

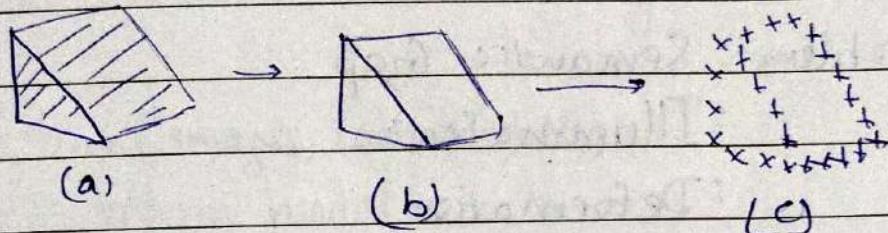
# "Lecture 1"

Computer Vision - interpreting visual data.

Majority of data (70%) is visual  
Problem: hard to analyze.

Computer Vision  $\rightarrow$  useful in all fields

Block World (Laasya Roberts)



Input Image  $\rightarrow$  2D sketch  $\rightarrow$  2.5D sketch  
 $\downarrow$

3D diagram

Overview

Primary focus  $\rightarrow$  image classification

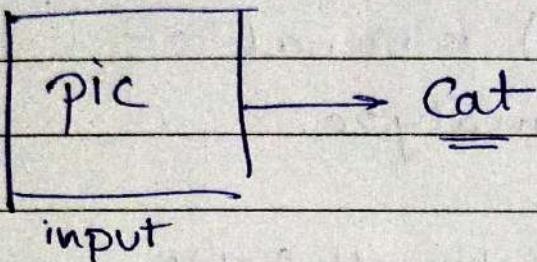
Other "  $\rightarrow$  object detection, action

classification, image captioning

Start with CNN

## "Lecture 2"

### Image Classification



computer sees a grid of numbers  $\in [0, 255]$

$$w \times h \times RGB \rightarrow 800 \times 600 \times 3$$

Problem : Semantic Gap

: Illumination

: Deformation

: Occlusion

: Background Clutter

: Inter-class variation.

Image classifier :

```
def classify-image(image):  
    return class-label
```

There is no obvious way to recog a cat

## Attempts

Pic → Edges → corners → classify  
(Poor at classifying)

## Data-driven approach:

- ① collect dataset of images & labels
- ② Use ML to train a classifier
- ③ Evaluate classifier on images

def train(images, ~~targets~~ labels):  
 return model

def ~~etc~~ predict(model, test-img):  
 return test-labels.

## First Classifier: Nearest Neighbour

\* train → memorise all data/label

predict → predict label of most similar  
training image.

Distance metric for comparison:

$$L_1 \text{ distance } d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

test img      training img       $d_1$

32	2
7	91

 - 

14	2
17	82

 = 

18	0
10	9

 $\Rightarrow [37]$

K nearest neighbor

→ instead of copying label from nearest neighbor we take majority vote for K pts

$$L_2 \text{ distance } d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

→ use for random coordinate systems.

Hypoparameters

→ best value of K to use

→ best distance to use

These are choices we set before.

try all, see which is best.

## Setting hyperparameters

- X #1: choose hyperparameters that work best on data
- X #2: divide into train & test and choose best which works on test.
- ✓ #3: Split into train/test and validation, choose hyperparameters on val.
- ✓ #4: Cross-validation - split data into fold.

(only for  
small data)

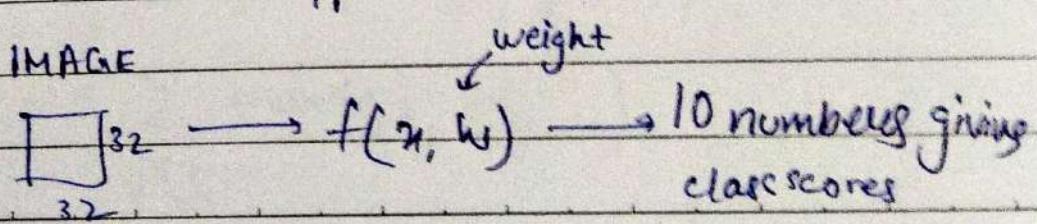
f1	f2	f3	f4	f5	test
f1	f2	f3	f4	f5	test
f1	f2	f3	f4	f5	test

NOTE: WE DONT USE K NEAREST NEIGHBORS

- i) Curse of dimensionality
- ii) Pic may have diff colors  $\rightarrow$  large dist, but same pic.

## Linear Classification

Parametric Approach:



# "Lecture 3"

## To Do

- ① Define a loss function that quantifies unhappiness
- ② Come up with a way of efficiently finding parameters that minimize loss fx" (optimization)

Suppose: 3 classes (cat/can/frog)

A loss function tells us how good our classifier is given a dataset,

$$\{x_i, y_i\}_{i=1 \text{ to } N}$$

Loss over the dataset is a sum of loss over examples

$$L = \frac{1}{N} \sum_i L_i(f(x_i, w), y_i)$$

## Multiclass SVM Loss:

Given an example  $(x_i, y_i)$  where  $x_i$  is the image and  $y_i$  is the label.

and let scores be  $s = f(x_i, w)$

$$L_i = \begin{cases} 0 & s_{y_i} > s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

NO.

Date

foreg,	cat	3.2	[Correct]	cat	1.3
	car	5.1		car	4.9 [Correct]
	frog	-1.7		frog	2.0

$$L_{\text{off}} = \sum_{j \neq y_i} \dots$$

$$\begin{aligned}
 &= \max(5.1 - 3.2 + 1, 0) + \max(0, -1.7 - 3.2 + 1) \\
 &= 2.9 + 0 \\
 &= 2.9
 \end{aligned}$$

$$\begin{aligned}
 L_{\text{car}} &= \max(1.3 - 4.9 + 1, 0) + \max(0, 2.0 - 4.9 + 1) \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

NOTE: Mean also works

Can't use square as it is non-linear

Supervised formula:

$$L(w) = \frac{1}{N} \sum_i L_i(f(x_i; w), y_i) + \lambda R(w)$$

Data Loss
Regularization

## Regularization (Penalty)

$\lambda \rightarrow$  Regularisation strength

$$\begin{array}{l} L_2 \text{ Regularisation: } R(w) = \sum_k \sum_L w_{k,L}^2 \\ L_1 \quad : \quad R(w) = \sum_k \sum_L |w_{k,L}| \end{array}$$

$$\text{Elastic Net } (L_1 + L_2) : R(w) = \sum_k \sum_L \beta w_{k,L}^2 + |w_{k,L}|$$

## L2 Regularization (wt Decay)

$$\pi = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [.25, .25, .25, .25]$$

$$w_1^T x = w_2^T x = 1$$

$$R(w) = \sum_k \sum_L w_{k,L}^2$$

## SoftMax Classification (Multinomial Logistic Regression)

Scores = unnormalized log probabilities

$$P(y=k|x=x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}, \text{ where } s = f(x_i, w)$$

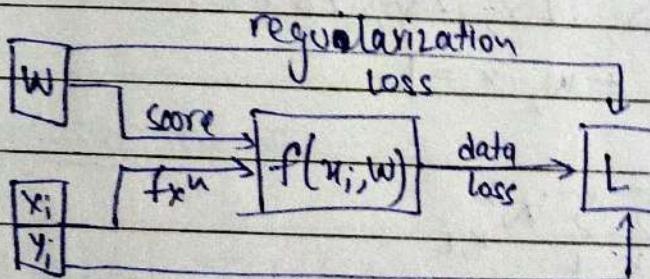
Want to maximise the log likelihood to minimise the negative log likelihood.

$$L_i = -\log P(Y=y_i | X=x_i)$$

$$L_i = -\log \frac{e^{\xi_j}}{\sum_j e^{\xi_j}}$$

foreg,	cat	<u>3.2</u>	24.5	0.13
	cat	5.1	$\xrightarrow{\text{exp}}$ 164.0	normalize 0.87
	car	-1.7	0.18	0.00

$$L_i = -\log(0.13) \\ = 0.87$$



Strategy #1: Random Search for  $w$

**BAD**

Strategy #2: Follow the slope  $\left( \frac{df(w)}{dw} \right)$

<u>Current w</u>	<u>(w+n)</u>	<u>grad w</u>
0.34	0.34 + 0.0001	-2.5
-1.11	-1.11 + "	0.6
0.78	0.78 + "	0

[BAD]

We need  $\nabla_w L$

then compute the analytic gradient  
 (this is faster)

Stoch.

Stochastic Grad Descent

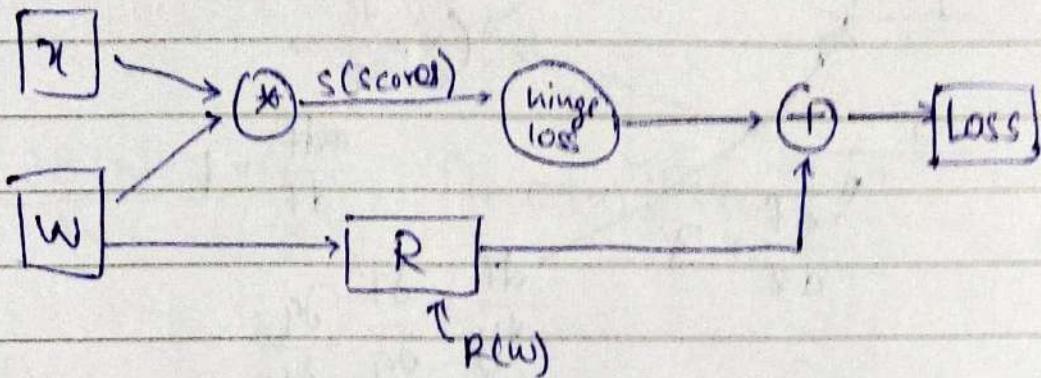
$$L(w) = \frac{1}{N} \sum_i L_i(x_i, y_i, w) + \lambda R(w)$$

$$\nabla_w L(w) = \frac{1}{N} \sum_i \nabla_w L_i(x_i, y_i, w) + \lambda \nabla_w R(w)$$

# "Lecture 4"

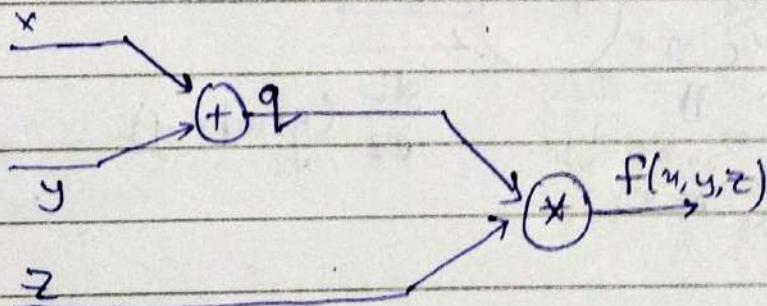
## Computational Graphs

$$f = Wx \quad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



## Backpropagation

$$f(x, y, z) = (x+y)z$$



$$q = x+y \quad \frac{\partial q}{\partial x} = 1 \quad \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z \quad \frac{\partial f}{\partial z} = q$$

$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$

No.

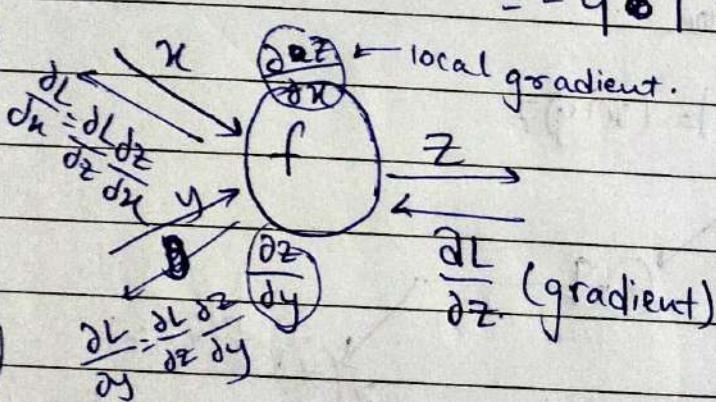
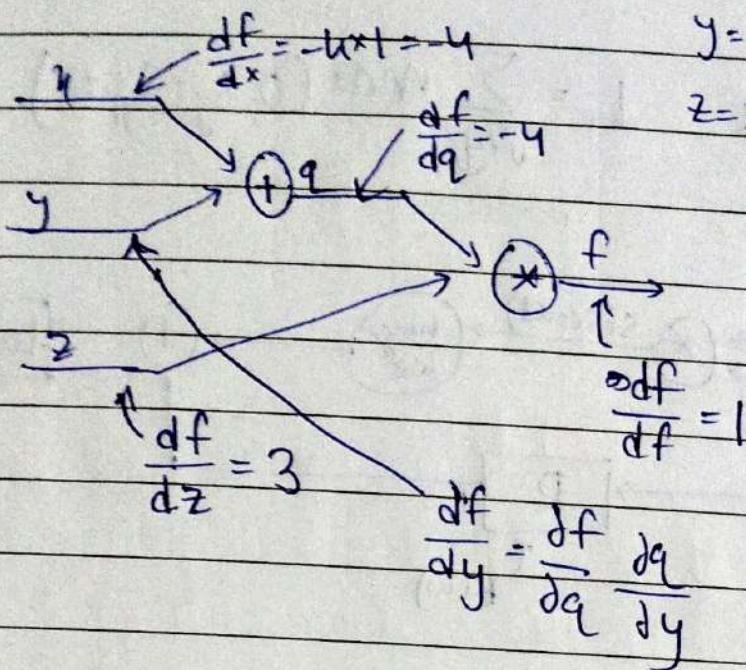
Date

Start from back

$$\lambda = -2$$

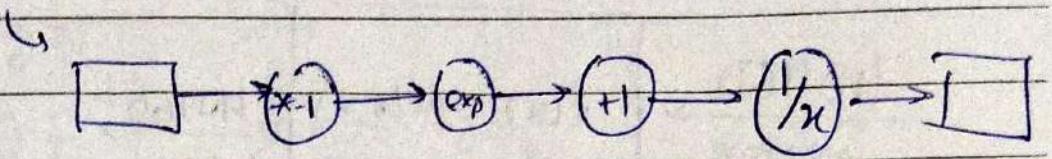
$$y = 5$$

$$z = -4$$



Another example.

$$f(x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Sigmoid gate

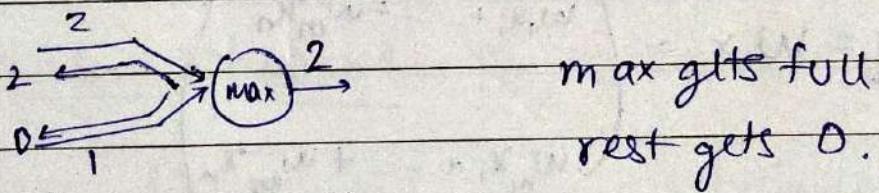
$$\sigma(u) = \frac{1}{1 + e^{-u}}$$

$$\frac{\partial \sigma(u)}{\partial u} = \frac{e^{-u}}{(1 + e^{-u})^2} = [1 - \sigma(u)][\sigma(u)]$$

Patterns in backward flow:-

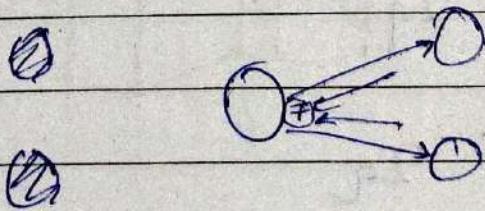
add gate: gradient distributor

max gate:

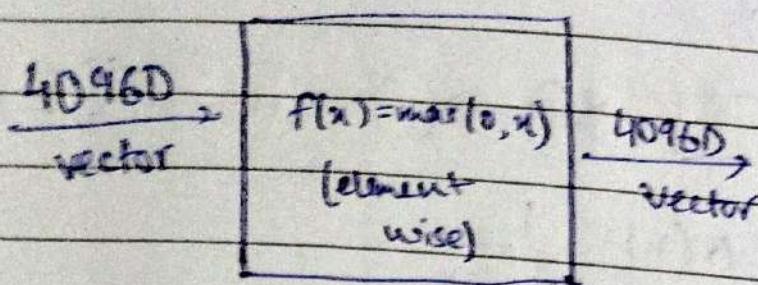


mul gate: gradient switcher

gradients add at branches:

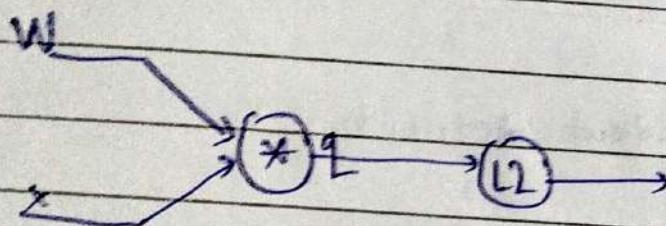


## Vectorised operations



### Example

$$f(q) = \|w \cdot x\|^2 = \sum_i (w_i x_i)^2$$



$$q = w \cdot x = \begin{pmatrix} w_1 x_1 + \dots + w_n x_n \\ \vdots \\ w_n x_1 + \dots + w_n x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$\nabla_q f = 2q$$

$$\frac{\partial q_k}{\partial w_{ij}} = l_{k,i} x_j$$

$$\begin{aligned}\frac{\partial f}{\partial w_{ij}} &= \sum_k \frac{\partial f}{\partial q_k} \frac{\partial q_k}{\partial w_{ij}} \\ &= \sum_k 2q_k (l_{k,i} x_j) \\ &= 2q_k x_j\end{aligned}$$

Similarly,  $\frac{\partial f}{\partial x_i} = 2q_k w_{k,i}$

$$\nabla_x f = 2W^T \cdot q$$

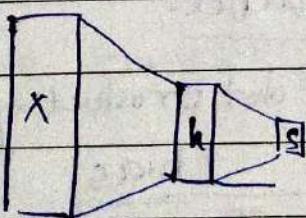
## Neural Networks

Before:  $f = Wx$

[single layered]

After:  $f = W_2 \text{Max}(0, W_1 x)$

[double layered]



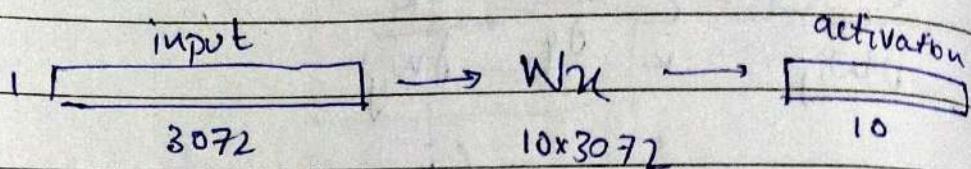
NO.

Date

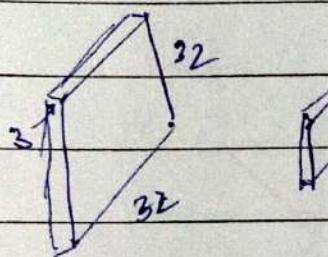
## "Lecture 5"

### Fully Connected Layer:

$33 \times 32 \times 3$  img  $\rightarrow$  stretch  $3072 \times 1$



### Convolution Layer

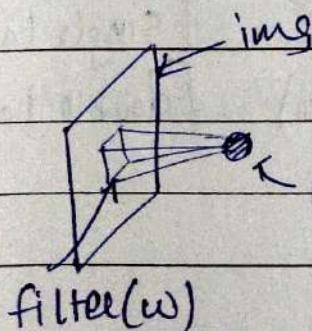


$5 \times 5 \times 3$  filter.

convolve the filter with  
the image

depth will be same

but, spatial area goes down

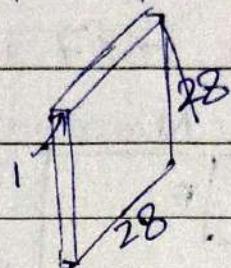


1 number:

dot product of filter & img  
+ bias

$$[w^T x + b]$$

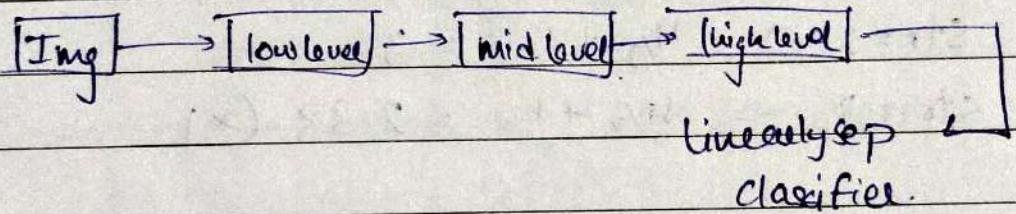
convolve  
over all loc



activation map

As many filters, that many act maps.

ConvNet: sequence of conv layers, interspersed with activation functions.

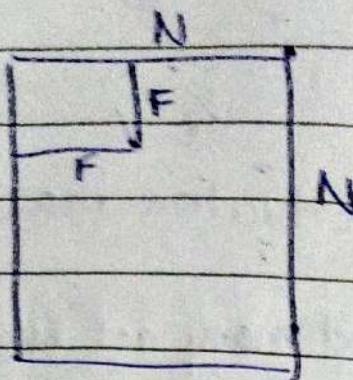


We call layer convolutional coz it is the conv. of two signals.

$$f[x, y] * g[x, y] = \sum_{n_1} \sum_{n_2} f[n_1, n_2] \cdot g[x-n_1, y-n_2]$$

NN will have conv, RELU, conv, RELU, Pool structure.

Let there be  $7 \times 7$  input and  $3 \times 3$  filter



$$\text{Output size: } \frac{N-F}{\text{stride}} + 1$$

$$\text{str} = 1 \Rightarrow \frac{7-3}{1} + 1 = 5$$

$$\text{str} = 2 \Rightarrow \frac{7-3}{2} + 1 = 3$$

$$\text{str} = 3 \Rightarrow \frac{7-3}{3} + 1 = 2.33 (\times)$$

Commonly, we have a padding of 0s so its  $8 \times 8$  in reality!

Convnet, shrinks slowly as fast shrinking doesn't work much well

② Input :  $32 \times 32 \times 3$

1D  $5 \times 5$  input with stride 1, pad 2 filter

$$\begin{aligned}\text{Output size: } & N - F/\text{str} + 1 \\ & = (32 + 2 \times 2 - 5)/1 + 1 \\ & = 32\end{aligned}$$

$$\begin{aligned}\#\text{ of parameters} &= 5 * 5 * 3 + 1 \\ &= 76 \text{ params}\end{aligned}$$

[760]

### ConvNet

→ accepts a vol of  $W_1 \times H_1 \times D_1$ ,

→ hyperparameters

- ↳ No of filters  $K$
- ↳ spatial Extent  $F$
- ↳ stride  $S$
- ↳ zero padding  $P$

→ produces a vol of  $W_2 \times H_2 \times D_2$  where;

$$W_2 = (W_1 - F + 2P) / S + 1$$

$$H_2 = (H_1 - F + 2P) / S + 1$$

$$D_2 = K$$

→ parameter sharing, so

$$\text{wts per filter} = F \times F \times D$$

$$\text{Total wts} = F \cdot F \cdot D \cdot K \text{ and } K \text{ biases}$$

### Common

$$K = 2^n$$

$$\rightarrow F = 3, S = 1, P = 1$$

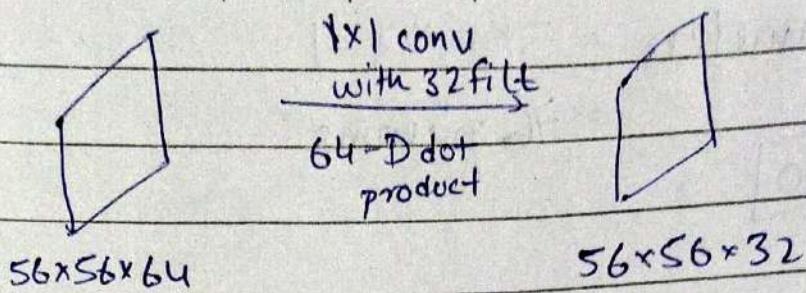
$$\rightarrow F = 5, S = 1, P = 2$$

$$\rightarrow F = 5, S = 2, P = ?$$

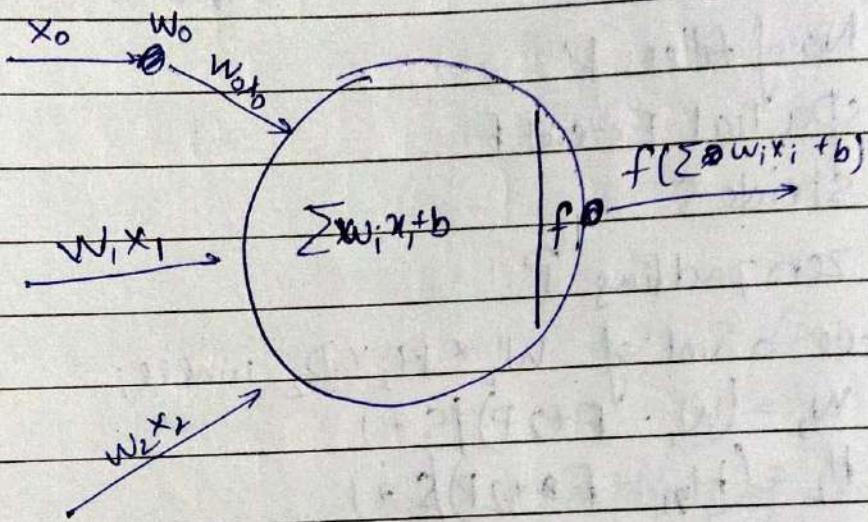
$$\rightarrow F = 1, S = 1, P = 1$$

NO.

Date



## Brain/Neuron View



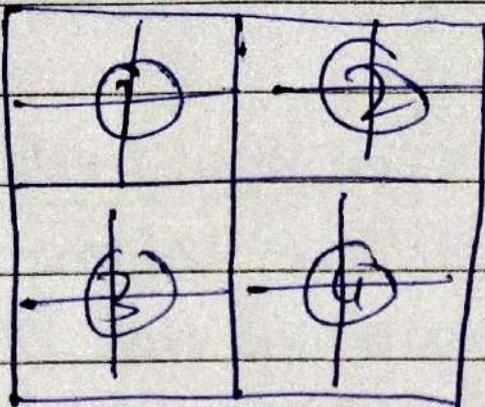
## Pooling Layer

- makes the representations smaller and more manageable
- operates over each activation map independently

NO.

Date

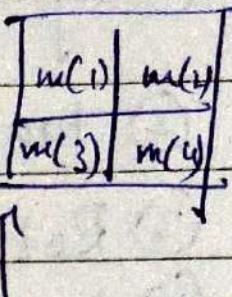
## Max pooling



maxpool

$2 \times 2$  filter

strides 2



max(max())

# "Lecture 6"

## Activation function

some examples

$$\textcircled{1} \text{ Sigmoid } \Rightarrow \sigma(x) = \frac{1}{1+e^{-x}}$$

$$\textcircled{2} \text{ tanh } x$$

$$\textcircled{3} \text{ ReLU } \Rightarrow \max(0, x)$$

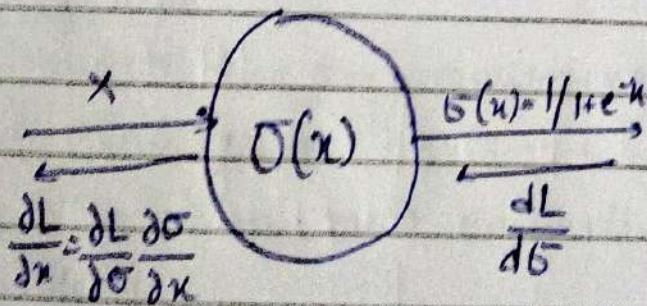
$$\textcircled{4} \text{ Leaky ReLU } \Rightarrow \max(0.1x, x)$$

$$\textcircled{5} \text{ Maxout } \Rightarrow \max(w_1^T x + b_1, w_2^T x + b_2)$$

$$\textcircled{6} \text{ ELU } \Rightarrow \begin{cases} x & x > 0 \\ \alpha(e^x - 1) & x \leq 0 \end{cases}$$

## Sigmoid

- Squashes numbers in range  $[0, 1]$



gradients on w?

↳ all positive or all negative

problems with sigmoid:

- ① saturated neurone kill gradient
- ② not zero centered
- ③  $\text{exp}()$  compute expensive.

$\tanh(x)$

- squashes no from  $[-1, 1]$
- zero centered
- but still kills gradients

ReLU

- doesn't saturate
- efficient
- converges faster than earlier
- more plausible than sigmoid.
- but not zero centered.

Leaky ReLU

↳  $\max(0.01x, n)$

→ will not die

NO.

Date

## PReLU

$$\hookrightarrow \max(\alpha x, x)$$

parameter.

## ELU

$$\hookrightarrow \begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

→ adds robustness to noise

## Maxout Neuron

→ generalizes ReLU and Leaky ReLU

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

problem: double the parameters

## In Practice:

- ① Use ReLU
- ② Use Leaky ReLU/maxout/ELU
- ③ use tanh (back)
- ④ never use sigmoid.

## Preprocessing of Data

- ① zero centre them
- ② normalise them. (std deviation)

Step 1:

- ① PCA the data
- ② Whiten the data

for images:

- subtract the mean image, or
- subtract per channel mean.

Weight initialization:

Q What happens  $W=0$ ?



First idea: small random nos

$$(W = 0.01 * \text{np.random.rand}(D, H))$$

This is okay for small networks

Xavier initialization

$$(W = \text{np.random.randn}(f_i, f_o) / \text{np.sqrt}(f_i))$$

NO.

Date

He et al

divide by 2

$$(w = np.random.rand(f_i, f_o) / np.sqrt(f_i/2))$$

Batch Normalization

To make each dimension unit gaussian apply:

$$\hat{x}_k = \frac{x^k - E(x^k)}{\sqrt{Var(x^k)}}$$

add after Fully connected or ConvNet layer.

first find mean and var for each dimension  
and then normalize

Then allow the network to do this.

$$y^k = \gamma^k \hat{x}^k + \beta^k$$

Batch Normalization:

i/p : values of  $x$  over mini-batch  $B = \{x_1, \dots, x_m\}$   
 parameters :  $\gamma, \beta$ .

o/p :  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

$$\mu_B \leftarrow \frac{1}{m} \sum x_i$$

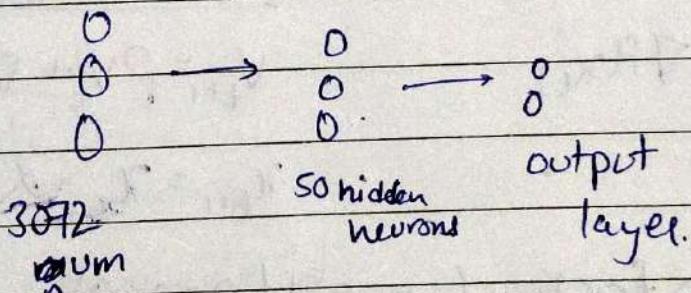
$$\sigma_B^2 \rightarrow \frac{1}{m} \sum (x_i - \mu_B)^2$$

$$x_i' \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$y_i \leftarrow \gamma x_i' + \beta = BN_{\gamma, \beta}(x_i)$$

② Choose an architecture

Let's choose 50.



double-check if the loss is reasonable.

✓ do cross validation

## "Lecture 7"

Loss function has "high cond" number.  
very slow progress along shallow dimension,  
jitter along steep direction.

Momentum (slow progress).

Q: What if loss fx<sup>"</sup> has a local minima or saddle point?

A: Grad descent gets stuck.

Gradients come in minibatches, so they can be noisy

SGD

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \alpha \nabla f(\boldsymbol{x}_t)$$

SGD + momentum

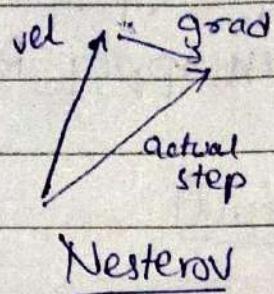
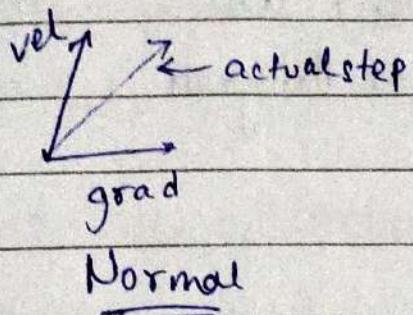
$$\boldsymbol{v}_{t+1} = \rho \boldsymbol{v}_t + \nabla f(\boldsymbol{x}_t)$$

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \alpha \boldsymbol{v}_{t+1}$$

$\rho \Rightarrow$  friction ( $\rho = 0.9 / \rho = 0.99$ )

- no noise
- better progress

## Nesterov Momentum



$$V_{t+1} = \rho V_t - \alpha \nabla f(x_t + \rho V_t)$$

$$x_{t+1} = x_t + V_{t+1}$$

change,

$$x_{t+1} = \rho x_t - \alpha \nabla f(\tilde{x}_t)$$

$$\begin{aligned} \textcircled{\$} \quad \tilde{x}_{t+1} &= x_t - \rho V_t + (1+\rho) V_{t+1} \\ &= \tilde{x}_t + V_{t+1} + \rho (V_{t+1} - V_t) \end{aligned}$$

## Adagrad

while True:

$$dx \rightsquigarrow$$

$$\text{gradsq} = dx^* dx$$

$$x_{\text{new}} = x - \text{learning rate} \times dx / \sqrt{\text{gradsq}} + 1e-7$$

## RMS prop

$$\text{gradsq} = \text{decay\_rate} * \text{gradsq} + (1 - \text{decay\_rate}) * dx^* dx$$

## Adam

$$\text{firstP} = \beta_1 * \text{firstP} + (1 - \beta_1) * dx$$

then apply RMSProp or Adagrad

## Adam (full)

### Adam half

Second moment

$$\text{first unbias} = \text{firstP} / (1 - \beta_1^t)$$

$$\text{second } " = \text{secondP} / (1 - \beta_2^t)$$

do adagrad / RMSprop on unbias

### I Use

$$\beta_1 = 0.9 \quad \beta_2 = 0.999$$

$$\alpha = 1e^{-3} / 1e^{-5}$$

## Learning Rate Decay

### ① Step decay

half per few epoch

### ② Exp decay

$$\alpha = \alpha_0 e^{-kt}$$

### ③ 't decay

$$\alpha = \alpha_0 / (1 + kt)$$

## First Order optimization

- ① Use gradient from linear approx.
- ② Step to minimize the approx.

## Second order optimization

- ① Use grad & hessian to form quadratic approx.
- ② step to the minima of approximation.

In practice,

→ Use ADAM

→ try out LBFGS

## Model Ensembles

- 1) Train multiple independent models
- 2) At test time, average their results
- 3) Use snapshots of single model during training.

## Dropout Regularization

- we drop random neurons to 0
- forces the network to have redundant rep.
- prevents co-adaption of features.

NO  
Date

Basically it is training a large ensemble of models.

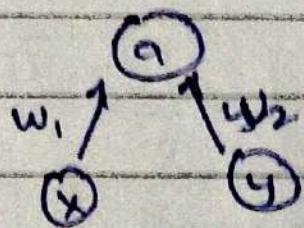
Each binary mask is one model!

Dropout makes our output random

$$y = f_w(x, \gamma) \quad \text{random mask}$$

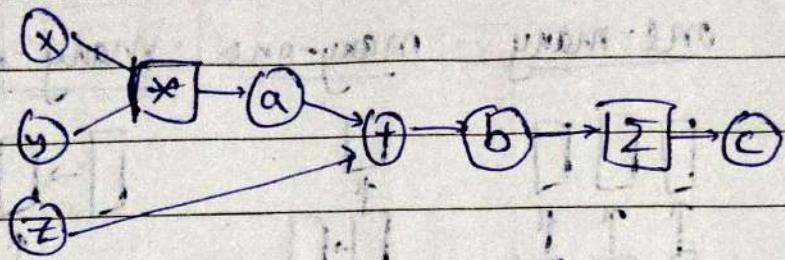
$$y = f(u) - E_2[f(u, z)] = \int p(z) f(u, z) dz$$

Consider a single neuron,



$$\begin{aligned} E[a] &= w_1x + w_2y \\ &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0w_1y) + \frac{1}{4}(0w_2y) \\ &= \frac{1}{2}(xw_1 + yw_2) \end{aligned}$$

Data Augmentation

"Lecture 8"Comp Graph Example

~~$N, D = 3, 4$~~

$n = \text{randn}(N, D)$

$y = \text{"}$

$z = \text{"}$

$a = n + y$

$b = a + z$

$c = \text{sum}(b)$

numpy $\rightarrow$  no GrPU

$\nabla c = 1 \cdot 0$

 $\rightarrow$  compute grad  
on own

$db = dc * \text{ones}(N, D)$

$da = db$

$dz = db$

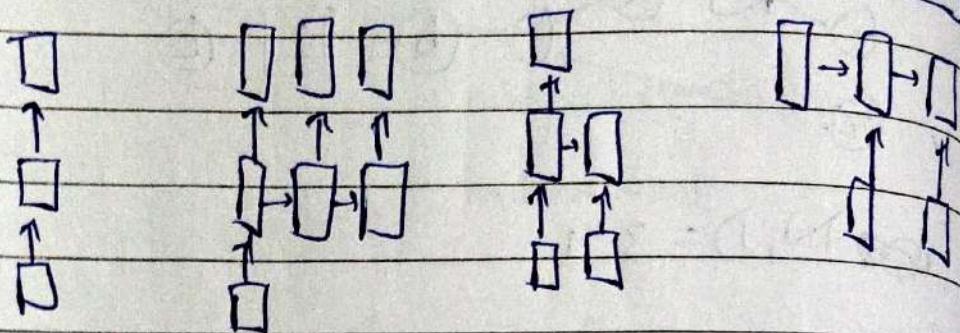
$dn = da * n$

$dy = da * n$

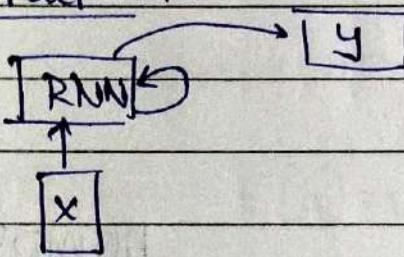
## "Lecture 10"

Process sequences of recurrent NNs

One-one    one-many    many-one    many-many



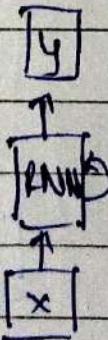
Recurrent NN



Recurrence relation:

$$h_t = f_w(h_{t-1}, x_t)$$

Vanilla RNN



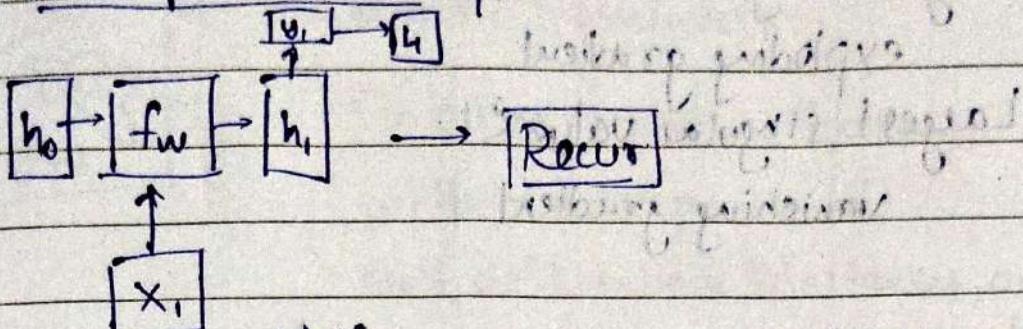
$$h_t = h_{t-1}$$

$$h_t = f_w(h_{t-1}, x_t)$$

$$h_t = \tanh(w_{nn} h_{t-1} + w_{nx} x_t)$$

$$y_t = W_{ny} h_t$$

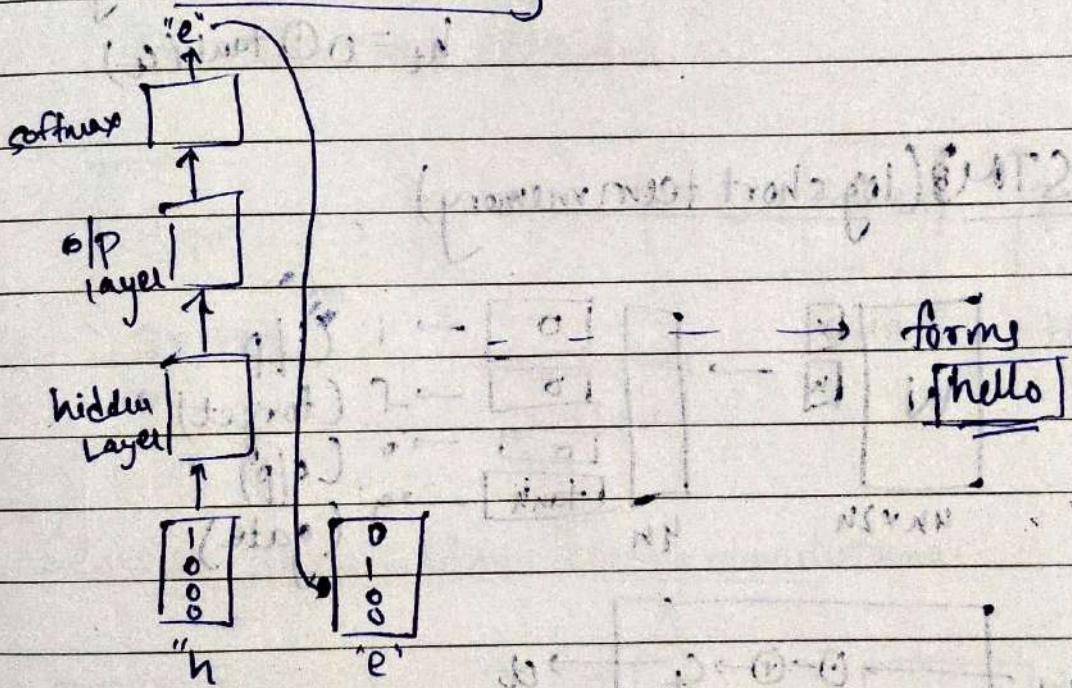
## RNN: Computational Graph



Use same  $W$  for every case for:  $x_i$ .

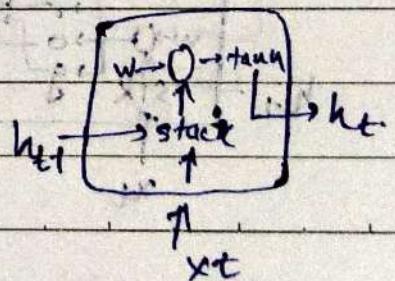
$$L = \sum L_i$$

## Character Level Lang Model



## Vanilla RNN gradient flow

$$\begin{aligned} h_t &= \tanh(W_{hh} h_{t-1} + W_{xh} x_t) \\ &= \tanh(W(h_{t-1}, x_t)) \end{aligned}$$



No

Date

Largest singular value > 1: exploding gradient

Largest singular value < 1: vanishing gradient

### Vanilla RNN

$$h_t = \tanh(w(h_{t-1}, x_t))$$

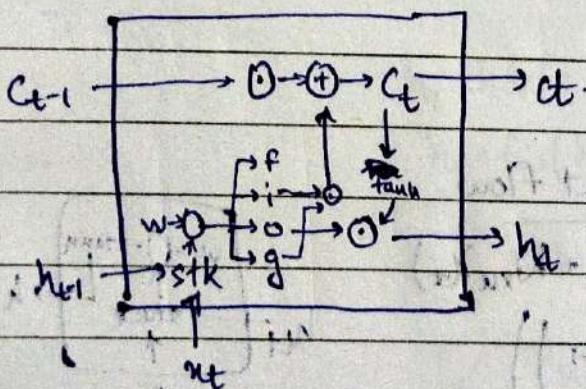
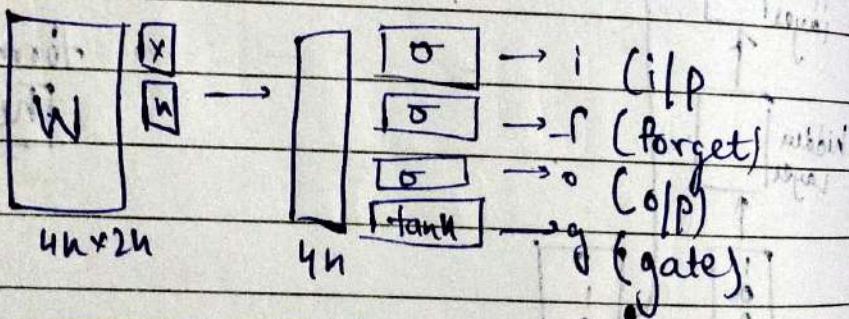
### LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} o \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} w \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

### LSTM (long short term memory)

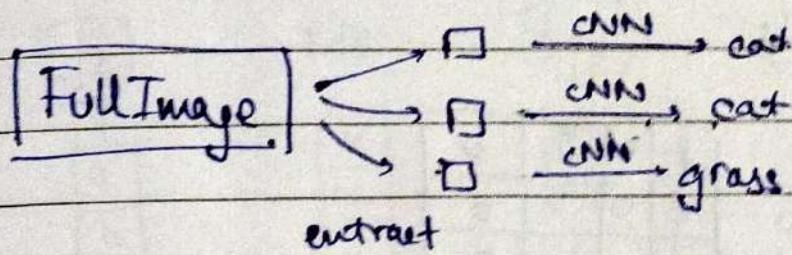


"Lecture 911"

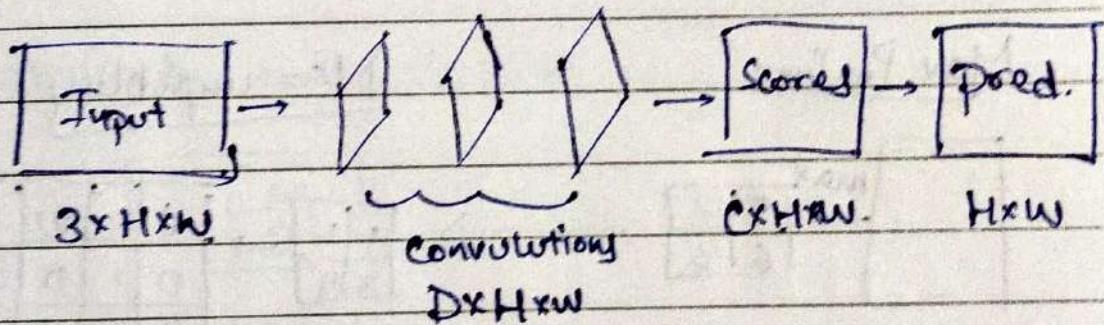
Semantic Segmentation

label each pixel in the image  
with a category email

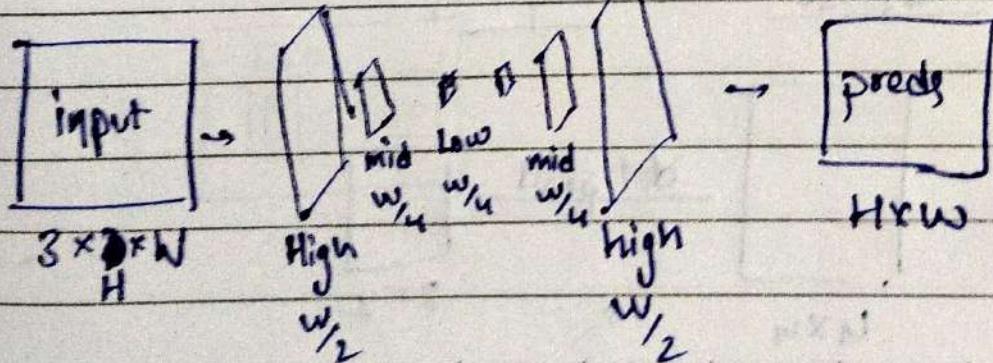
• don't differentiate instances, only  
care about pixels.



Problem: Not efficient



Semantic Segmentation: Fully convolutional



downsample then upsample

NO.

Date

## In-Network Upsampling: Unpooling

NN

1	2	1	1	2	2	...
3	4	1	1	2	2	...
		3	3	4	4	...
		3	3	4	4	...

Bed of Nails

1	2	0	0	0	0	0
3	4	0	0	0	0	0
		3	0	4	0	0
		0	0	0	0	0

Max Pooling

1	2	3	4
5	6	7	8
9	10	11	12

Max unpooling

1	2
3	4

0	0	0	0
0	1	0	6
0	0	0	0
3	0	0	4

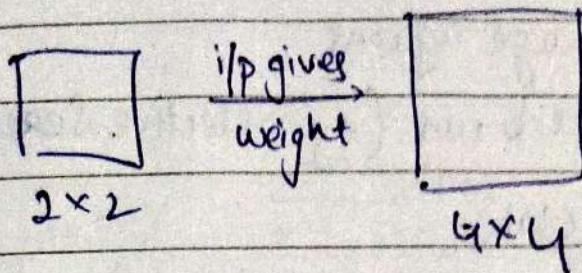
Transpose conv : sampling


4x4

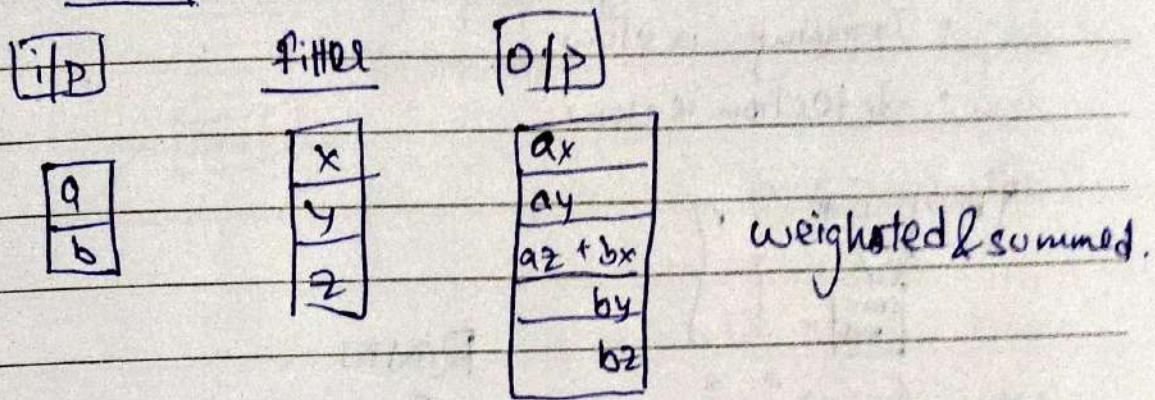
dot prod


2x2

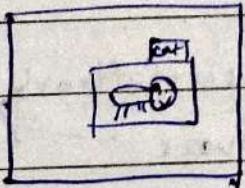
## upsampling



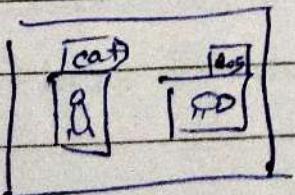
## 1-D example



## Classification and localization



## Object Detection

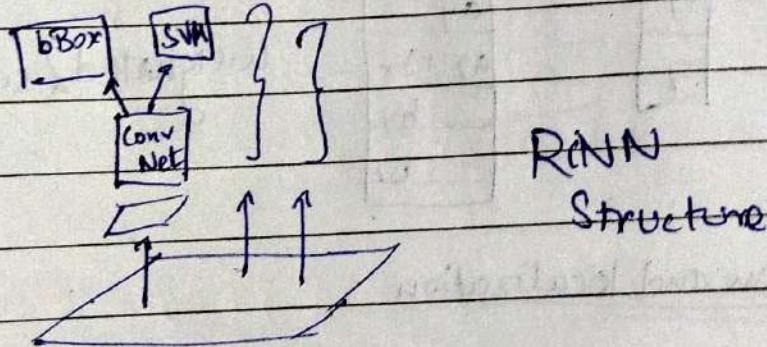


## Region Proposal

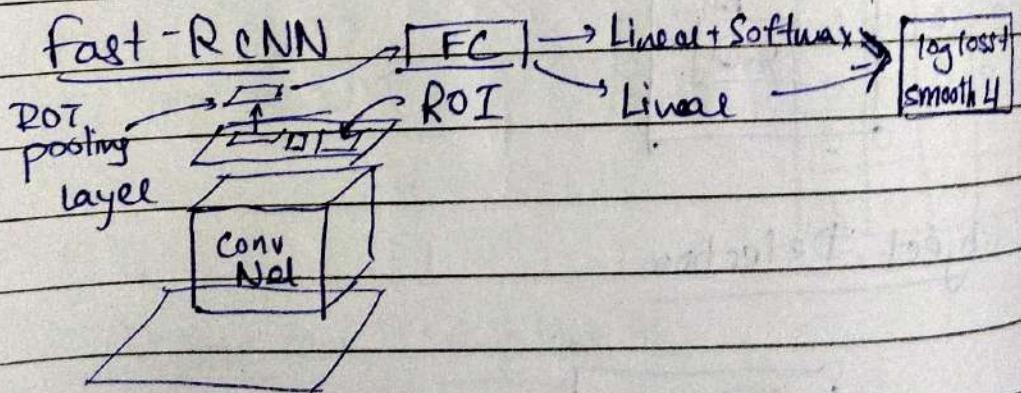
- find blobby image regions
- relatively faster to run (e.g. Selective Search)

## Problems w RCNN

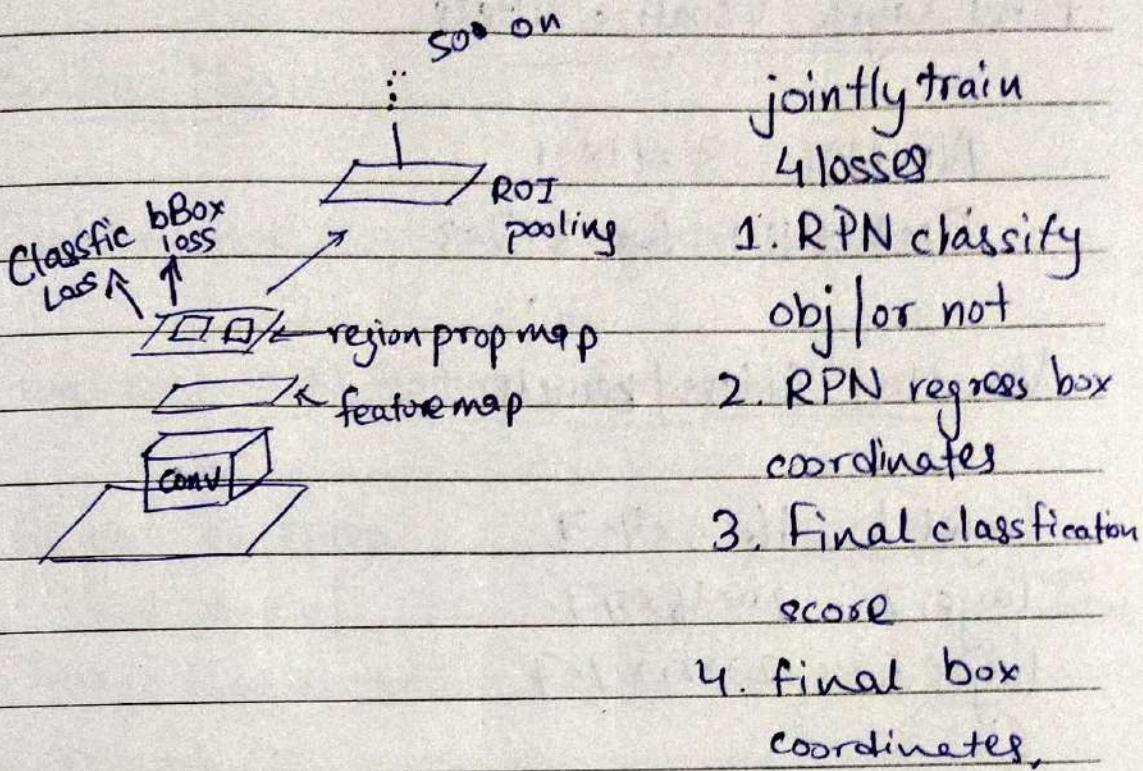
- Ad hoc training objectives
- Training is slow
- Detection is slow



RCNN  
Structure



## Faster RCNN



## Classification w/o proposals (YOLO)

- within each grid cell
  - ↳ makes a final box with 5 no ( $d_y, d_x, d_t, d_w, conf$ )
- predict scores for each class

## Instance segmentation

↳ Mask RCNN

NO. \_\_\_\_\_  
Date \_\_\_\_\_

## "Lecture 12"

### [ConvNet Structure]

First Layer : Visualize filters

AlexNet :  $3 \times 11 \times 11$

ResNet 101 :  $64 \times 3 \times 7 \times 3$

Visualize filters / kernels

Layer 1 :  $16 \times 3 \times 7 \times 7$

Layer 2 :  $20 \times 16 \times 7 \times 7$

Layer 3 :  $20 \times 20 \times 7 \times 7$

Last Layer : FC7 layer

↳ 4096-dim feature vector for an image.

↳ applies NN.

↳ dimensionality reduction.

visualize  $13 \times 13 \times 128$  map as  $128 \text{ } 13 \times 13$  grayscale images

Maximally Activating patches

## Occlusion Experiment

## Saliency Maps

Intermediate features via guided backprop

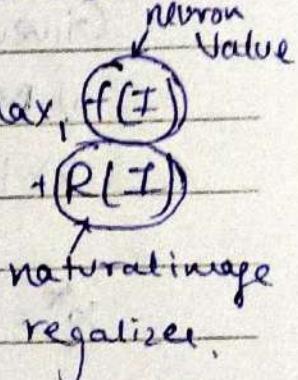
Visualizing CNN features : Gradient Ascent

Guided Backprop:

find part of  
image that neuron  
responds to

Use grad ascent.

$$I^* = \text{argmax}_I f(I)$$



1. Initialize your image to zeros

Repeat

2. Forward into compute current score

3. Backprop to get grad

4. make small image

$$\text{score for class } c \Rightarrow \text{argmax}_I S_c(I) - \lambda \|I\|_2^2$$

$$\text{Simple regularizer} \Rightarrow \lambda \|I\|_2^2$$

Better regularizer  $\Rightarrow$  Gaussian  $\rightarrow$  clip pixel to 0  
 $\hookrightarrow$  clips small grad  
to 0.

NO.  
Date

## Adversarial Examples

- (1) Start from arbitrary image.
- (2) Pick arbitrary class
- (3) Modify image to maximise class
- (4) Repeat till they are fooled

## Feature Inversion

Given a CNN feature find new image:-

- 1) Matches given feature vector
- 2) looks natural

$$x^* = \underset{x \in \mathbb{R}}{\operatorname{argmin}} L(\phi(x), \phi_0) + \lambda R(x)$$

$$\text{where, } L(\phi(x), \phi_0) = \|\phi(x) - \phi_0\|^2$$

$$R(x) = \sum_{i,j} \left( (x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)$$

. Texture Synthesis : Nearest Neighbour.

NO.

Date

- Neural Texture Synthesis: Gram Matrix.
- Neural Style Transfer.
- Fast Style Transfer

## Lecture 13"

### Unsupervised Learning

Data:  $x$  (no label)

Goal: Learn some hidden structure of the data

Eg: Clustering, dimension reduction, feature learning, density estimation.

### Generative Models

from training data, generate new sample

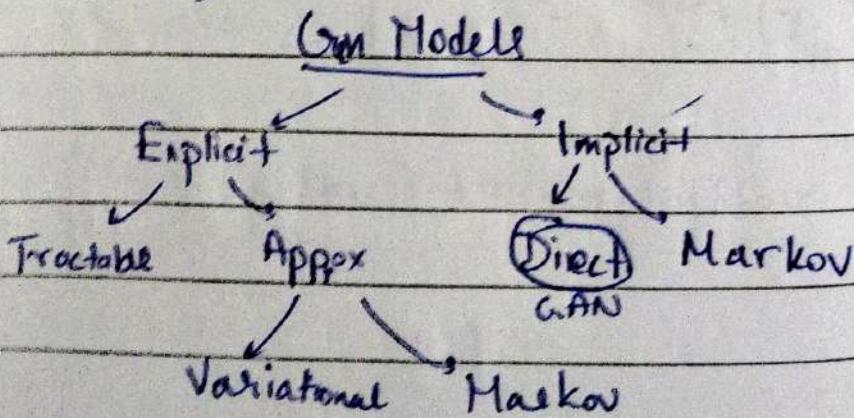
$$P_{\text{data}}(x) \longrightarrow P_{\text{model}}(x)$$

Addresses density estimation, a core problem.

#### Flavors:

- i) explicit density estimation
- ii) Implicit density estimation

### Taxonomy



## Pixel RNN & CNN

fully visible belief network:

$$P(x) = \prod_{i=1}^n P(x_i | x_1, x_2, \dots, x_{i-1})$$

### Pixel RNN

Generate image pixels starting from corner.

use RNN(LSTM)

### Pixel CNN

Still generate image pixel starting from corner.

↳ train using max likelihood.

Cone → Sequential gen → slow

Pros → can compute likelihood / expected likelihood  
→ good sample.

## Variable Autoencoders (VAE)

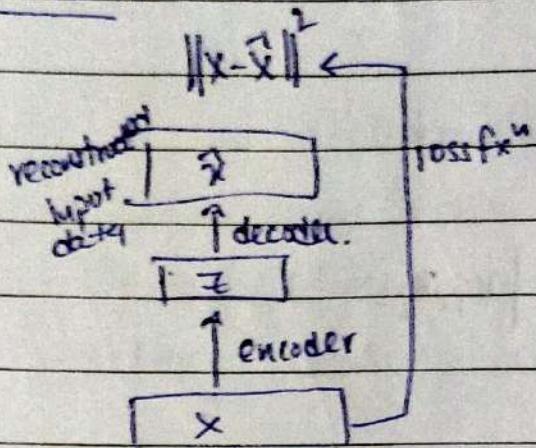
$$P_\theta(x) = \int P_\theta(z) P_\theta(x|z) dz$$

↳ can't optimize directly.

NO.

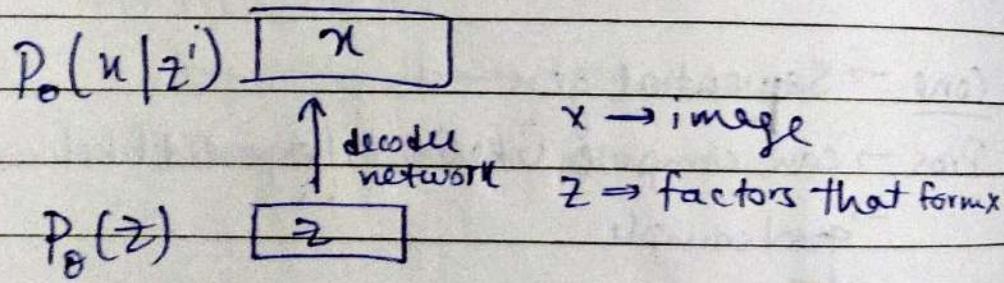
Date

## Autoencoder



Encoder can be used as supervised model

## Variational autoencoder



How to represent this model?

- Choose  $p(z)$  to be simple
- ~~$p(x|z)$~~   $p(x|z)$  is complex  $\Rightarrow$  neural network

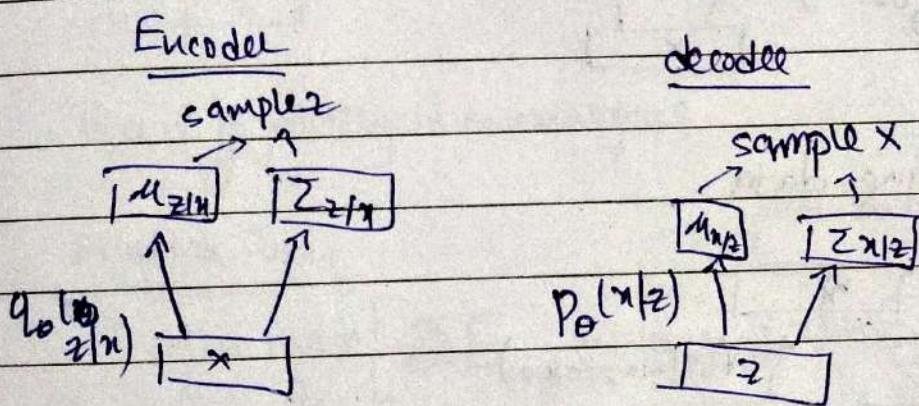
How to train this model?

Learn model parameters to maximise likelihood of training data.

$$P_\theta(x) = \int P_\theta(z) P_\theta(x|z) dz$$

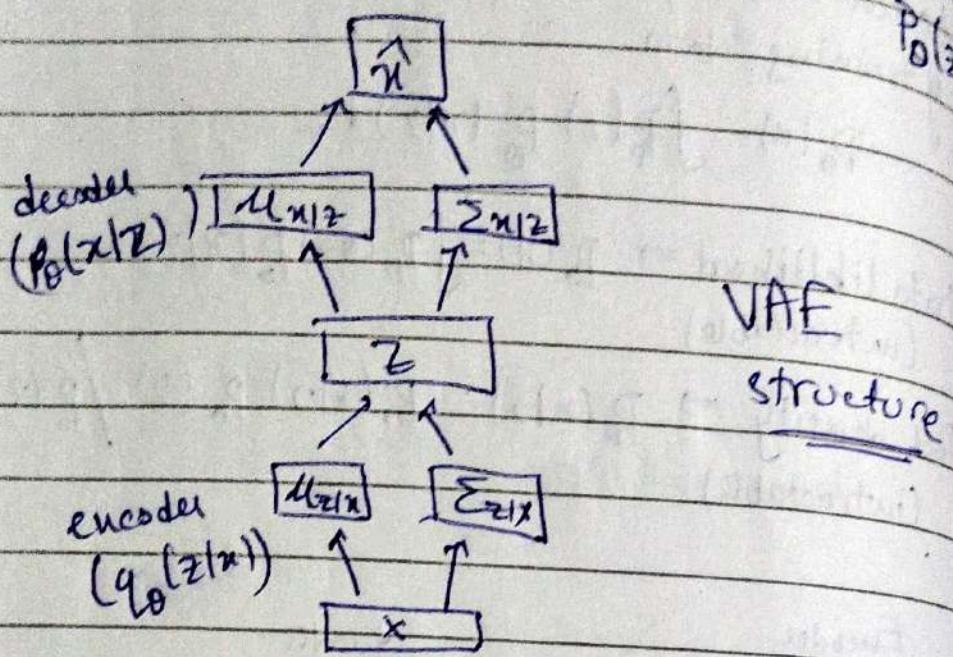
$$\text{Data likelihood} \Rightarrow P_\theta(x) = \int P_\theta(z) P_\theta(x|z) dz \quad (\text{intractable})$$

$$\text{Post. density} \Rightarrow P_\theta(x|z) = P_\theta(x|z) P_\theta(z) / P_\theta(z) \quad (\text{intractable})$$

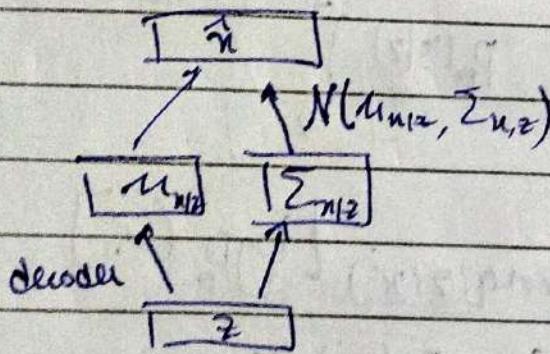


$$\begin{aligned}
 \log P_\theta(x^i) &= \mathbb{E}_{z \sim q(z|x^i)} [\log P_\theta(x^i|z)] \\
 &= \mathbb{E}_z \left[ \log \frac{P_\theta(x^i|z) P_\theta(z)}{P_\theta(z|x^i)} D_{KL}(q_\phi(z|x^i) || P_\theta(z)) \right] \\
 &= \mathbb{E}_z [\log P_\theta(x^i|z)] - \mathbb{E}_z \left[ \log \frac{q_\phi(z|x^i)}{P_\theta(z)} \right] \\
 &\quad + \mathbb{E}_z \left[ \log \frac{q_\phi(z|x^i)}{P_\theta(z|x^i)} \right] \\
 &\quad \rightarrow D_{KL}(q_\phi(z|x^i) || P_\theta(z|x^i))
 \end{aligned}$$

$$L(x^i; \theta, \phi) = E_z \left[ \log_p(p_{\theta}(x^i | z)) - D_{KL}(q_{\phi}(z|x^i) \| p(z) \right]$$



Generating data



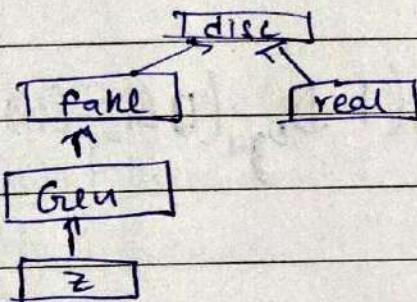
Vary  $z_1/z_2$  for degree of smile.

Cons: maximised lower bound of likelihood  
samples blurry

## GANS

→ Training (2 player game)

- Generator: try to fool disc by making real img
- discriminator: check if img is real or no.



## GAN structure

Training jointly in minmax game.

## MinMax func

$$\min_{G_g} \max_{D_d} \left[ \underbrace{E_{x \sim p_{\text{data}}} \log D_d(x)}_{\text{real Data}} + E_{z \sim p_g(z)} \log \underbrace{(1 - D_d(G_g(z)))}_{\text{generated Data}} \right]$$

$\rightarrow \theta_d$  wants to maximize obj s.t.  $D(x) \approx 1$   
 $\& D(g_i(x)) \approx 0$

$\rightarrow \Theta_g$  wants minimize obj s.t  $D(G(z)) \sim [fooling]$

## Alternate b/w

### ① Grad asc on disc

$$\max_{D_\theta} [E_{\text{empirical}} \log D_{\theta}(x) + E_{z \sim p(z)} \log(1 - D_{\theta}(G_{\phi}(z)))]$$

### ② Grad descent on gen

$$\min_{\phi} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta}(G_{\phi}(z)))$$

↳ updated obj

$$\min_{\theta} E_{z \sim p(z)} \log D_{\theta}(G_{\phi}(z))$$

### Algorithm

for no of training do

for k steps do

- Sample minibatch of m sample from  $p_g(z)$

$p_g(z)$

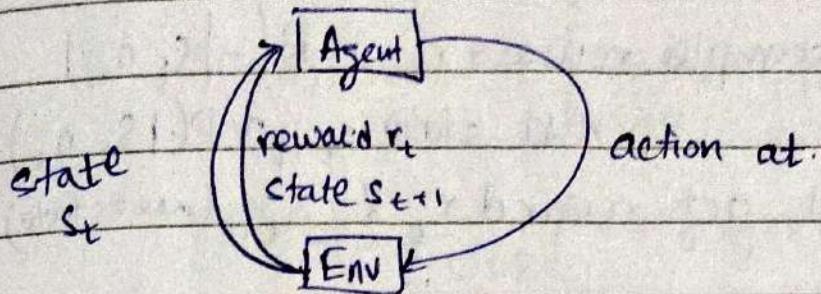
- Sample minibatch of m from  $p_{\text{data}}(x)$

s stochastic descent

- Sample m noise samples from  $p_g(z)$

- Update gen by stochastic descent

end for.

"Lecture 14"Reinforcement LearningEg

- (1) Cart Pole problem
- (2) Robot Locomotion
- (3) Atari Games
- (4) Go

How to write RL in math?

Ans Markov Decision Process

defined by,  $(S, A, R, P, \gamma)$

$S \rightarrow$  set of possible states

$A \rightarrow$  set of " actions "

$R \rightarrow$  dict of reward to  $(S, A)$  pair

$P \rightarrow$  dict to next state of  $(S, A)$  pair

$\gamma \rightarrow$  discount factor.

At  $t=0$ , env samples  $s_0 \sim p(s_0)$

Then, for  $t=0$  until done:

- Agent selects  $a_t$
- Env samples reward  $r_t \sim R(\cdot | s_t, a_t)$
- " " next state  $s_{t+1} \sim P(\cdot | s_t, a_t)$
- Agent gets reward  $r_t$  &  $s_{t+1}$  (next state)

A policy  $\pi$  is a func from  $S \rightarrow A$  that tells the actions to take

Obj: find  $\pi^*$  that max cum. disc. reward

$$\sum_{t \geq 0} \gamma^t r_t$$

Formally,

$$\pi^* = \operatorname{argmax}_{\pi} \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | \pi \right] \text{ where, } s_0 \sim p(s) \\ a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)$$

Value fx<sup>n</sup>

↳ expected cum. reward from policy from  $s$ .

$$V^\pi(s) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi \right]$$

Q value fx<sup>n</sup>

↳ expected cum. ~~rewards~~ from a <sup>1<sup>st</sup> state</sup>  $s$ .

$$Q^\pi(s, a) = \mathbb{E} \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

Bellman Eq<sup>n</sup>

Optimal Q value fx<sup>n</sup>  $\rightarrow Q^*$

$$Q^*(s, a) = \max_{\pi} E \left[ \sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

$$Q^*(s, a) = E_{s' \sim \mathcal{E}} [r + \gamma \max Q^*(s', a') | s, a]$$

optimal policy  $\pi^*$  corresponds to taking best action acc to  $Q^*$ .

Value iteration algorithm

$$Q_{i+1}(s, a) = E \left[ r + \gamma \max_{a'} Q_i(s', a') \right] | s, a$$

$Q_i \rightarrow Q_a^*$  as  $i \rightarrow \infty$ .

Q-learning

$\rightarrow$  use func approximator for a-v func.

$$Q(s, a; \theta) \approx Q(s, a)$$

$\rightarrow$  deep q learning (if we use DNN)

We need to find  $Q^*$  satisfying Bellman eq<sup>n</sup>.

## Forward Pass

Loss fn:  $L_i(\theta_i) = E_{s,a \sim p(\cdot)} [y_i - Q(s, a; \theta_i)]^2$

where,  $y_i = E_{s \in E} [r + \gamma \max_{a'} Q(s', a'; \theta_{i+1})]$

## Backward Pass

$$\nabla_{\theta_i} L_i(\theta_i)$$

$$= E_{\theta(s, a \sim p(\cdot); s \in E} [s + \gamma \max_{a'} Q(s', a'; \theta_{i+1}) - Q(s, a; \theta_i)] \nabla_{\theta_i} Q(s, a; \theta_i)$$

## Net architecture (Atari)

$$Q(s, a; \theta)$$

[FC4]

[FC256]

[32 4x4 S2]

[16 3x3 S4]

$$S_t \rightarrow 84 \times 4 \times 4$$

Train the Q network by exp relay

- continually update memory table of transitions  $(s_t, a_t, r_t, s_{t+1})$  as eps all played
- train Q network on minibatches of transitions instead of consecutive samples. (random)

### Policy gradients

$$\Pi = \{ \pi_\theta, \theta \in \mathbb{R}^m \}$$

$$\text{Let } J(\theta) = E \left[ \sum_{t \geq 0} r^t r_t | \pi_\theta \right]$$

$$\theta^* = \underset{\theta}{\operatorname{argmax}} J(\theta) = ?$$

Reinforce Algo

$$J(\theta) = E_{r \sim p(r; \theta)} [r(T)]$$

$$= \int_T r(T) p(T; \theta) dT$$

where  $r(T)$  is reward of trajectory  $T = (s_0, a_0, r_0, s_1, \dots)$

$$\nabla_\theta J(\theta) = \int_T r(T) \nabla_\theta p(T; \theta) dT$$

$$\Rightarrow \nabla_\theta p(T; \theta) = p(T; \theta) \frac{\nabla p(T; \theta)}{p(T; \theta)} = p(T; \theta) \nabla \log(T; \theta)$$

NO.

Date

$$\text{So, } \nabla_{\theta} J(\theta) = E_{T \sim p(T; \theta)} [r(T) \nabla_{\theta} \log p(T; \theta)]$$

$$\text{we know, } p(T; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$$

$$\log p(T; \theta) = \sum_{t \geq 0} \log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t)$$

$$\nabla_{\theta} J(\theta) \propto \sum_{t \geq 0} r(t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

→ if  $r(T)$  high, push up prob

→ " " low, pull down "

### Variance Reduction

$$\textcircled{1} \quad \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left( \sum_{t' \geq t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

$$\textcircled{2} \quad \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left( \sum_{t' \geq t} \gamma^{t-t'} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

### ③ Baseline

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left( \sum_{t' \geq t} \gamma^{t-t'} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Actor-Critic Algorithm

$$\Lambda^{\pi}(s, a) = Q^{\pi}(s, a) - V^{\pi}(s)$$

for  $i=1, 2, \dots$  do

$$\Delta \theta \leftarrow \theta$$

for  $i=1, 2, \dots m$  do

for  $t=1, 2, \dots T$  do

$$\Lambda_t = \sum_{t' \geq t} \gamma^{t-t'} r_{t'} - V_\phi(s_{t'}^i)$$

$$\Delta \theta \leftarrow \theta + \Lambda_t \nabla_\theta \log(a_t^i | s_t^i)$$

$$\Delta \phi \leftarrow \sum_i \sum_t \nabla_\phi \|\Lambda_t^i\|^2$$

$$\theta \leftarrow \alpha \Delta \theta$$

$$\phi \leftarrow \beta \Delta \phi$$

end for

Recurrent Attention Model (RNN)

↳ uses Reinforce Algo