formation of reinforcement learning application is called Markov Decision process (MDP), it selfers that the future only depends on current shall & not on anything that might have accurred prior to getting to the current state.





State -action	value	function	eramp	



R(E) is also called as the immediate seward.

## Random (sto chastic) environment

-> Due to some standomness in envisionment the system may not perform as commanded Ex on commanding the stoven to perform a particular action, it does so 90% of threes but 10% of times -> How in stochastic we have multiple steturns, maximizing selection describ make some cay steturn is a standom number instead we maximize the ang-value of own of discounted stewards



### Example of continuous state space applications

-> Up until now the source could be in one of six partiens, these are discrete state, say now that relate (funch) can be any where on a long this is continuous state space.



Shut of teach is given by it position in my dram, it as ion -takion, it's well in may dram Gits angular well ô ie how fast it can team. State of helicopter includes its position in 2, y, z dro", its 2011, pitch, your & their of one sate of changes



∴ in a continuous state MDP a state is a vector comprising of many numbers which could take on any range of values.

## Lunar lander



Leasning the state -value function





If he neural network is haved peoply then for any shell is All four Q values can be calculated & he max one among them can be performed. Guel: to learn a policy IT flust, given & picks at 17(5) 80 as to measuringe seturn (\* 0.1985 for this. [this is given, generally for these high when it is chosen]

> Appsionch: We will use bellmank en to create a training set with lots of examples of a Gy, Guse supervised learning then to born a map from a toy re from state -action pair to target value of Qlég

$$Q(g,a): R(a) + 1 \max_{a} Q(g',a) = f_{ab}(a) \approx y$$




Refer last Moon rover lab for more info on reinforcement learning, several papers have been mentioned in labs to study from, refer those at times....do download <u>all lab files</u>

Philometican funct. Looky where max  $(d \geq 2)$ ELU  $f(u) \cdot \int x \quad if \quad 2>0$   $(d(e^{\tau}-1) \quad if \quad \tau \leq 0$ SELU  $g(u) = \lambda \quad f(u)$  will self remove by recent O stel down  $(f(u) \cdot \nabla E(u) \cdot 1)$  bet reall should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet reall should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  bet maded should be  $f(u) \cdot \nabla E(u) \cdot 1)$  be the should be a should be  $f(u) \cdot \nabla E(u) \cdot 1)$  be the should be a should be should be a should be a should be

The second secon





27	-	-	Conductional ADM	Statement ( ) Street,	sfor 1-D data
	1000		Convolution of the	The second s	- ( )
	1000		- Roundian NN	A 10 10 10 10 10 10 10 10 10 10 10 10 10	10 C
-			Theathing		-
-		_		successive successive in Amount	

Structured data: database with colourns, news etc. Unstructured data: - . ing, audro or text to translate

# Why is deep leasning taking off?



# Module -2





# Logistic Regression

$$f(\omega,b,x) = \frac{1}{1 + e^{-(\omega x + b)}}$$

Logistic seguession cost function

$$\mathcal{L}(\hat{\mathcal{G}}, \mathcal{G}) = -(\mathcal{G} \log \hat{\mathcal{G}} + (1-\mathcal{G}) \log (1-\mathcal{G}))$$
  
 $\mathcal{J}(\mathcal{W}, \mathbf{b}) = \frac{1}{\mathcal{D}} \sum_{i=1}^{\infty} \mathcal{L}(\hat{\mathcal{G}}^{(i)}, \mathcal{G}^{(i)})$ 

100 - + + + + + + + + + + + + + + + + + +	All read and set.	10-10-11
11台上	1 K	1111
:420	- ALTER AND	
0.000	and the second	a martin

-> Python broadcasting

$$\frac{1}{1+e^{-a}},$$

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

Vertraiging Logistic its section Vertraiging Logistic its section Vertraiged Vertraiged Vertraiged Vertraiged Vertraiged Vertraiged

# Beadwesting in Fythe

# Note on numpy/pythan vectors

a · np. sonder · sonder (5) 
$$\Rightarrow$$
 5 num of gaussine  $\left[-, -, -, -, -\right] \longrightarrow$  here as a  $T$  of a ·  $\sigma^{-1}$  is just and   
a ·  $\cdots$  (5,1)  $\Rightarrow$   $\left[\left[-, -, -, -, -, -, -\right]\right] \Rightarrow$  here  $a \neq a^{T}$  here a  $\sigma^{-1}$  is an average

Tip: we assert func. I python to get notified who somethin is indeed wrong.

# Newral network



Vectorizing across multiple examples

# Activation functions

 $\tan h = \frac{e^{z} - e^{-z}}{e^{z} + e^{-z}}$  took is almost always superior to sigmaid function cay its contexed at 0 G mean comes near to 0 which make barning fails

tank cannol be used in output bayer of binary class cas we need output tober 0 to I there

Problem with sigmoid & tank & when Z is very small or harge than gradient is clar to zone which shows down the harming a left

(2 0.0) orsome small

$$\frac{\operatorname{Re}_{u}}{\operatorname{Re}_{u}} a = \max(0, 2) \quad \text{Looky Rolu 5} \quad \dots \quad \\ \operatorname{Re}_{u} \operatorname{Hord}_{u} \operatorname{Rolu}_{u} \operatorname{Ro$$

Why do we need non-linear actuation functions

->If we activit all neurons with linear function than the model connort learn anything nove complex than a simple linear sequencien, have we peoply non-linear functions as activation func. -> Linear func, is used very samply when we do compression (in hidden layors) nove about this later -> If the model we build up sequession then we can use it only in the subput layor to get the derived cutput.

Grundient descent for rewal networks -> keepdims = True = provers pythen from giving output of shape (11, )



Kandom instealisation

-In logistic/linear stepstession we can initialize all weights to zero but we cannot do the same in drep coz greationt descent won't work. Why?

$$\begin{array}{c} \rightarrow & \text{Reason.} \\ \text{Lets say there are 1 hidden layer, with two nervous.} \\ & \pi \rightarrow \stackrel{O}{\longrightarrow} \stackrel{O}{\longrightarrow} \stackrel{O}{\longrightarrow} \text{output} \\ & \omega^{1} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{1} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \omega^{2} = \begin{bmatrix} 0 & 0 \end{bmatrix} \quad b \cdot \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Now since you't & you's same their updations are al ways grand became hence they want be able to calculat anything complex, it just same This is known as Symmetry peddom

generally vous small can if we are using symria ax torth taking large value makes greatent very small have makin learnin very small if output ax hidden layers don't involve tank ar signarid then this is not much of an visue.

# Module -4

Deep L-Layer neural network





### Why deep sepsesentations



The faither layers automatically compute & analyse complex stuff.

XOR =	nes	usig	logn	neus	ons	⇒ fast &	easy		
			single h	ayer	⇒ ¢	slow & co	mplex	Ę	in efficien

## For word & back word prof

-> in forward prop itself we cache is store a few values so that they can be used on backprop, we cache 2, wh





# Paxametexs vs Hypexpaxametexs

w, b are parbometers

Hypexparameters are devided by us, which ultimately control the parameters.

E. a, # itera tion, # hidden layers, # units, choice of activation func., segularisation, mini-batch size, momentum et.

# Louise - 2: Hyper passameter tuning

# Train/Dev/Test sets

- -> Pervicuely the split would be F920 ox 80/20 or 60/20/20, but now where there is abundance of data we can have fruer examples in day/test like maybe 10,000 if there are say million trainin examples. Many times split is like 98/1/1 or 99.5/025/025 et.
- -> DevEtest set must come for some dataset not different ones for en all are pics clicked from smort phone.
- -> Normally when ppl say they just have train /test set then it means they achaily have train/devset but they call it train /test.

# Bias/Vavience



## Basic principle for mochine learning



In NN we don't have to warry able beakaff as getting more data I varience arbeat affectin bias much & choosing biggen NN with paper regularisation I bias almostarily at a choosing biggen NN with paper regularisation of fative varience.



$$T(\omega^{(i)}, b^{(i)}, \dots, \omega^{(i)}, b^{(i)}) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\left(\hat{g}^{(i)}, \hat{g}^{(i)}\right) + \frac{\lambda}{2n} \sum_{i=1}^{L} \left| i_{i} \mathcal{L} 0 \right|_{F}^{2} \qquad \left| \left| \omega^{(1)} \right|^{2} + \sum_{i=1}^{n} \sum_{j=1}^{n} \left( \omega^{(1)}_{i,j} \right)^{2} + \sum_{j=1}^{n} \sum_{j=1}^{n}$$

-> As we want by them gradient value? In sigmaid on tan h a almost linear a reduces complexity hence reduces our fifting

# Depout segularisation

We copy NN then wondomly xemove few neurons & train it to test if overfibting is reduced.



> matrix with sondon values \_\_\_\_\_\_ probability that a numor

Ty hat say the as 50 with \$ 20% charge of degoing remains so 10 remains are dropped but weight of later layers are frained based on all remains two sole as to not defined s has any the as so that from the course of any program to Two time because Ins. I comparative because Ins. complex due to this -sat each pass of gradent descent we grow out different est of functions not he came are -sat here time we due to be discovered in the particulation of the came are weights : a3/= keep-pool

-> The inputs get climinated sandomly hence the nearon connot sely on a partialay feature ... mathes I is regularized

-> keep peop is different for different layers, he layer having more neuro	a house lower heepprob Eurice verses 4 in layers where we do not warry about oursfitting hop-prob as i	L. Input layer is generall not penali -sed or if read its done with high
. The doconcide of this is we have even more hyper parameters now.	- One down side of disport-is cut func is not well defined anymae, so are of the publ sol" is to have of disport is answer that barningsame is normally than implement deprove	keep_pado >0.9
$\rightarrow \mathbb{T}$ is generally used in computer vision $\rightarrow$ The computation of cos	there becomes harder as input are drapped at random.	
Othes xegularisation methods		

# Data augmentation [ vefer prev course pg no 24]



initially geoph will underfit G on long training its gone overfit hence its stopped somewhere at middle Here we follow tradeoff between bias & vouvence

 $\chi = \frac{\chi - \mu}{\sigma^{-2}}$  be sure to normalize test of devisets  $\int_{\sigma}^{\infty} \frac{\psi_{\sigma}}{\psi_{\sigma}} dv$  with  $\mu_{\text{train}} \in \sigma^{-2}_{\text{train}}$ .



of Un-normalized inputs take a lot of iterations to conversige & bounce around a lot too



Las if Lis longe this is very huge Ewill explode

Numerical approximation of gradients

Gradient checking

Take 
$$W_{1}^{(2)}$$
  $U_{1}^{(2)}$   $W_{2}^{(2)}$   $U_{2}^{(2)}$   $U_{2}^{(2$ 

Take 
$$dW^{(1)}_{j}$$
,  $db^{(1)}_{j}$ ,....,  $dW^{(1)}_{j}$ ,  $db^{(1)}_{j} \notin$  xehape into a big vector  $d\theta$ 

Gread check  
for each r:  

$$d \Theta_{approx}[i] = \int (\Theta_{i}, \Theta_{2}, ..., \Theta, + \varepsilon, ...) - J(\Theta_{i}, \Theta_{2}, ..., \Theta, -\varepsilon, ...)$$
  
 $2\varepsilon$   
check whethes  $d\Theta^{(c)} \simeq d\Theta^{(c)}_{approx}$   
to check it...  
 $\frac{\left| | d\Theta_{approx} - d\Theta | \right|_{2}}{\left| | d\Theta_{approx} - d\Theta | \right|_{2}} \approx 10^{-7} \rightarrow g_{stat}$   
 $\frac{\left| | d\Theta_{approx} - d\Theta | \right|_{2}}{\left| | d\Theta_{approx} - d\Theta | \right|_{2}} \approx 10^{-7} \rightarrow \phi k \rightarrow hat$   
 $\left| | x | \right|_{2}^{2} \sim \left| | |x | |_{2}^{2}$   
 $determinant$   
 $0^{-3} \rightarrow de bug$ 

I mplementation of grad check

→ cig its very slow to cak grad check. → link for which layer 6, when d0 is differing 6, by to debug > ic d0 in cludes segularisation term too → cag comparing cast becomes difficult when resonderly subset fromums are despect so check and than implement dispect

lexy society it so happen!

that implementation of gradient descent is consert when wells are clase to zero, but as we sure gradient descent wells brane bigging it desire becames inaccurat after that. To check this is not the case sure aread check at wells in the state of the test of the state.





# iterations

Batch

Too much time

per iteration

Chossing mini batch size  $\rightarrow$  Size of mini-baldh = m : Bald reddent descent :  $(X^{\xi_1}, Y^{\xi_1}) = (X, Y)$  $(\chi^{\sharp \iota 3}, \chi^{\sharp \iota 3}) \cdot (\chi^{\iota 0}, g^{\iota 0})$ 

In practise mini-batch size is somewhere in between 16m.

Stochastic Looses benifit fuerbaisation 08 we are iterating & computing & updation at each example

Choosing mini bath size

Take account of available CPU & 61PU If data set is too small then (m < 2000) use batch. Typicals izes are in process of 2 => 64, 128, 256, 512

Mini - batch fast V Fastest learning in general V

Of course we can create a hyper-posameter & pick out the best mini-butch size too

Exponentially weighted moving are suge

 $\beta \cdot \beta \cdot \theta_{n-1} + (1-\beta) \Theta_n$  if B is small then more weight is given to current value hence the graph will be noisy (yellow)

Using this we are going build hyper parameter which chooses but B.

Undexstanding exponentially weighted averages

$$\begin{split} & \stackrel{\Phi_{1}}{\rightarrow} \Theta_{n} \quad \forall_{1} \in \Phi_{n}^{*} + \langle l - \beta \rangle \Phi_{1} \\ & \stackrel{\Phi_{2}}{\rightarrow} \beta \left[ \left( l - \beta \right) \Theta_{1} \beta \frac{\beta}{\beta} + \left( l - \beta \right) \Theta_{2} = \left( l - \beta \right) \left[ \Theta_{2} + \beta \Theta_{1} + \beta^{2} \Theta_{n} \right] \\ & \stackrel{\Phi_{3}}{\rightarrow} = \beta \cdot \Psi_{2} + \left( l - \beta \right) \Theta_{3} \cdot \left( l - \beta \right) \left[ \Theta_{3} + \beta \Theta_{2} + \beta^{2} \Theta_{n} + \beta^{2} \Theta_{n} + \beta^{2} \Theta_{n} \right] \\ & \stackrel{\downarrow}{\vdots} \\ & \stackrel{\Phi_{n}}{\rightarrow} = \left( l - \beta \right) \left[ \Theta_{n} + \beta \Theta_{n-1} + \beta^{2} \Theta_{n-2} + \beta^{2} \Theta_{n-3} + \dots + \beta^{2} \Theta_{n} + \beta \Theta_{0} \right] \quad - \qquad (1)$$

. As power I weightage to older data points is reduced

of secent datapoints taken into consideration is given by <u>1</u>

As value of  $\beta$  increases graph changes more gradually hence less noisy graph.

uning a more of more of more of more of the more of th

orth moment greatient

Bias-consection



Look at graphs at left, three -90 =0 & have at the start the average is not calculated properly, have we need some conserving so instead of 19t use equals RHS in (1) to  $\underline{-9t}$   $1 - B^{t}$  $\rightarrow if t is less 1 - B^{t}$  is small here  $\underline{-9t}$  will be large

Grenewally B= 0.9 fox NN

-> Many a times people ignose initial phase & just check hates phase, if you wanna consect we can implement this.

> Eventually both the graphs combine (as shown in graphs on sight) doing the juit provide bias coveries

 $g_{ko} = \beta_{1} \cdot g_{kopenv} + (1-\beta) dw$ meet commen value for  $\beta_{1}$  to 0.9in general we don't apply bias correction reg after just <u>(0 recessions</u> our also usuallike usasmed up updation  $U = W - \alpha \cdot g_{av}$ Initial  $rg_{av} = 0 \in g_{ab} = 0$   $b \cdot b - \alpha \cdot g_{b}$ many at me we just use  $rg_{av} \cdot \beta_{1} \cdot g_{1} + dw$ many at me wight use  $rg_{av} \cdot \beta_{1} \cdot g_{1} + dw$ many at me wight use  $rg_{av} \cdot \beta_{1} \cdot g_{1} + dw$   $rg_{av} \cdot g_{1} \cdot g_{1} \cdot g_{1} + dw$   $rg_{av} \cdot g_{1} \cdot g_{1}$  RMS prop

in vertical dis" we wanna slow down but in crease learning in horizontal dis"

$$\begin{split} & \mathcal{B}_{d,w} = \beta_{2} \, \mathcal{S}_{d-pperv} + (1-\beta_{1}) \, \mathcal{J}(v)^{2} & \stackrel{b}{\longrightarrow} \mathcal{J}(w) \quad \text{is small} \quad \text{ince warms} \quad f \text{ th } G \quad \text{dh } \text{is large } G \quad \text{we warms} \quad f \text{ th } G \quad \text{dh } \text{is large } G \quad \text{we warms} \quad f \text{ th } G \quad f \text{ th } G \quad \text{we warms} \quad f \text{ th } G \quad \text{we warms} \quad f \text{ th } G \quad f \text{$$

initialization: vau = 0, vab = 0, saw = 0, sab = 0

most common value for \$2= 0.999

Hdam optimization algorithm : adaptive moment estimation algorithm

-> Its a combination of momentum & RMS prop

 $\rightarrow Bns connection is performed on momentum & RMS prep. i.e. <math>n^{g} = \frac{n^{g}}{1-B^{t}} = \frac{8}{1-B^{t}}$ 

Learning sate decay

- ) If we slowly soluce the learning sate this helps the algo., this is called learning sate decay

-> Inhuition, initially we can take bigg ex stops but as we slowt to go towards minima lower knowing rati is better as we can escillate in a tighter region



-y as and when sequered

 $\rightarrow$  normally default values of  $\beta$ ,  $\beta$ ,  $\epsilon$  is used  $\xi$  K is fund

#### The peoblem of local optimo

-sales a local optima, this is mose common kird & is called a sould be point.

the gradient & zero have too

#### > all dobs have one local optima

-> There are very intuitions about lower dimensional space but as wego to higher dimensions the intuition decreases. More often in high dimensional space we stop at saddle point rath -ex than global goty

-> Problem of plateaus



Plakan is a segion where destructive is clase to zero for a long time, then it stays on plataw for a very long time and there descends off of it

How to choose the values of hyper-paramaters?

-> We can choose some no of standom points in a n-dimensional space where n is the no of hyperparameters

-> Coorse to fire search

In the space say a particular point, did well then we goom into that segion and take more samples from there and go an experience this

### Using an appropriate scale to pick hyper parameters

-> Sampling unifermity doesn't make much sense, we have to sample acress an appropriate scale.

Let's say we have to choose met between 50 to 100. her uniformly searching is fire.

To search this way on log scale we implement it this way. ( we show 0.000 1 mm. 0.001 mm. 0.1 mm. 0.1 mm. 0.1 mm. hats say south between a & b Some associated standom values from these stanges are taken then start = log 10 a & end : log 10

×E [stort; end] > hyp.parameter - 10 & check to take best volue

En3 Let's say we have to choose B

B between 0.9 to 0.999

B= 0.9000 to 0.9005 Fr both soughly toke any own 10 orangles As B gets closer to 1 it's sensitivity increases very highly say B= 0.999 to 0.9995 .. We need to sample very 2 3 does in 1 avgever 1000 example 1 ang over 2000 examples

# Hyper pastameters tuning in practise : pandas vs Caurian

-> Some times hyper ponentines can get stale, at these times they need take selected as re-malusted ofter centrain period of time

-> Ho w people go about searching for hyper pasameters

Babysit one model : Panda strakegy where when detend is very large or there are inited CPU or GPU resources, its moviewed individually & hyperpresenters are updated

Traning many models at a time : Couriar shakeyy -> Many moduly with different sets of hypersparsameters are chosen a trained possallyly of the best one of pricked

The type of approach we choose depends on computational seconsures & the amount of data we have

TVENTUR New activations in a returner -> Bath normalization (created by two separates: Sergey Loffe & Christian Szegedy), make by perparameter search peoblem vory easire, make WN more related to the to chose from longe surge Quril also help to train deep network more easily.

<sup>-&</sup>gt; We generally normalize the input so that computation becomes easier, can we similarly normalize outputs of her each layor?

→ Yes, now have is a choice to whether making the finally output all as z<sup>CU</sup> which gave into activation fore, i.e. g<sup>(D)</sup>(z<sup>(D)</sup>) = a<sup>(D)</sup>, there are many anguneris on the bot we goverbally normalize z<sup>(D)</sup>
<u>Implementation</u>

Late say we have z of a extrat by  $z^{(i)}$  .....  $z^{(i)}$  of some logar

$$\mu = \frac{1}{m} \sum_{l=1}^{n} z^{(l)}$$

$$\pi^{-2} \cdot \frac{1}{m} \sum_{l=1}^{n} (z^{(l)} - \mu)^{2}$$

$$\overline{z}_{mm}^{(l)} \cdot \overline{z}^{(l)} - \mu$$

$$\sqrt{y^{-2} + \varepsilon} \quad \varepsilon^{z}$$

When we do not wont the mean as zero & vanience as 1 them .....

 $\widetilde{Z}^{(i)}$  . ( $\widetilde{z}^{(i)}_{new}$  +  $\beta$  . (i) are learnable pasarrobes which are updated using some algo.

if i to a point & point then Z = Z (i)

for other value of TGB the mean & variance of normalised data charges.

En in signarial various longer than I might be perfixed or non-zono nam to take advantage of its non-linear rateur Note: B have & B in exponential aug. is different.

Fifting batch normalization into a neural network

$$\times \xrightarrow{\mathsf{M}^{(0)},\mathsf{b}^{(1)}} Z^{(1)} \xrightarrow{\mathfrak{g}^{(1)},\mathsf{f}^{(1)}} \stackrel{\mathfrak{g}^{(2)}}{\Longrightarrow} \xrightarrow{\mathfrak{a}^{(2)}} g^{(2)} \xrightarrow{\mathfrak{g}^{(2)},\mathsf{g}^{($$



Pasameters: W<sup>CO</sup>, b<sup>CO</sup>, W<sup>CO</sup>, b<sup>LO</sup>, ...., W<sup>LO</sup>, b<sup>CO</sup> β<sup>CO</sup>, τ<sup>CO</sup>, β<sup>DO</sup>, τ<sup>CO</sup>, β<sup>DO</sup>, τ<sup>CO</sup>, β<sup>CO</sup>, τ<sup>CO</sup>, β<sup>CO</sup>, τ<sup>CO</sup>, β<sup>CO</sup>, τ<sup>CO</sup>, β<sup>CO</sup>, τ<sup>CO</sup>, β<sup>CO</sup>, τ<sup>CO</sup>, β<sup>CO</sup>, τ<sup>CO</sup>, δ<sup>CO</sup>, δ<sup>CO</sup>

For mini batch gradient descent

$$\begin{array}{c} X \xrightarrow{\mathbb{R}^3} & \underbrace{\mathcal{M}^{(0)}}_{\mathbb{R}^3} & \mathbb{Z} \xrightarrow{\mathbb{C}^3} & \underbrace{\mathbb{R}^3}_{\mathbb{R}^3} & \mathbb{Z} \xrightarrow{\mathbb{C}^3} & \mathbb{Z} \xrightarrow{\mathbb{C}^3}_{\mathbb{R}^3} &$$

Passometers . W<sup>CLD</sup>, b<sup>CLD</sup>, B<sup>CLD</sup>, Y<sup>CLD</sup>

$$z^{(1)} \cdot \bigcup^{(2)}_{\alpha} (k) + b^{(1)}_{\beta}$$
Now we such  $z^{(2)}_{\alpha} = z^{(2)} + \mu_{\sqrt{y^{2} + 2}}$  ( $\mu$  is subtracted if we all b while calculating  $z$  or not does not matrix at all
$$\frac{1}{\sqrt{y^{2} + 2}}$$
We can compute this as
$$z + \bigcup^{(2)}_{\alpha} (k - 1)$$

$$\frac{1}{\sqrt{y^{-2} + 2}}$$
we can just versione  $b^{(1)}_{\alpha}$ 

$$\frac{1}{\sqrt{y^{-2} + 2}}$$
we can just versione  $b^{(1)}_{\alpha}$ 

$$\frac{1}{\sqrt{y^{-2} + 2}}$$

$$\frac{1}{\sqrt{y^{-2} + 2}}$$
we can just versione  $b^{(1)}_{\alpha}$ 

$$\frac{1}{\sqrt{y^{-2} + 2}}$$

$$\frac{1}{\sqrt{y^{-2} + 2$$

#### Why does Batch norm worsk?

#### Input distribution

-> If input distribution varies from kts say shall 1 to state 2. A medid trained for inputs of pattern state 1 than it will made particularly not do so well on date from state 2 even if grownd true function is same, if grown I have function also changes than its even worse. This moving from state is state is called <u>covariat</u> clift.





Why batch norm works?

- -> The data from pases layou vary is covariate shift occurs as model loone due to 812 atlast moon & vanience will be in a find starge which helps the model to loon more quickly & effici -criently. It limits how much the data change also to prev layous. It is stability if prevides laker layout hostand ground to shond on.
- -> The dependency of one layer on prive layer & & layer are more fur. It wraters the coupling between what early layer pass metres has to do with later layer parameters

# BN also has a segularisation effect

House Service as repub	attivites;
and worth A division	
Charles in the local division of the	And Party of Street, or other
Research of the local day of the local operation of the local day of the local day of the local days o	and a set of a set of a set of a
Real Party States	as clines-

→ just like disposit which adds rosse into keyon by aliminating neuron at sandom this ensures hater layou do not depend on pseurious layous much. In assurilar, way Bri adds additive noise since pick or 2 of minibatch is calculated so there is a digit noise introduced hence this beings about a digit segularization effect.

-> As size of mini bookh 1 the noise inhodulud I cay more examples are considered at a time hence the segularisation effect is J

No smallization : collapse numbers between 061 Standardization: make mean 06 ubusience 1.

.

#### Batch norm at test time

- how to efficiently compute RG+ "ever entire minibatch. It's done using exponentially weighted average

 $\xi$  at test time we compute  $Z_{nown}$  using exponentially weighted angolf  $\mu\,\xi_0-^2$ 

in frame works they have more sobust ways to calculate plant required for that time

#### Soft max Registersion.

C + Hasses output larger has C remains  $\rightarrow \sum_{i=1}^{N} \sum_{j=1}^{N} \sum_{i=1}^{N} \sum_{i=1}^{N$ 

Islaining a soft max classifier

What is hardmax func. ? The largest number is assigned as one and all else are assigned as zero.

## Deep learning frame works



## Tensorflow

-> If we implement four peop in sight way, it calculates back peop automatically.

> just like a classic if removables what was done in formula peop then it plays it back to calculate derivative



it makes computation graph here to find all draivatives.

Structuring ML Projects

# Module -1

Suppose our medel is done but we wanna further improve it. How do you choose what to do?

## Os they onalization

-> Its always better to have hyper parameters which perform different things when twenhed. It better if hyper parameters and interpretent & twenhing one shouldn't affect others.

→ When madel doent fit train set > bigger network different algo, diff optimizer ek · · \_ dev \_ > Regularisation, bigger thaining set · · · tral . > Bigger dev set · · Perform well in rient world > change dev set as cast func.

-> Early stopping is discouraged car it tures to do 2 things at some tree is handle fitting of train and dov set.

#### Single number evaluation metric

-> Suppose you have percision & second new combine both to a single metric to evaluat the model. Evaluation is easing if there is only one number / variable to check

#### Satisficing and optimizing meterics

-> Its not always easy to combine metaics into a single number, then we do this.

-> Say there are N meterics, then we choose one of them as eptimizing and all others as saturficing.



ne say we pick model which has sumtime is each of these will have a these hold only alwave which it will be accepted below 100 ms with Lighut accuracy than,

dev/ test set

accuracy is optimizing metric & suntime is satisficing metric

#### Train/dev/test distributions

- -> Dev & fust sets must come from same superset of data
- -> The data given must be somilar to what we expect in future

#### Size of dry and test sets

-> 70-30 as 60-20-20 splits are at for small size date but date is very large days anythin as smaller of her and another like 1% 61% for days sites sets are also five

#### When to change dev / lat sets and meterces?

-> Let's say three are ture algos A 4 B ; A has been survey a performing better on sceners but due to remain cancer (interthind) Bis encidered better , in there schulteres that is a need to charge metric system or charge

#### why human level performance?

-> Bayu aptimal cance: the bat preside cover you can have on a function mapped from x to y. -> It is normally chee to human hit paformance E cat closeficition. Harvas are partily good at closefiging cats That why after human level supervising madels little hard.

 Provideble biss
 E1

 Human 25 Bayes
 1%

 Training enner
 8%

 Dev evener
 10%

 Station of the set of th

Understanding human level performance

Bayes over + Avridable bias Training + y + difference Dev + Vanderce

Suspassing human level performance



Is these are all tooks where marking has

access to huge databases

Humans are more better at Nahusal perception took such as speech recognition, compater usion et though some madels, perform better than humans have too ... its handor for a markine to perform too well on a natural processing task.

Bayes essor or human level performance

#### Improving your model performance



# Module - 2

#### Carrying out error analysis

-> whenever we try to implement some perceible as to improve accuracy of model soughly per-calculate by here and it could improve the model and then decide whether it's worth paraming

#### Cleaning up in cossectly labelled date

-> Depending on how much excer is being caused by increasedly labelled data we can service as fire it.

- Developtist set must be from some distribution, train out can be from a different one.

#### Build your freet sys quickly and then itersole

-> Building a ok system then weeking on it and improving is better than building the perfect one at first ge.

#### Training and testing on different distribution

-> What data we may ensure in seal would emile to these must be at maximum in dev/test set

#### Bras & vowence with mismatched data distributions

-+ Lett say hain examples is 1% adre examples is 10% but thain data à dev data is different, example thain data has high searing but dru set has bury ring too -> Now use do not know if wantere 3 high as algo works well its just that dev set is hand.

-> train-dev export: a small portron train set is carved out as train-deviset than the model is not trained on this.

Baya excer: Frain crease: Train deu ouror: Jodesa mamatah deu ouror: Jodesa mamatah deu ouror: Jodesa en ourofitting hat euror: came form some défuturion

#### Addressing date mismatch

-> identify manually how data differs, try to synthesize date

- There is a vish of overfitting rule the small out of dubts that we are synthesising Es All mags may look like cars to us but machine may just about to string ruje there as care.

#### Transfer learning

We take a number have be trained for a similar took and allor it to withe news repub. We allor the output layor(s) or add new layors as required & train them (free furring) if required their prev layors to (pre-training)

when is transfur learning used \$→if both tacks have some input hype → do to fir the task it was originally trained for was vory abundant & for now took the data 15 leas → Low level fasters fatures learnt in original model is useful for new model too

# Multi-task learning

-> Ex computer vision model which detects multiple of at same time, its ganna have I neuron for only of, statlen than having separate neurol networks for each object. -> even if data is missing labele it can be haved for whatever labels are available. { this is not orfinar this is logistic sequession for all neurons i.e coch neuron has output between 0 & 1 how combining all outputs gives 1

# -> when is multi-took learning good:

# - (22)-

#### End to end deep - learning

-> Ex sport wargerition took it can receiving several fatures then phonome then woods than transcript it, now this sequeres us to hand design various things. -> What if we want own machine to figure at all this on its own, then this is called and to only learning this can be done only when there is date in high anomals -> E. far requiring first zoone in and compatible for the term matches it with database

-> For higher complex tasks we need larger amount of data.

#### When to use end to end bouning?



→ Many a time a combinat P hand design and end to and kerning.

# Convolutional Neural Network

# Module

# Computer vision

 $\rightarrow$  The data input have is very large say a 64x64 ing then input is 64x64x3 which is a  $\rightarrow$  say a 1000 x1000 ing then its 3 million inputs, the NoN will have too many parameters and is seally difficult to  $\rightarrow$  Hence we use convolution operation.

## Edge detection example for convolution operations

→ Bay we a 6x6 rmg in grayscals input is 6x6x1. → we we a filler as keard of any 3x3. Hun when will be high motion

111121 (2)	the clements are multiplied new wire and summed c	P
CODE - INT	ostange box \$ 1x1 + 0x1 + 1x1 2.	
1000	÷ +	
	i <= 0×1 + 1×1 + 0×0 + +	*
	0X1 + 0x0 + K1 1	

 $\rightarrow$  why do we consulture a matrix?  $\Rightarrow$  it helps identify failures like edges et. E for simple filler which detects vertical edge in gasy stale image

the filter is moved continously to convolve a matrix

883H -	123	-188
121211		1.000
- B - F		- 11

Vestical edge filler has bright pixels on left 6 dasher ones on right

#### More edge detection





Musi of the times de in film one tobated as parameters if one deposit, this may detect more complex edges at voreas angle too

# tadding

→ mxn ing convolved with fxf fitter gives n-f+1 x n-f+1 matrix

+ the can only convolve a few times before ing gets too small

-> 2<sup>rd</sup> downside : pixels at edges are used very les compound to pixels in middle, hur a lot of natives overlap. .: we are throwing owney a lot of data at edges.

 $\rightarrow$  To whe both of the probs, we pad the ring with say p pixels on all sides, then the dimension becames  $n+2p-f+1 \times n+2p-f+1$ 



There are normally two choices to part: Valid & Same convolutions

I milding pod so that input size is common as output size i.e.  

$$m \circ m + 2p - f + 1$$
  
 $P \circ \frac{f}{2} = f$  is usually odd in computer usual

#### Studed convolution

-> Normally we just move one roow one colourn but in studied we move a row or a colourn alread

$$\begin{array}{rcl} nx & n & fxf & = \left\lfloor \frac{n+2p-f}{\delta} + 1 \right\rfloor & \chi & \left\lfloor \frac{n+2p-f}{\delta} + 1 \right\rfloor \\ & \text{pathing } p \\ & \text{strute : } s \end{array}$$

Crass correlation us convolution

what we are doing here is called exect-correlation in mathematical literature (Granked convolution in ML literature)

Technically when convolving mathematically

the filler is flipped vertically & horizontally than operation is performed

by doing this associationly property holds i.e (A+B)+C = A+(B+C)

But we do not need those for ML hence we do not do this

En if films is 1 2 3 h 5 6 mathematical 9 8 7 7 8 9 mathematical 6 5 h 3 2 1

Convolution over volumes

Suppose these are multiple channels like RGB et then how do we represent?

One layer of a convolutional network

Say we have f number of march files, how many trainable pasameters are present?  
(mal + j) f  
hotation of conductional layon l  
f<sup>(1)</sup> = filter size 
$$s^{(1)}$$
, shide Each filter  $a: f^{(1)} \times f^{(1)} \times n_c^{(1-1)}$  weights:  $f^{(1)} \times f^{(1)} \times n_c^{(1-1)} \times n_c^{(1)}$   
 $f^{(0)}$ , pabling  $n_c^{(1)}$ ,  $n_c$  of filters a chirabians:  $a^{(1)} \to n_{\mu}^{(2)} \times n_c^{(1)}$  bias:  $n_c^{(1)} > often taken at  $1 \times 1 \times 1 \times n_c^{(1)}$   
 $r_{\mu}^{(1)} \times n_{\mu}^{(1)} \times n_c^{(1)} \to n_{\mu}^{(2)} \times n_c^{(1)} + 1$  if using batch/minibatch of asse m then  $P_{\mu}^{(1)} \times n_{\mu}^{(1)} \times n_c^{(1)}$   
 $output: n_{\mu}^{(1)} \times n_{\mu}^{(1)} \times n_c^{(1)}$$ 

A simple convolution network example

#### tooling layer

This seduces the size of matrix by pooling various submatrices. 2 types of pooling:

Say there is 4x + matrix & a pooled with 2x2 filter



Max peol - takes the max all of submatrix > most commonly used Ang pool - - - ang of all ele in sub matrix > not often used

Hyperparameters while pooling				input: n <sub>H</sub> x n <sub>w</sub> x n <sub>c</sub>
t: Hites size s: stude	no	paxa metexs	to boun	output: n <sub>H</sub> 'xn <sub>w</sub> 'xn <sub>e</sub>
Max ex eve pooling				

Convolutional neural net work example



Say we try to process by ANN then

for mapping 32X32X3  $\frac{1}{n_c+6}$  28X28X6  $\Rightarrow$  parameters: 5X5X3X6+6 = 4556 3072 X H70H  $\Rightarrow$  parameters: 14450688 : two many parameters for a single layer which is preventing annull (AN N) image as 32×32.

Benefits of using convolutions



s the is converted with only I submatrix all allow all have no effect on this This makes madel has posed to overflit This makes madel has posed to overflit There have no a submatrix over if we shift a few pitch a cat to a cat the all haves loved features are some this can be humited very officently.

We build a CNN & use gradient descent to optimize pasameters to reduce J.

$$\frac{\mathcal{T}_{*}}{m} \sum_{i=1}^{m} \mathcal{L}\left(\hat{g}^{(i)}, g^{(i)}\right)$$

Module - 2

Why look at case studies?

Due to large no of hyperparenters its difficult to build a model from screetch here we use architectures similar to popular ones which have been built. Er LeNet-5, Alexnet, VG.G.

ResNet, Inception

#### Classic net weeks

LeNet -5	<ul> <li>Shi that hime signed Aanh ware used and not ReLa</li> <li>Computation was show horce filters baland at panillalar channels only</li> <li>con-linearity is inheduced after pooling layous</li> </ul>	; 60k paza meters
AlexNet	s had been reprint normalization is it would normalize the filter's output	≈ 60 M million parameters To anneal as we remered n.n.l.5
VGG-16	$\approx$ 138 M pasta meters	

VGG 19 wooks almost same as UGG 16 honce many use VGG 16

## Residual Networks (ResNets)

In this we stack sesicilial blocks to form the network

## Problem with going deep

when we initialise the neural network all weights are standomly initialised

When data goes alread by the time it searches the last layor its seduced to soundom noise due to multiplications with soundom numbers (initially)

The heas which is computed at end through bock peop to update gradients, new at later layous data was sceambled too much hence updating gradients there doesn't matter much Now the update to initial layous also won't man much because their gradients have been sceambled (i.e. w. w. x.D)

Hence deeper networks take long time to learn.

How to solve this ?

→ We must find a way so that more meaningful data reaches lake layers & make their inputs morningful → · · · · · · · loss gradients to arrive at early , to - - update -

We establish skip connections where data goes through & around blocks

We combine output & data at end giving it two paths to fillow

ow do we combine? sp (1) we take both tensors & add them element wise as concatenati them

### Ad vantages

-> each block sugments the date. Now the work of each block is to figure out what it can add to input softwe than figuring out what the input is

-> Even if we are adding data instead of concatorating initially the weights are contend around zero so that initially we are passing data selectively unchanged

→ .: Each block has simpler things to learn, has access to better info

### -> Shostor gradient paths

- intrating weeks greating are nearring here we can get good update. From shores party of the as the model lows more morning ful weaks are to form longer party. Here has always if Here has always if Here has no

-> Madularity is its easy to add blocks.

#### Concerns

#### Shape mismatch

if only ANDI is involved is use have to match dimensions, so we need some way to stephage things. if CNN is also involved is no match, there are constraining on no Einy . We cannot keep analonating else the orthistion tensor will get too by & there will be an explosion of postameters there if there are too many blocks we greefe addition over consideration & there are savely ship american use which we addition continuesly.

#### Shape no thing

→To simplify row blocks we can use set 6 valid publing so that My GMo is preserved →To adjust depth we can use 1×1 Film with stepsend depth so that we can adjust the available depth





Now even if the learning is insignificant have the output will still be equal to a<sup>(C)</sup> itself have not making model any waves Next of lon, the pseulern is the learning become unuse as we go depar this is suicided have

We is vertex multiplied to match shapes, it might be learnable or just a fixed matrix to pad

Lets take on example  $X \longrightarrow Big NN \longrightarrow a^{CR2}$ 

X -> Big NN-----

Now,  $a^{(l+2)} = g(z^{(l+2)} + a^{(l)})$ 

say g() is ReLu



Network in Network and IXI convolutions

if we want final output to how depth it than 
$$\frac{\omega w}{\omega}$$
 do use use that?  
we doe of fitters  $\rightarrow$  to t or J n.  
 $\rightarrow$  even if we decide to keep  $n_c$  some this stilladds non-tinenranity hence will learn something.

#### Inception network motivation



#### Mobele Net

 $\rightarrow$  This is computationally less heavy, hence we may shun this on devices like mobile phase too.

→ kay idea: Normal us depthwise separable convolutions

ie cast = #filter passametres x #filter pasitions x #filters

$$n_u \times n_\omega \times n_c + f \times f = n \times n \times n_c$$



The cost seduced is genesially : 
$$\frac{1}{n_c} + \frac{1}{p^2}$$



## Advantages of mobile not v2

" it has exp" & depth use the model knows sicher functions these and then due to projection layer the menory taken is also soluced have deployable on small devices.

How to scape up as down a model based on device its deployed on?

Alter resolution of image, depth of NN, ar make layous wider.

ue have to take can of trade off & on planent this. [Can have a look at paper mentioned in notes for further info]

Freeze a few layers & train others or depends .....

- -> Rotating, Sheating, Local watering
- -> Apply PCA on RGB channels and alker them
- -> When data is too huge then data is token from hardbick, distortions are applied 6 this is treaked as botch and preserved

when there is little data hand engineering is smally impostant & hade like through domning plays major sole



Lets say we have to classify into

2- cour 3- motorcycle 4 - back ground this output whether the object is present or absent & various landmarks. Ex we need baselies of eyes, mouth position of ears ct.

i Lon In y

n land masks





A sliding window exaps parts of ings of model hats if abject a pascent in it; then some is done for bigging size windows

Downside : computation cost



All bounding books predictions are seen the one with highert one is prival 6 all other boxes which show IoU higher than 0.5 or 0.7 are eliminated



#### Achor boxes

Say two objects then, how do we deal with it?



to form ing. Here it classifies pixels to cart Exnot cat

l sanspase convolution

How do we increase the size of matrices again ?





pet 1st els. is 2, maltiply with all ele mfriter & place it on output; where there is no padding write the value there

		2				
				•	1	
2	•	2				
_						
Ы	- R	2				
-			J			

Then second ele. I . shide is 2 shift filter accordingly

	2	1	1.1										Γ
		_		۰	·	•	1		•		•	1	
2	0	1		4	<b>*</b> *e	2	Λ.			2	2	A.	1
-													
0	2	1											1
-	-	-											

i.e it knows there is a calt in the stegriou

go on doing same to fill matrix

takes exact region where cat is present

#### U-Net aschitecture

it utilizes skip connections from earlier layous. The pasent layou has high level contextual info 6 spatial info is taken from skiponwertion



ve are adding as can caterating

here so dimensions vory occordingly

What is face - second nition ?

#### One shot learning

Lease from just one example to secondly presen again  $\frac{Similarily fore}{d(ing 1, ing 2)} = degues of diff to img <math display="block">\mathbb{T}f = d(ing 1, ing 2) \leq \mathcal{V} \qquad \Rightarrow \text{ same provision}$ 

>2 » diff.

#### Signese network

An encoding is created out of an image through the model & d( ) is plotted

$$d(\mathbf{x}^{(i)}, \mathbf{x}^{(i)}) = \left| \left| f(\mathbf{x}^{(i)}) - f(\mathbf{x}^{(i)}) \right| \right|_{\mathbf{x}^{(i)}}$$

\* # # + |

→ Duxing tstaining if A.PGN are chosen sometimity then has fire is aveily solufied f model may not have much.

→ While testiving model we need multiple ing of a Person but once its fully trained to addrew prople to database one-shot learning can be used.

$$\int_{1}^{\infty} \int_{1}^{\infty} \int_{1$$

#### Face verification & binary clasf"

Ercade the input ing. if preading motether do before output 1 else O > alternative to triplet lass

$$\begin{array}{c} \hat{\mathfrak{g}} \cdot \boldsymbol{\nabla} \cdot \left( \sum_{k=1}^{20} \left| f(\boldsymbol{x}^{(i)})_{k} - f(\boldsymbol{x}^{(i)})_{k} + b \right) \right. \\ \\ \hat{\mathfrak{g}} \cdot \left( \frac{f(\boldsymbol{x}^{(i)})_{k} - f(\boldsymbol{x}^{(i)})_{k}}{f(\boldsymbol{x}^{(i)})_{k} - f(\boldsymbol{x}^{(i)})_{k}} \right)^{2} \Rightarrow \mathcal{N}^{2} \quad \text{similari} \\ \end{array}$$

Database encoding can be percomputed & stored at a place to make things faster.

### Neusal style teamsfer

Content - (c) Grenestative - (Gr) Starle - (s)

what are deep convinces learning?

Initially it ares things like only as later on it should recognizing complex shapes

Cast function

$$\begin{array}{c} J(G) = & J_{control}(C,G_1) & + \beta J_{shyle}(S,G_1) \\ & & & & \\ & & & \\ & & & \\ & & & & \\ & & & \\ & & &$$

how similiar is content of oxiginal & generated img is?

-> Initialize G randomly

 $\rightarrow$  Use gradient descent to minimize  $J(G) \Rightarrow G_1 = G_1 - \frac{2}{3G} J(G_1)$  is use one updating pixel values

#### Content cost function

-> Say a hilder layer I is computing content cells, this layer cannot be at very early shape cay it will by to replicate actual ing & cont be to late either else it will beer spotial info

-> Use pretrained conv net Ex VG161 & compute a (23(c) & a (23(6)), if they are similar then ings are similar

$$\Rightarrow \overline{J}_{comband} \left( \zeta_{3} G_{1} \right) = \frac{1}{2} \left| \left| \alpha^{\left( \zeta_{3} G_{1} \right)} - \alpha^{\left( \zeta_{3} G_{1} \right)} \right| \right|^{2}$$

Style cost func

Say we are using layer is activation to measure style. Style is correlation between activations across channels

If they have similar values then both of those features tend to occur together

G(1) =  $\sum_{i,i}^{\infty} \sum_{j,i}^{\infty} a_{ijk}^{(L)}$  (L) . These matrices are will be large else small (Unnormalised cosselation). These matrices are called gram matrices are

$$\int_{\frac{3}{4}}^{(10)} \int_{\frac{3}{4} \text{ gressifive ing.}} \int_{\frac{1}{24 \text{ gressifive ing.}}}^{(10)} \int_{\frac{3}{4} \text{ gressifive ing.}} \int_{\frac{1}{24 \text{ gressifive ing.}}}^{(10)} \int_{\frac{3}{4} \text{ gressifive ing.}} \int_{\frac{3}{4} \text{ gressifive ing.}} \int_{\frac{1}{4} \text{ gressifier ing.}} \int_{\frac{1}{4} \text{ gressifier$$

hypespasameter

Why sequence models?

1

Notation

Example: Name entity secognition

- x: Having Pottes & Hestmoine Gitanger truented a new spell  $\chi^{\chi_1>}$   $\chi^{\chi_2>}$   $\chi^{\chi_2>}$

 $y^{(i)} = i^{th} \text{ output}$  $y^{(i)(i)} = i^{th} \text{ ele. of } i^{th} \text{ output}$  $\overline{Ty}^{(i)} = \text{ length of } i^{th} \text{ output}$ 

thou do ue separant usade in serience? - We come up with vocabulary (dictionary) normally commonstal apprications have 50k - 10sk words , bigga make may have it in millions

 $\rightarrow$  now we use one hot anceading to separate the world

→ .: each x<sup>42></sup> is a one hot encoded vector

-> What if there is a new work which is an acountred? We create a new token called unhown word <UNK>.

Recursent Neural Network Model

why not a standard network?

→ Inputs & outputs can be diff lengths in diff examples

-> Doesn't share features learned across diff positions of test.



RNN is also sepresented this way

y. 1 if it post of name else O

- Here it sections data from input & prev activation to predict output. -> RNN only secreme data from previously scanned words and not the words that come after it, Bi-distributional RNN (BRNN) does it this way
- -> why is it important to service data from woods after ? Ex Teddy came to home. Totady is name
  - I want today bear. ? - rot a -

... on mature multiplication of these we get Wax X<sup>112</sup> + Waa a<sup>16-12</sup>

Back propagation through time



# Different types of RNNs

Name entity secondition is many to many aschitectuse, it many of Grany i/p Servisional cleaft is many to one.

Music generation is one to many.

Machine translation is many to many many has leparts an encoder & a decoder

#### Language model and sequence opnesistion

held say there are discriteries : The apple & pair salad. - 5,

S2, is casen't but they are permaned cracitly the same, how does markine differentiate? By Language model. > It outputs the pseudolity of likeliness of that centence

P(y<sup><1</sup>), y<sup><1</sup>,...., y<sup><1</sup>) is the output

To build such a medial we need a donge data set ox a Corpus → Fixet we takenize the sontence [KEGS> » and of sontence taken]

→ Next we build the model



: RNN learns to predict one world at a time.

 $\mathcal{L}\left(g^{\text{AD}},g^{\text{AD}}
ight)_{\text{F}} = \sum g^{\text{AD}}_{t} \log g^{\text{AD}}_{t}$ 

 $\mathsf{P}(y^{(s)},y^{(s)},g^{(s)},\ldots,g^{(s)}) = \mathsf{P}(y^{(s)}) \mathsf{P}(y^{(s)}) \xrightarrow{} \cdots \mathrel{\mathsf{P}}(y^{(s)}) \xrightarrow{} y^{(s)}, y^{$ 

#### Sampling Novel Sequences

After we train a model, to get stough idea of what it has learned we can somple a sequence



-> the output from a particular, with is chosen at retardorn from represendent then fait borows whith made to pendict -> When to Stop? Inf discour cases there only its generated also, with some fried <u>on of wards</u>

-> If we do not want RUNK> to be a part of of them we can code it so.

Character level have grage madel Instead of works it was helper as if Eagle. Its samply used E is more complex than normal madels

Recap: in back peop we use chain sub multiply many desirations, it many smallvalue have as we go back further the gradients here getting smaller

In sentences a word which came long ago influences word which come long fine office , there, we need to find a way to overcome vanishing gradient problem Explading gradient will just give NoN values & a efficient sol is dip gradients.

i.e if gradients go above a certain theshold then we neumalize it in some way

$$\frac{1}{3} \frac{1}{3} \frac{1}$$

The cost, which already ate ....., was full

GRU with has a new variable called 'c' which study for menory cell.

$$C = a^{2r}$$
every time step we case write  $c^{4r}$  as  $\left[ c^{4r} - g \left( u_{\perp} \left[ c^{4r} + z^{4w} \right] + b_{\perp} \right) \right]$ 
then we have a gate  $(\Gamma)$  & has value be  $0 \in 1$ 

$$\Gamma_{u} = \nabla \left( u_{\perp} \left[ c^{4r} + z^{4w} \right] + b_{\perp} \right)$$
is stands for apolate gate
gate devices by how much we are going to update the cell value

$$c^{(4)} = \int_{u} * c^{(4)} + \left( 1 - \int_{u} \cdot c^{(4-1)} \right)$$

Tu generally is O as I i.e there is high as very very low update

i if C<sup>333</sup> is not updated it almost eemains constant, hence addressing vanishing gradient problem.



$$\int_{K} C^{44} = g\left(\omega_{c}\left(\Gamma_{8} * c^{4+1}, x^{4+2}\right) + bc\right)$$

$$\omega = \int_{\alpha}^{\alpha} = \nabla \left( \omega_{\alpha} \left[ c^{(4+1)}, x^{(4)} \right] + b u \right)$$





Vaxiations » peephole connection



We can just take LIST M units & join many fem of make a model.

## Bi disectional RNN

-> Lets us take info from after E before

- -> Info from start goes to next & from the end it comes to start too.
- -> Disadu. » it reads to have complete sequence of data available to start processing, most applications do have complete sequence available beforehard itself except things like speech to text where entries sequence may not be present, these use more complex versions of BRNN.
- -> Normally the go to choice for NLP is LSTM with BRNN



Lo uni diver for word connection

## Deep RNNs

-> Normally toomany RNN layous assort stackard on top of each other cry they are computationally very expressive as it temperal dimension is high.

- → A few RNN layers followed by normal ANN type of networks are pretty common.
- > Here instead of RNN, GRU or LST M with RNN or BRNN can also be used.



#### Word representation

-> Make model understand analogies like man: woman :: king: queen

→ Bie visius ly we septemented woods as one-hot-encoded vectors. Is these any disade to it? Howe we cannot astabilish any selationship between woods

Er king & queen are more solubed than king and apple. Er I will have assange \_\_\_\_\_\_ jure seems to fit both & if some how modine how apple & asange schold then it \_\_\_\_\_\_ apple \_\_\_\_\_\_ could learn this more easily.



Learning word control in the all words before Three are control tanget point A or a few words before A or a few words before & after A or a nearby 1 word all of three work well

Word 2 vec

Ship- grams

-> Come up with few context-target pairs to create supervised word problem

- choose a vandom workt as contest & choose a vandom workt as target 6 tay to peolicit target have use are just taging to learn workt embeddings 6 not make on model goed a peolicity (maybe in a small worker)

$$\mathcal{L}(\hat{\mathbf{G}}, \mathbf{y}) = -\sum_{i=1}^{n} \mathbf{H}_i \log \hat{\mathbf{H}}_i$$

 $\frac{Peoblem s_{computation}}{\sum_{j=1}^{n} e^{Te_j}} \quad \text{is very heavy especially if these are many worked}$ 

Sol<sup>n</sup>: using a hie souchial softmax classifier isk

in practice hierarchial classification common woods at the G which man doep down the trap

We can choose 't standardly but that we will end up mapping woods which accuse too frequencity which is not good. There are many ways in which woods are chosen carefully and mapped

#### Negative Sampling

Pick a contract target from sentence a label it as I that our tre enample & pick other words at random from compus & assign you to then ; say we do this to generall k negative enamples of we have very large data then choose small h, maybe 2-5

- - - small - - - - loggeh - - 5-20

then use supervised model to lease to map from a tot

# $P\left(y^{=1}\left|c,t\right)=\pi\left(\boldsymbol{\varphi}_{t}^{\mathsf{T}}\boldsymbol{e}_{c}\right)$

now we train only three & examples instead of all 10k azamples in one-iteration

-> how do we choose -ve examples ?

Sample according to emperical fug of woods i.e of to how ofter words appear, but this ends up taking woods like a.a.n. the et

- or randomly sample using i which is other end of extreme

► Global vectors for word separatorian  $X_{ij} = \stackrel{\text{model}}{=} f_{ine}^{j}$  appears in context of i  $f_{\pm}^{j}$   $\frac{M_{odel}}{I}$ minimize  $\sum_{i=1}^{loc} \sum_{j=1}^{loc} f(X_{ij}) (\Theta_{i}^{T} e_{j} + b_{i} + b_{j}^{'} - \log X_{ij})^{2}$   $f(X_{ij})$  is weighting tran  $f(X_{ij}) = 0$  if  $X_{ij} = 0$   $f + f_{x}$  would like this, the make sure its net too high  $f_{i} \in F_{x}^{j}$  - durian (appear saidy) it net too how Roles of  $\Theta_{ij}$  has an symmetre  $e_{w}^{j} = \frac{e_{w} + \Theta_{w}}{I}$ 

we can't guagantee that the individual embeddings are interprettable

#### Sentiment clasf"

Recognize would & any there e, vectors & sun soft max on it to pseulicitemotion



also as pont viewiew coz waedzond appears so many fing a good alternative

Problem this will classify last comment

1111

Debiasing word embeddings

Gender , ethnicity bias which is learnt must be debrased

is to use RNN for this

How do we do this?

They there is gender bise, we what for words in which bias is vight like the star , they, they, the any them of find there - grand is this is bias dis

-> All other die" are non biss die"; all the words which need to be debiard are projected to non-bias die" > this is called Neutralization

 $\rightarrow$  Equilize points is point like  $e_{log}$  (equil put them at separal difficm non-line dist




### Basic Models

encoder m/w

# -> Say we wanna build a machine translation model, howdo we build it?

-> Bay we have an image we have to generate a caption or comment on that





-> The cotch in these both is we want the most likely answer, not the standom answer

which gp I wood at a time

decides m/w

icking the most likely sentence

-> if we have to pick met also not pick it geterabily? it doesn't wark out well can mare after will be labelled meet likely own if do the best translation

its called as conditional larguage mode cay second half is just a long, model but indeed of feding gammar standom vertees are feel it a vector which captures input enterce



Beam seasch

-> This has a hyperpresenter called beam width (10), it always picks top B world & evaluates next word for each prewholity & next out of all these again picks top B.

-sayones ciscled in blue are picked, then it eliminated sept. as starting wood.

(B·3) wh B·3) wh B·3) wh B·3) B·

 $P(y^{\alpha}, y^{\alpha}|x) = P(y^{\alpha}|x) P(y^{\alpha}|x y^{\alpha})$ it evaluates 30k possibilities in one zen

& B copies of model will be made to evaluate them.

Kefinements to Beam search

 $\rightarrow our flue was any max TTP(y^{(1)}| 2, y^{(2)}, ..., y^{(1-3)})$ , now all there are less than one, if we have a long sentence then computer will base track of definal values & suffer from sounding off peoplem. How do we exercise this?

 $\Rightarrow$  instead of hopping track of product we take log & keep track of products  $\Rightarrow$  and  $\sum_{t=1}^{t} \log P(y^{th})(x,y^{th}), \dots, y^{th})$ 

 $+ \operatorname{Finally} \text{ we also noticalize it by multiplying } \underbrace{1}{\mathsf{T}_{y}^{\mathsf{X}}} \text{ where } \mathsf{X} \text{ is hyperparameter defined by us, if } all 0 > no neormalization <math>\mathfrak{g}$  we choose a balanced value:  $a(1) > \operatorname{Complete} u$   $b_{1}$  tween  $0 \in 1$  like  $0 \neq or$  $avg \max\left(\frac{1}{\mathsf{T}_{y}^{\mathsf{X}}} \sum_{t=1}^{\mathsf{T}_{y}} \log P\left(\frac{it}{y} \mid x_{1} y^{(t)}, \dots, y^{(t+1)}\right)\right)$ 

Esses analysis on Beam search

- -> What if beam search is making a mislake, erris it own madel that's making the mistake. How do we identify?
- $\rightarrow$  Lats say there is a sentence which is translated by human (y\*)

	And apple of the second	r evalua
and the second s		if
all such as a second se	Martin .	

- evaluate what X of excess is due to Broom search is what X due to RWN

if peop is in Beam search we can t Beam width or allow normalization

- - - RNN we have to evaluate our architecture & apply vorvious ideas in course 3

Bilingual evoluation understudy

- > it gives high score if the output is close to human translation. This ack as an alternative to having a human evaluate the output
- say we have 2 references & both are correct, we have to score our output.
- -> MT output = machine translation ofp

So, placision of our op is  $\frac{1}{2}$ , this is a bad scoring size, hence we clip the world count to max no of occurrences in both sonkince the world 'the' occurs twice in fract sentence : we clip it to 2

how many woords match words in either of ref embences.

we call this modified precision.

 $\rightarrow$  clipped counts , now its demonstrated considering 1 word, we can consider n-words too (n-gram) ( called Linighton )

-> How do we score ? Frist we compute clip counts

- → Say there is a seally long sentence, then its hand to privalist the sentence in such cases the machine gives bad off, how to modely the?
- -> there is the machine translating at present? It looks at a sentence, memorizes it entirely, then translates. Now, how does human doil? We look at sentences in parts then translate it





A separate RNN computer the off is its viscious weights from french waves too which specifies how important is that wood to approach that specific output

s how important ist' to generate t

## →BRNN



-> hese attention weights ase computed using a small NN



BRNN



Aix pressure us time graph & audio data false back ofp ?

 $\rightarrow$  Aspectro gram of how intense sound is with frequency is prepared as pre-precessing step.

- How do we go about building speech to text age?

We can do this by attention maded itself, where it looks at a certain time from & outputs the characters / words



Say audio input was of freq. 100 herter is it 100 if per second . The characters are generated continuesly Ea special characters interduced '\_' a called blank

En The quick .....

Collapse all char not separated by '\_ into one . this becomes the q



RNN processes the data say ventrally. Trans formers processes data parallely

## -> Attention + CNN

ightarrow Self attention $ ightarrow$ say we have free workeds in sentence its computer S :	sepsentations for it.	w, w, w are weights
ightarrow Mulli head $ ightarrow$ for loop even self attention process		9/2 W 2 x 13
Self-Atlantion	also called scaled dot psoduct septreentation admonimeter is scaled to prevent it from exploding	κ <sup>63</sup> •ω <sup>k</sup> * <sup>65</sup> ν • ω <sup>ν</sup> * <sup>65</sup>
$P_{i}\left( \mathbf{a}_{j},K,V\right) \ast$ a thertion - based vector for dependentiation of a world.	ay∶quury, its like ashing whate happening these? h≥hay	
$P(a_{v}, K_{v}) = \sum e^{e^{a_{v}, K^{t}}} v^{$	VE Velue	
1 5 e <sup>(*,1,4)</sup> >	ark is like an usering the quastion	
Z, C	Now, embedding is done considering sussounding "	words.
	Value Jatte of A Device J A Device	. J. N. B. J. NP.

Each time we calculate self attention is called head

·	Second question : ie Self atbention for 2nd time	# head = 3 J	all find answers of hime of heads as
	Address (Assessed of the second secon	annellenn -	stocked toge thes & given as of
712.33	and the second s	·	
Pro Prode James marinella and	Annalise Annalisation	Friday Friday and	Though it appears like we are iterati
pos prese time and to an			to calculate these we can calculate all the

 $\omega_{i}^{\theta}, \omega_{i}^{h}, \omega_{i}^{h}$ E we compute Say highest softmax for W is highlighted with blue awrow



sun for n-times

> initially only <sas> is given, then sun n-times next word is generated

then reason next word both one given crits sun

m-fimes to predict next & so on.



head; Attention ( we Q, with, wiv)

pasablely.

Details



pig x200 are added lithey contain contextual semantic embedding & positional encoding info

 $\rightarrow$  Along with this we use sessibled network too to pass on politional info

- Transformes also uses a layer very similar to back norm, its purpose is to pass along partional of

-> ~ ~ ~ called adda norm, similar to back norm

-> Mark multihead attention, this is used mainly while training, has it has access to consect translation too what it does is it hides a world & prestored crowything before it predicted consect -y

## → Best, Best Distill are add versions of this model.



, along these positions which are dotted, the values are read for phi<sup>3</sup>