

Lecture - 3

Deep learning, Intuition

1. Problem with one-hot encoding?

If there are multiple objects in image, it will not work. We will need multi-hot encoding then.

The way you choose label is important & you should choose before starting anything.

2. Encoding:

Shower layer see pixel level information.
Deeper layer see more complex, high level features.
Thus every layer encodes.

3. Given an image classify as taken "during the day" (0) or "during the night" (1).

1. Data - How many images? Depends upon the complexity of task to learn. If need to classify more complex like dawn images, indoor-night images. You need high

Around 10,000 sufficient.

split?
↓
80% 20%

Bias?

Correct balance b/w classes.

2. Input: Image.

Resolution?

→ In terms of computation, higher pixels, higher parameters.

~~256 x 256~~ with human accuracy.
fine (64, 64, 3) and also in terms of
for human. Complexity of task

3. Output:

$$y = 0 \text{ or } y = 1$$

lost activation?

↳ sigmoid.

4. Architecture - shallow network should do the job pretty well.

5. Loss:
$$L = -[y \log(\hat{y}) + (1-y) \log(1-\hat{y})]$$

4. Face Verification: A school wants to use face verification for its facilities. i.e. Bym, (for validating the ID)

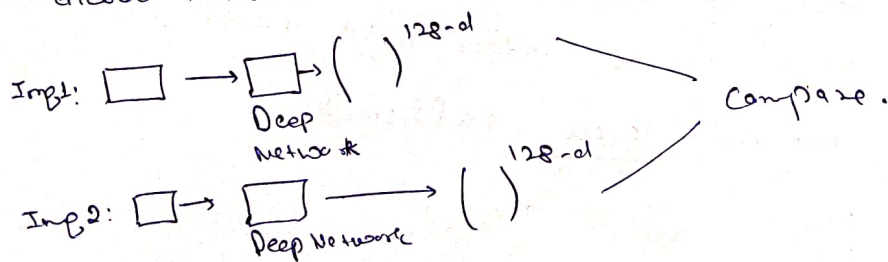
i) Data? - picture of every student labelled with their name.

ii) Input - picture of person standing
Resolution - for (64,64) it's harder to decide the deep features like nose, eyes etc.
So try over (412, 412, 3)

iii) Output: $y = 1$ it's you
or $y = 0$ it's not you.

iv) Architecture:

encode information about a picture in a vector



we gather all student faces encoding in a database
Given a new image, we compute its distance with the encoding vector.

v) Loss? Training? - first we need more data so that our model understands how to encode. Use public face datasets.

what we really want?

Similar Faces

↓

Similar encoding

Different Faces

↓

Different Encoding

Thus: Anchor

Positive

Negative.

Loss function:

$$L = \| \text{Enc}(A) - \text{Enc}(P) \|_2^2 - \| \text{Enc}(A) - \text{Enc}(N) \|_2^2$$

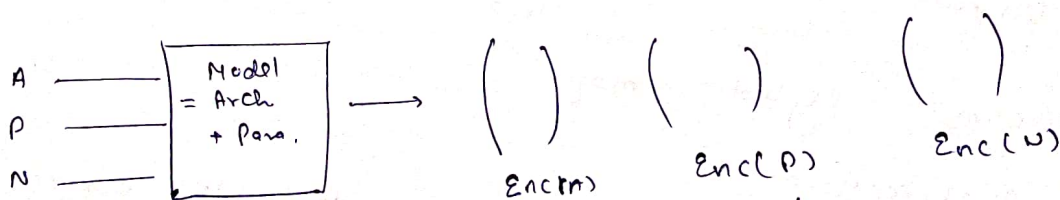
(A)

$$L = \| \text{Enc}(A) - \text{Enc}(N) \|_2^2 - \| \text{Enc}(A) - \text{Enc}(P) \|_2^2$$

(B)

$$L = \| \text{Enc}(P) - \text{Enc}(N) \|_2^2 - \| \text{Enc}(P) - \text{Enc}(A) \|_2^2$$

(C)



Gradients.

Loss
↓

$$L = \| \text{Enc}(A) - \text{Enc}(P) \|_2^2 - \| \text{Enc}(A) - \text{Enc}(N) \|_2^2 + \alpha$$

margin

to push the
networks to learn
more.

5. Face Recognition:

Goal: A school wants to use face identification to recognise students.

k-Nearest neighbours

on encodings.

Goal: You want to use face clustering to group pictures of same people on your smartphone.

k-Means Algorithm

on encodings ~~from~~

6. Art Generation: Neural style transfer.

Goal: Given a picture, make it beautiful.

1. Data: Let's say we have a data

2. Input: content + style Image

3. Output: styled Image.

4. Architecture:

We want a model that understands images very well - we load an existing model trained on imagenet.



When this image forward propagate, we can get information about its content & ~~by~~ style by inspecting the layers.

5. Loss?

$$L = |Content_c - Content_s|_2^2 + |Style_s - Style_c|_2^2$$

(A)

$$L = ||Content_c - Content_s||_2^2 + ||Style_s - Style_c||_2^2$$

(B)

We are not learning parameters by minimizing L .
We are learning an image.

7. Trigger word detection.

Goal: Given a 10sec audio speech, detect the word active.

Data: A bunch of 103 clips

Distribution?

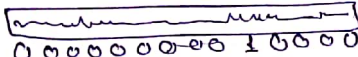
↓
all types of noise, accents

Input: $x = A$ 10sec audio clip

Resolution? (sample rate).

↓
asks researcher.

Output: $y = 0$ or 1

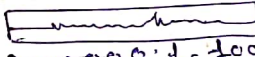
Output: - 

Pomarezo
↓
Afternoon in Italian

Should word 1 come before or after?

- After. You need to hear to label.

Naive: Add several ones after first 1.



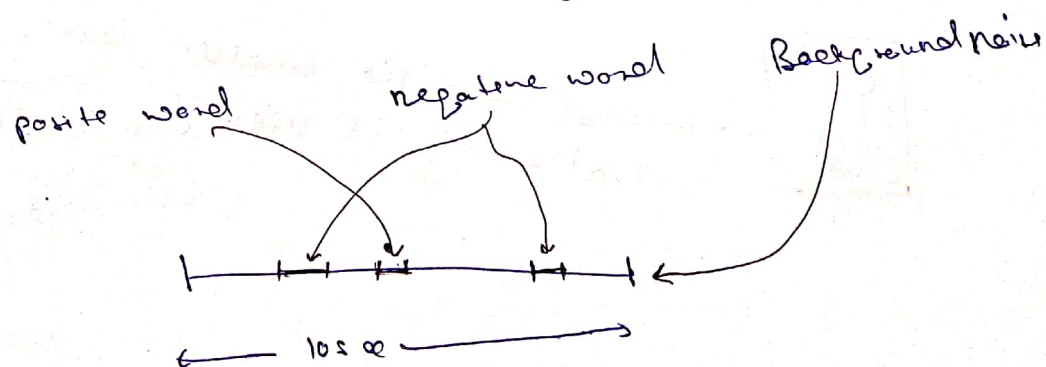
Last activation: Sigmoid (sequential)

Architecture: Sounds like it should be RNN.

Loss: $L = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$

What is critical to the success of this project?

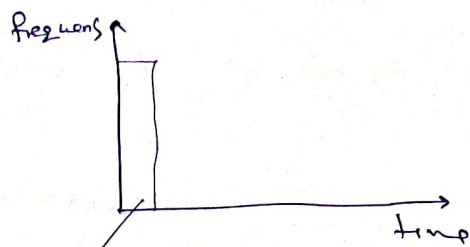
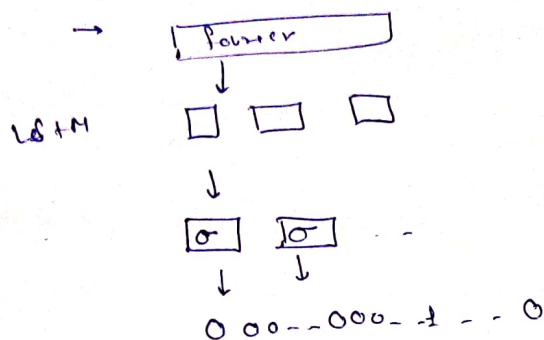
i) Strategic data collection / labelling process.



By this way preprogrammatic labelling & generation.

ii) Architecture Search & Hyperparameter tuning.

→ started with Fourier transform.



value of all the amplitudes of this frequency for every time step.

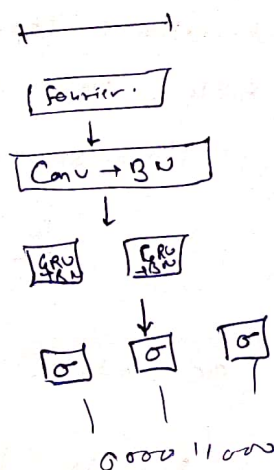
This is spectrum.

↓

But it doesn't work.

So an expert forced guided & finally.

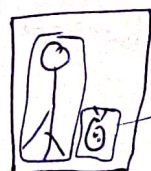
Find an expert for your problem



The most beautiful loss functions of 2015.

↓

YOLO. Bounding boxes.



penalize more the smaller boxes.
That's why it used $\sqrt{w_i}$ & $\sqrt{h_i}$