

CV2 Notes

Vision → To know what is where by looking.

Application: {CV → Art of perception by H/cs.}

- ① Computational photography (Increasing exposure and facial recognition).
- ② Biometrics (used as facial recognition for security purposes)
- ③ Optical character recognition (for recogn. chars on car numbers & hand-writing recognition).
- ④ Health (used in radiology and Cleo-II and III for doctors/ pathologists).
- ⑤ Gaming (recognizing hand gestures and postures along with location).
- ⑥ Shopping (Amazon Go).
- ⑦ Security and Tracking (Recognizing images and use it to track miscreants / law breakers).
 - Where do you draw the line?
- ⑧ Special Effects → in movies, thru CV, images/videos are animated & they have to track fine details on the face to track it to animations.
- ⑨ Insight → "Track stock prices of walmart dep. on no. of cars in front of stores".
CV could help us get data / insights into the data.
- ⑩ Self-driving cars → CV detection, human-figure & sign detection and make changes (RL).
- ⑪ Augmented Reality → "digitizing correct, location correct, shadow also considered".
(into real-time, look as real objects).
- ⑫ Space Exploration → "Mars Rover".
- ⑬ Why is CV so hard?
 - 1. firstly, computer sees image as numbers (as a matrix, 3D matrix of numbers) → Rep. digital images.

- ② Illumination → As lighting is different, the image changes based on illumination when a photo is taken. (a shadow could even change the image by a lot).
- ③ Occlusion → we live in a 3D world but try to look at the world through 2D. (having overlaps of background and foreground). → in a 3rd axis which can't be shown correctly / detected.
- ④ Class variation → "Diff kind of chairs", it has so many variations, so to categorize an object, it's really hard problem.
- ⑤ Clutter and camouflage → we can see camouflage but comp. (Evolution)
Can't see camouflage easily as the intensity graph formed is really hard to separate the object and background. (map) (or find a needle in a haystack).
- ⑥ Color → It depends on diff. kinds of intensities in an image.
(Colors seem different due to that).
⇒ [color of light * color of the object] might seem different.
- ⑦ Motion → "Blue occurs in an image" → due to which processing becomes diff. for an object. (motion detection).
- ⑧ Ill-posed problem → "Necker cube". (diff. 3D shapes could project the same 2D image).
→ how to solve this inverse transf. problem from 2D to 3D. (Accidental view point).
(diff. angles used for perception).
↳ mostly used for art!

Cambrian Explosion. → "sudden evolution of vision" (why?)
↳ sets off an evolutionary arms race where animals either evolved / died!

{ Evolution of biological eye } → Aperture → Exposure → Intensity control (camera)

→ About $\frac{2}{3}$ rd of the brain is devoted to visual processing. (2)
Loops are important! → (for complex processing).

Lecture 2

Image Denoising ⇒

film grain ≈ {black & pepper noise on an image.}

Ways:-

(i) Average many photos (to remove noise) → Problems:-

(a) what if you have only one?

(b) frame shouldn't be moved!

Image as functions ⇒ $F[x, y]$ → n_{x,y} coord → {returns the pixel values in the matrix of the image.}

Type of filters that could be used →

(i) Moving Avg. (Take an avg. of the neighborhood.)

(filter sizes shouldn't be very big) → otherwise you lose a lot of data!
→ Kernel
(This will blur the image).

Linear filter properties:

$$\left\{ \begin{array}{l} (1) \text{ filter(im, } f_1 + f_2) = \text{filter(im, } f_1) + \text{filter(im, } f_2) \\ (2) \text{ } c * \text{filter(im, } f_1) = \text{filter(im, } c * f_1) \end{array} \right.$$

(Avg. kernelization)

$$\Rightarrow F[x, y] * \frac{1}{9} \times G_i[n, y] = \text{Output Image.}$$

↑ convolution
↑ filter
↑ image

1	1	1
1	2	1
1	1	1

Defined

$$(f * g)[x, y] = \sum_{i, j} f[x-i, y-j] g[i, j]$$

rule of flip the Kernel and then apply it to make mathematics easy.

(for symmetric filters, this won't matter much!)

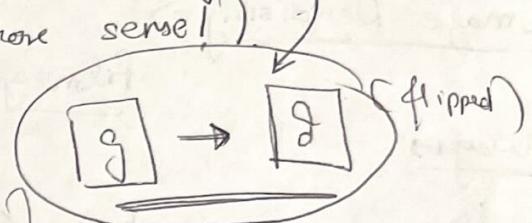
But for asymmetric filters → this makes much more sense!

→ But we generally do!

{ dec. indices in F and inc. in filter
ep. to flipping the kernel }

flip LR, UD

(left-right, up-down)



Translation filter

because we flip the filter and move in the image from bottom right, the image moves accordingly & an extra white/black would be added to the opp. edge.

0	0	0
0	0	1
0	0	0

more imp. to the right & a white col is added to LHS.

(bcz. of the flip)

$\frac{1}{9}$	1	1	1
1	1	1	
1	1	1	

→ "Blue filter".

(As you increase size of the filter, you lose more data!)

→ "No linear filter" can be used to rotate/transpose the image.

If we put a 0 in the averaging filter, we see it as an intensity loss (info. loss) and because of which the blue looks a little intensified.

Sharpening filter

0	0	0
0	2	0
0	0	0

$$- \frac{1}{9}$$

1	1	1
1	1	1
1	1	1

Ans.

(3)

Sharpening

$f_{\text{Image}} - f_{\text{Blurred}} = \text{Details lost (edges)}$

$f_{\text{Image}} + \text{Details lost} \Rightarrow \text{Sharpened Image}$

Essentially, edges are intensified & other details remain the same and noise will also be intensified. (In case there is noise).

Explanation

$$\left\{ \begin{array}{l} f_{\text{Image}} = D_2 + \text{Noise} \\ \text{Blurred} = D_1 \end{array} \right.$$

$$f_{\text{Image}} - \text{Blurred} = D_2 - D_1 + \text{Noise}$$

↓
Details lost

As noise is intensified, the edges look sharper to us (another way).

Convolution Properties:-

Commutative :-

$$F * H = H * F$$

Distributive :-

$$(F * G) + (H * G) = (F + H) * G$$

Associative :-

$$(F * H) * G = F * (H * G)$$

Scale

$$\text{filter}(A * f) = A * \text{filter}(f)$$

Ans

Shift Invariance :-

$$\text{filter}(\text{shift}(f)) = \text{shift}(\text{filter}(f))$$

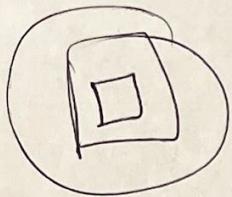
Cross-correlation

but not nice property

$$(f * g)[n, y] = \sum_{i,j} f[n+i, y+j] g[i, j]$$

Ans.

Box filter :-



1	1	1
1	•	1
1	1	1

{ Kind of shaped images are
averaged together. }

→ middle pixel not being considered.

(losing info. in the middle)

& blurring the image.

+ you will see multiple details.
(multiple glasses, multiple left eye, etc.)

Ans.

Instead we use gaussian filter

(this solves the problem with box filter)

(symmetric filter)

$$G_\sigma = \frac{1}{2\pi\sigma^2} e^{-\frac{(n^2+y^2)}{2\sigma^2}}$$



mean = 0

Constant factor at front makes volume sum to unity).

($\sigma \rightarrow$ determines the extent of smoothing)

Complexity?

↳ of filtering an $n \times n$ image with an $m \times m$ kernel?

$$\boxed{O(n^2m^2)}$$

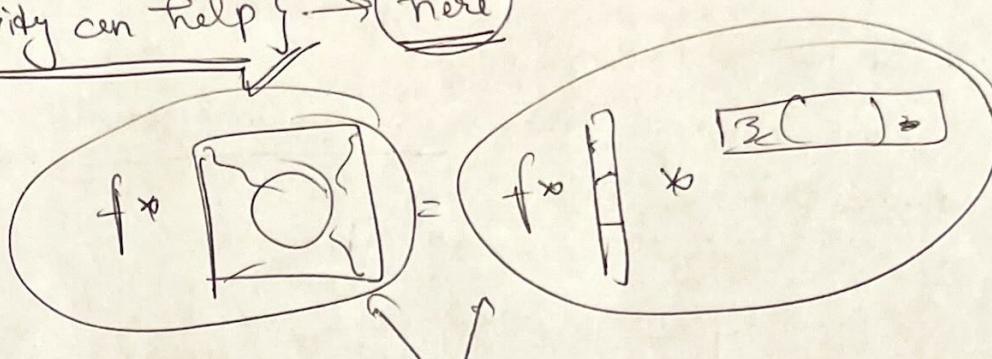
$$G(n, y) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{n^2+y^2}{2\sigma^2}\right)$$

$$= \left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{n^2}{2\sigma^2}\right) \right) \left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{y^2}{2\sigma^2}\right) \right)$$

Ans.

(4)

{Associativity can help} \rightarrow here



gives us
the same result

① Kernel separable then complexity \Rightarrow

$$O(n^2 m)$$

Ans.

Gabor

(cosine & gaussian)

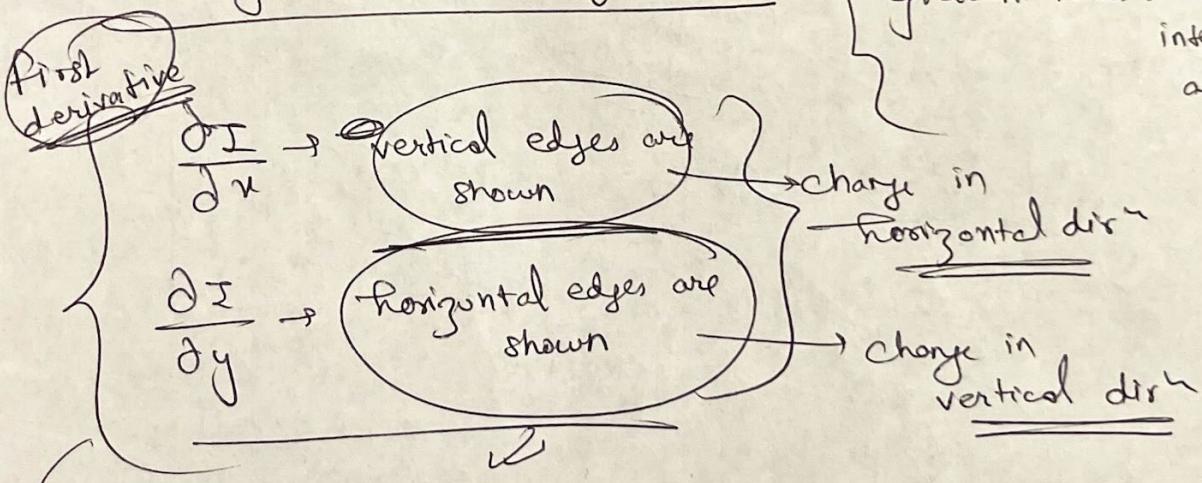
$$\Psi(u, y) = e^{-\frac{u^2 + y^2}{2\sigma^2}} \cos(2\pi\mu u)$$

biologically
related to
cortex receptive
field

Ans.

How do you take derivative of an image?

gradient tells \rightarrow how does intensity change along an axis.
(First it's derived)



→ How can we do this using convolution?

$$\text{if } I_{\text{image}} * [-1, 1] \Rightarrow \frac{\partial I}{\partial x}$$

for rate of change,
take second derivative

$$I_{\text{image}} * [-1, 1]^T = \frac{\partial I}{\partial y}$$

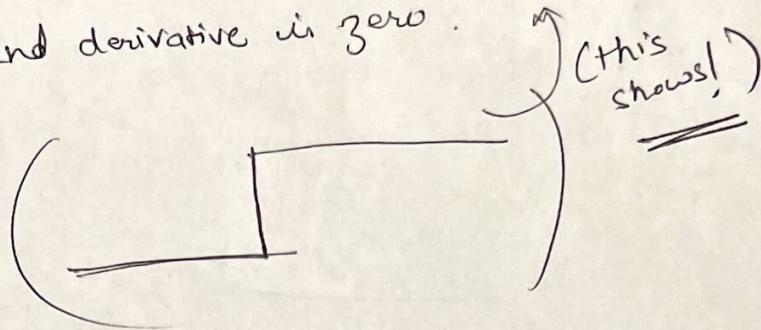
$$\frac{\partial^2 I}{\partial x^2}, \frac{\partial^2 I}{\partial y^2}$$

$$\textcircled{1} * [-1, 1]^T$$

Ans.

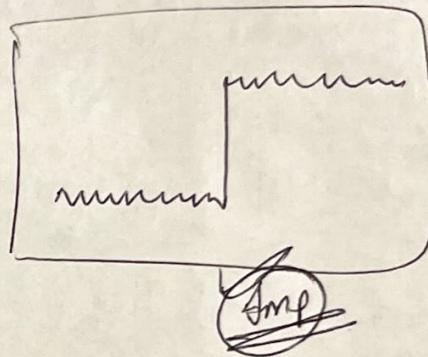
Edge?

- Change is measured by derivative in 1D.
 { — Biggest change, derivative has maximum magnitude. (large magnitude)
 } — Or 2nd derivative is zero.



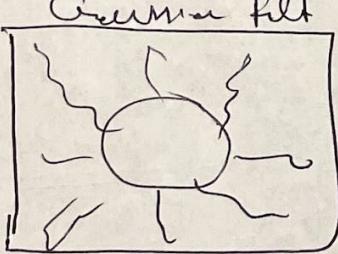
What about noise?

- Derivative is high everywhere.
 → Must smooth before taking gradient!



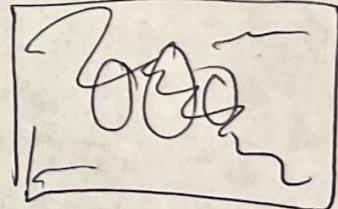
Hunting Noise

- filter with gaussian to smooth, then take gradients.
 → But, convolution is linear.



$$\star [-1, 1] \star [-1, 1] =$$

Transpose these
to find
horizontal edges ✓



(Charge in
hor. dir.)
Laplacian filter

(charge in
vertical dir.)
(edges will be
more prominent!)
(That's why!)

Doing first $[-1, 1]$ & then $[1, -1]$ will give blank
(nothing) → might only see noise.

Laplacian filter defined as sum of second order partial derivatives

$$\nabla^2 I = \frac{\partial^2 I}{\partial x^2} + \frac{\partial^2 I}{\partial y^2}$$

Ans.

5

finding boundaries

$$(f * g) \text{ depth}$$

$$\frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} > \lambda$$

why a threshold?

2 reasons -

- (1) To take coarse lines/edges rather than taking into consideration all the fine lines/edges.
- (2) for visualization purposes, we binarize the output image to get either black/white pixel values. (giving edges on a background of black).

finding things

If we take a filter with an image of an object, then -

$$f * (eye \text{ image in the filter}) \Rightarrow (f * g) \text{ (weird blurred image).}$$

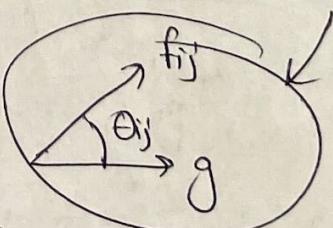
cannot be used for image detection!.

(a difficult problem!) → need neural nets for this.

Response for one window

$$(f_{ij}) \xrightarrow{D} f$$

$$f_{ij}^T g = \|f_{ij}\| \|g\| \cos \theta_{ij}$$



(norm of f and g gets multiplied which won't make sense!)

$$f * (g - \bar{g})$$

to find the filter!

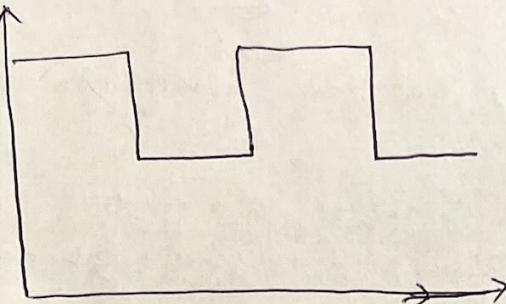
after thresholding

false true's and true true's

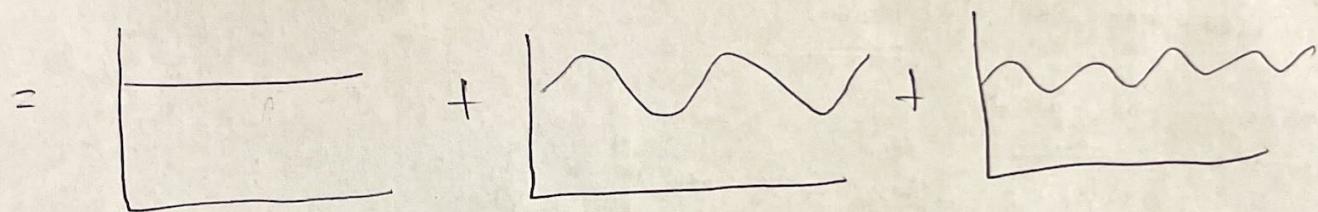
As

Fourier Transform

Any univariate func' can be rewritten as a weighted sum of sines & cosines of different frequencies.



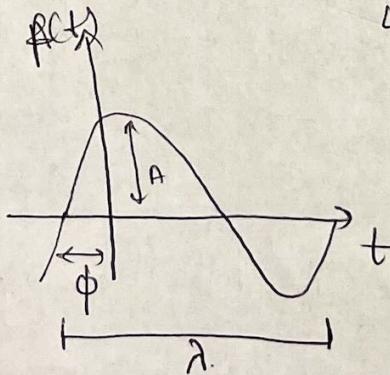
Q. How to build this signal using sin waves?



(can be applied to images)

+ ... when still ∞ , freq. increases of the waves.

Sinusoids \Rightarrow



$$\boxed{f(t) = A \sin(\omega t + \phi) = A \sin(\omega t + \phi)}$$

A = amplitude
 ϕ = phase
 f = frequency

$\boxed{\text{Signal} = \text{Amplitude} + \text{phase}}$
(vs. freq)
(vs. freq)

(also applicable to images)

freq. is implicit

black $\rightarrow 0$
white $\rightarrow 1$

treated as signal

black $\rightarrow -ve$
white $\rightarrow n \text{ large } (20)$
grey $\rightarrow zero$

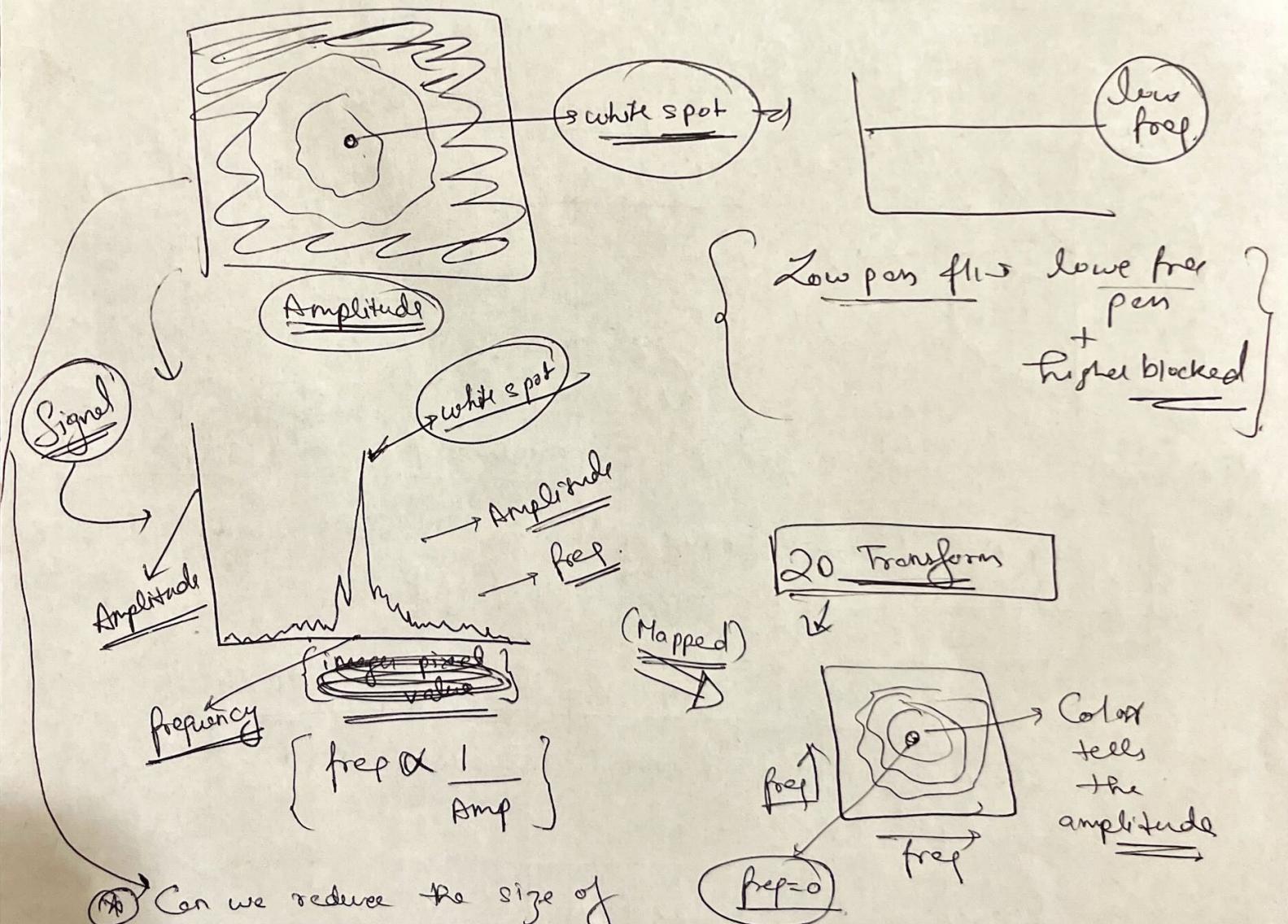
Q # Q is some to represent image and its Fourier components.

My (Most of the info. is in -the phase!)

(Amplitude has less info. about the image.)

(high freq. data is mostly stored in the phase → which give the details of the image).

{ Simplistic → basically mostly gives you the low freq. part of }
the data.



④ Can we reduce the size of the image using this?

As high freq. waves are mostly black (0) in terms of amplitude which doesn't provide us with a lot of data about the mag. So, maybe we can do so.

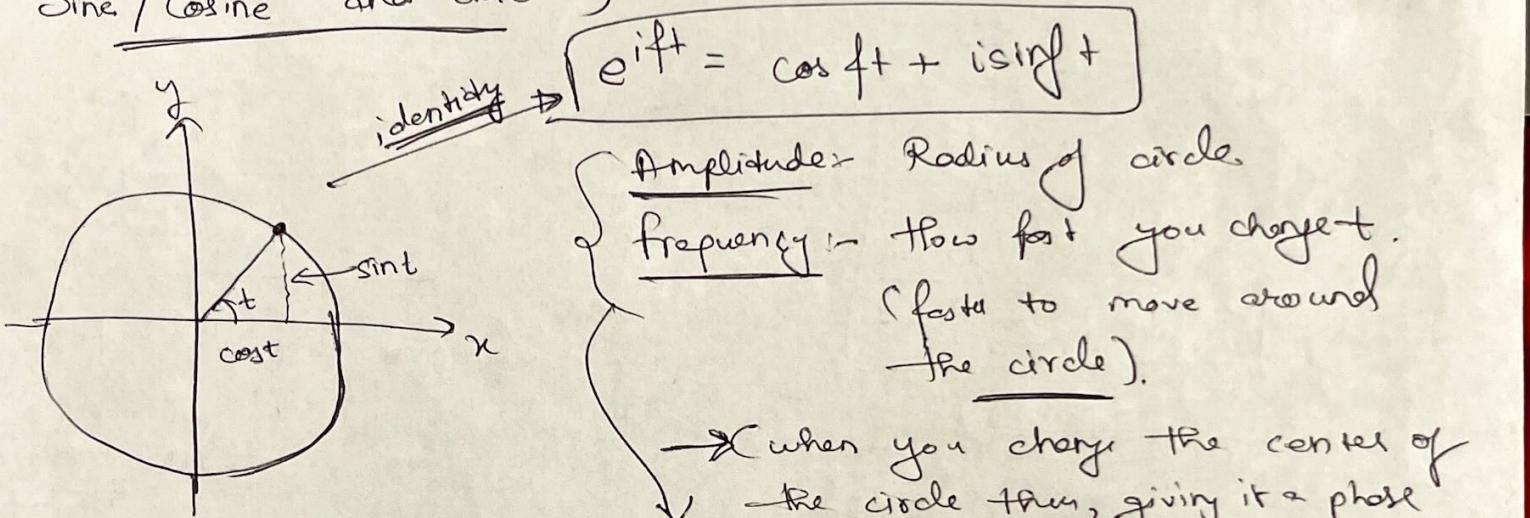
~~Ques~~ All the signals (components of fourier transform) can be mapped to amplitude & phase, which when combined would give us amplitude & phase for the whole image.

→ Convolution basically does this transformation for us.

→ CFT is also a linear transformation.

* How fourier transform works?

Sine / Cosine and circles



* Each circle is mapped to a particular freq. and amplitude.

$$g(t) e^{-2\pi ift}$$

maps $g(t)$ on to the unit circle with freq. f .

this is essentially the amplitude.

→ You wrap $g(t)$ around the circle with freq f , then calc avg pos. of $g(t)$.

$$G(f) = \int_{-\infty}^{\infty} g(t) e^{-2\pi ift} dt$$

when f move from $-\infty$ to ∞ .

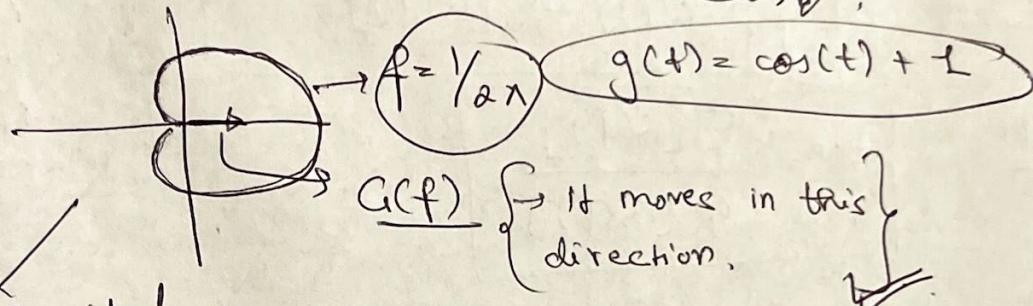
as if we move indefinitely.

$g(t)$ transformed to the freq. domain.

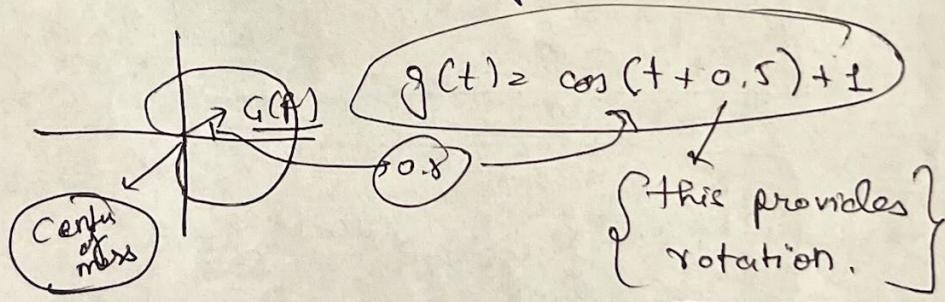
- Very low freq goes around the circle very slowly. (7)
 (some signal can be mapped around.)
- when freq. doubled, you would go faster.
- If we keep inc. f, then move till ∞ , you will generate diff. kind of patterns
- ↳ this can be mapped to image data
 → then (to generate diff. pixels).

$$f = \frac{1}{\lambda}, \frac{1}{2\lambda},$$

↳ it translates to lines (rather than covering areas), ✓.



when provided as phase,
 → $G(f)$ would move (shifts).

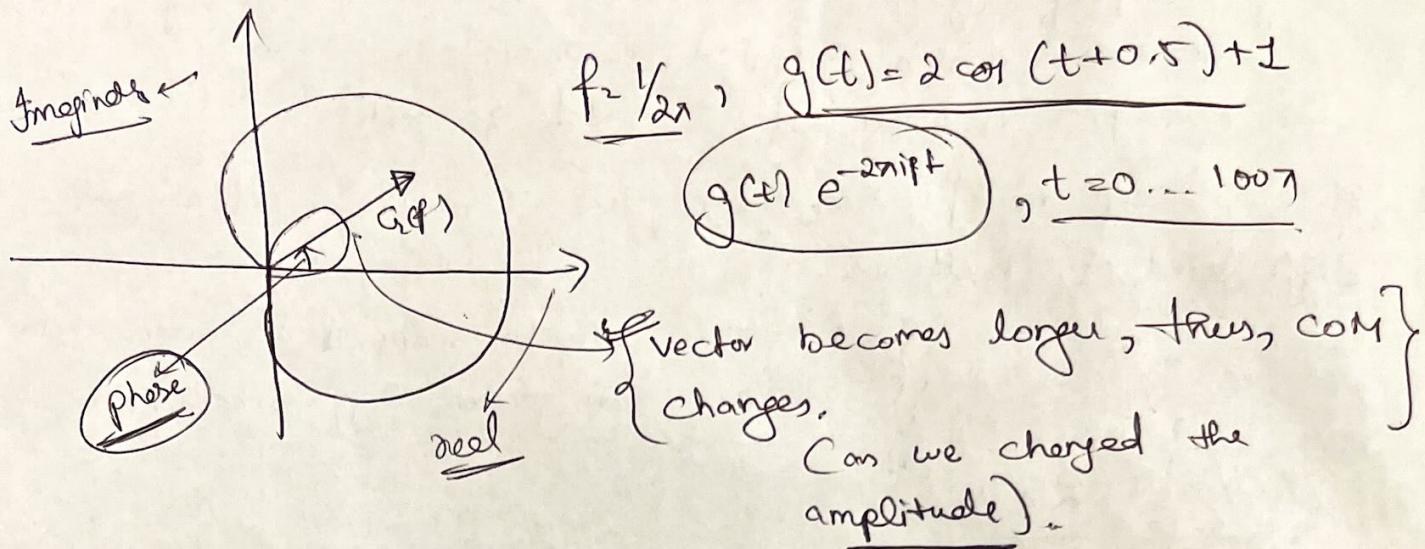


→ Amplitude → radius of the pattern.

→ frequency → how fast you move along the pattern.
 (this would change the avg. position in the pattern)

→ Phase → this rotates the pattern. (rotation of the whole pattern)

~~COM is only affected by Amplitude.~~

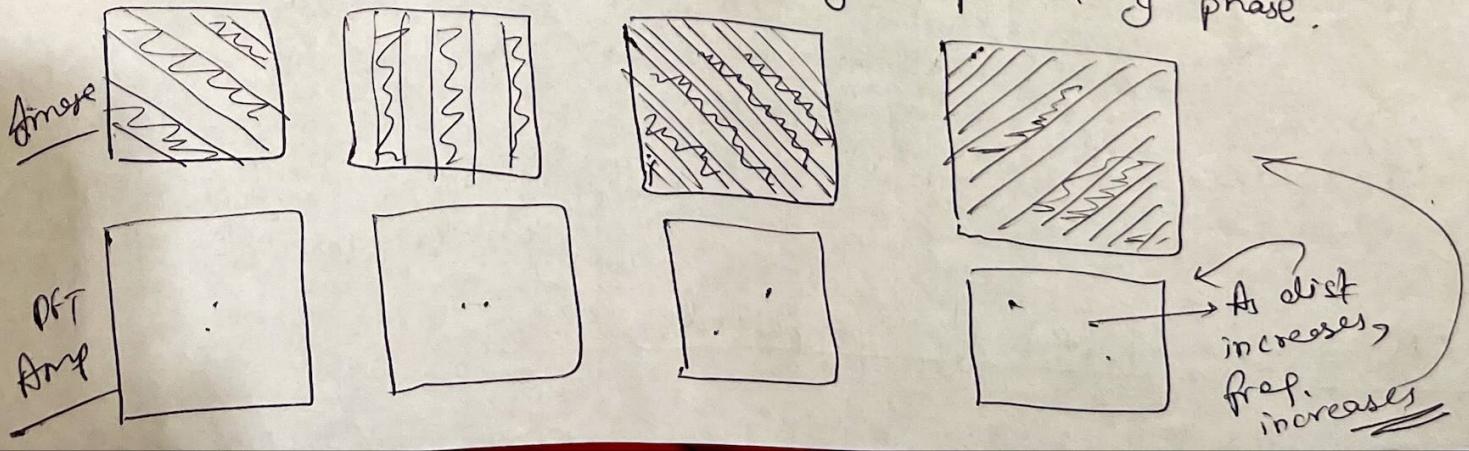


Amplitude → $\sqrt{R[G(f)]^2 + I[G(f)]^2}$
Phase → $\tan^{-1} \frac{I[G(f)]}{R[G(f)]}$

(Real & Imaginary Components)

Inverse FT → $g(t) = \int_{-\infty}^{\infty} G(f) e^{j\omega t} dt$
FT → $G(f) = \int_{-\infty}^{\infty} g(t) e^{-j\omega t} dt$

→ You can reconstruct an image using amplitude & phase.

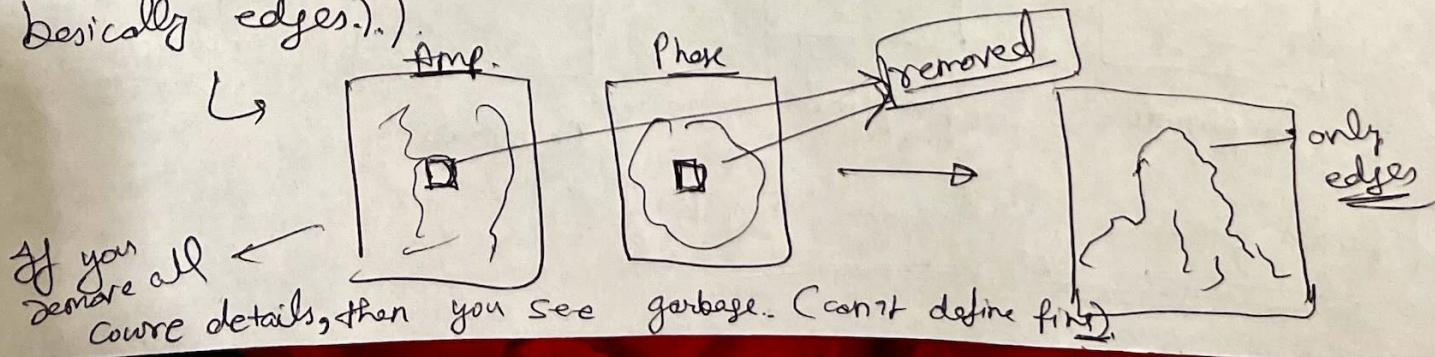


- ⇒ FT has peaks at spatial frequencies of repeated structures (for repeated info. in the image, spatial freq. data changes in a repeated manner). (8)
- ⇒ Acc. to grad changes, you get similar patterns in the amplitude part in accordance with the frequencies (freq. of change).
- ⇒ Might be both in the x- & y-directions.
(Depending on sharp changes in the corresponding directions).
- { Ex:- Brick wall image example. }

{ we can even throw away some of the values that are very close to zero or zero and only store the frequencies that contributes a lot. (using a low-pass filter).
Could hold a threshold that could work as a filter, then, giving us compressed image. ↑
(How JPEG works?)

→ when we convolve and arg.-out some of the pixels data, that removes the high freq. data (mimics a low-pass filter). → this is actually convolution mapped to Fourier transform.
~~As more convolution map (larger), gets more blurred)~~

→ for edge detection, we basically use high-pass filter. (that lets go high freq. details (only the fine ones) → which are basically edges.).



→ Convolution in time space is multiplication in freq. space

$$g(x) * h(x) = \mathcal{F}^{-1} [\mathcal{F}[g(x)] \cdot \mathcal{F}[h(x)]]$$

→ Conv. in freq. space is multiplication in time space

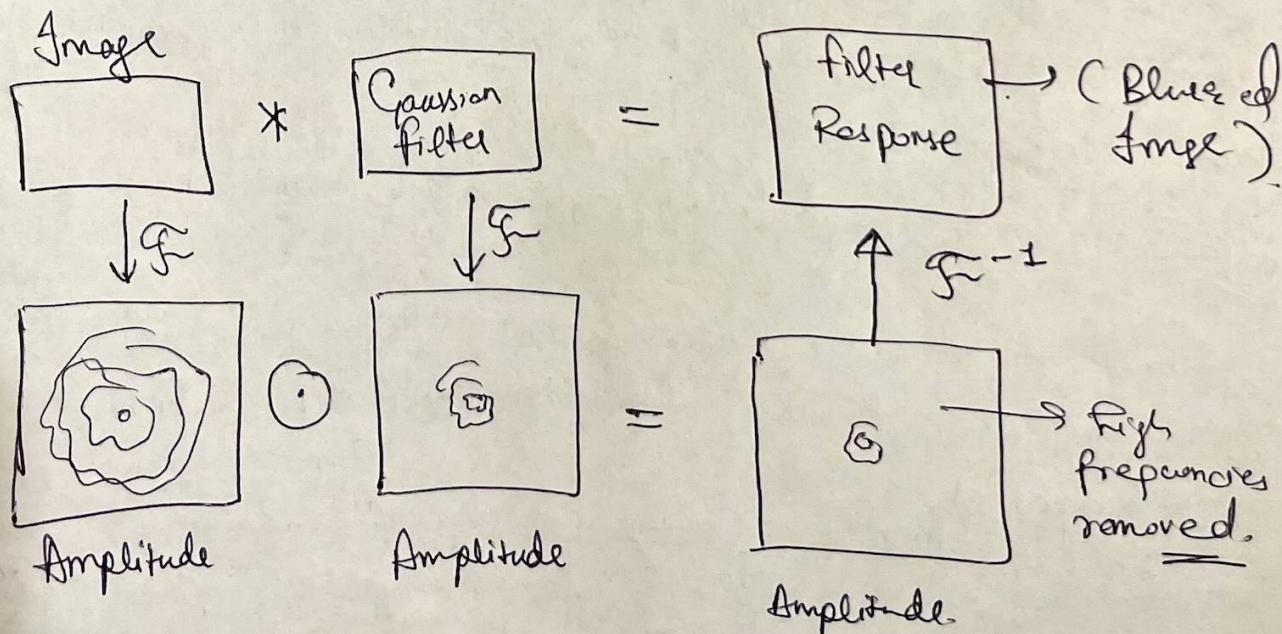
$$\mathcal{F}[g(x)] * \mathcal{F}[h(x)] = \mathcal{F}[g(x), h(x)]$$

Which is more computationally efficient?

Depends on how fast, Fourier transform can be done.
[from time \rightarrow freq. or freq \rightarrow time]

- If filter is large, then taking Fourier transform and then taking inverse will be faster.
- If filter is small, then direct convolution will be faster.

Steps:



[Process gives us the same thing!]

Mean filter \rightarrow Low-pass filter \rightarrow less blurred image than Fourier transform from gaussian.

Gaussian \rightarrow Low-pass filter \rightarrow blurred image \leftarrow Fourier transform gives

Laplacian \rightarrow High-pass filter \rightarrow edges (edge detection) ⑨

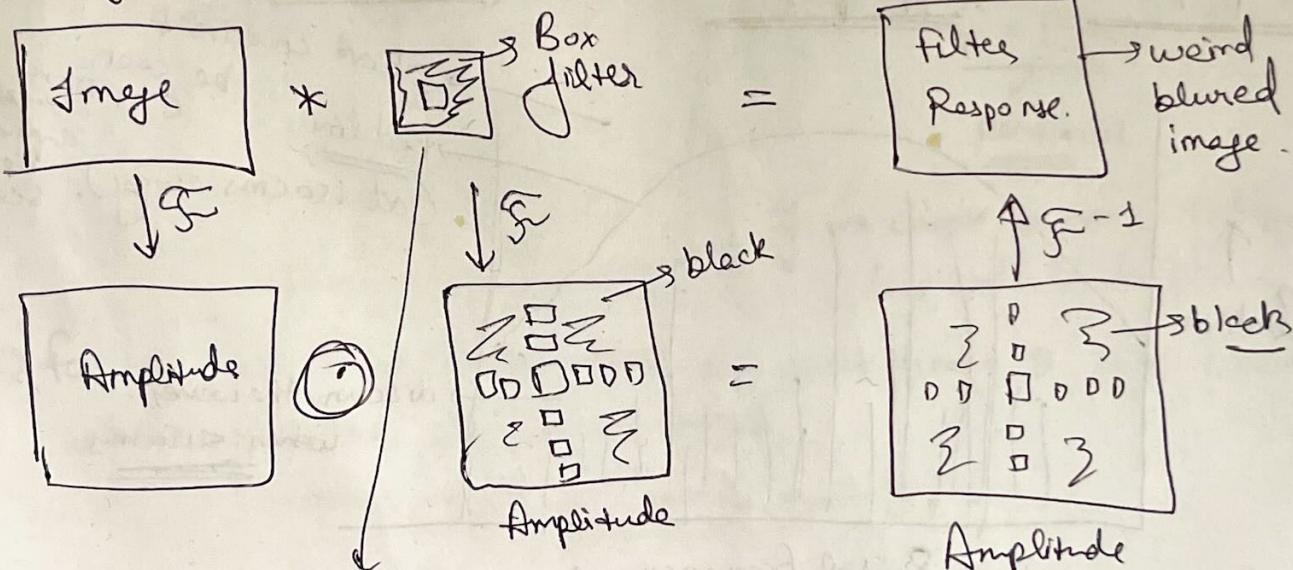
Sobel \rightarrow High-pass filter \rightarrow only vertical edges detected (some edges might be lost)

Sobel-y \rightarrow High-pass filter \rightarrow only horizontal edges detected.

Box filter \rightarrow High-pass filter \rightarrow edge detection

but for Scharr-X \rightarrow we will mostly all vertical edges

\rightarrow FT of gaussian is gaussian.

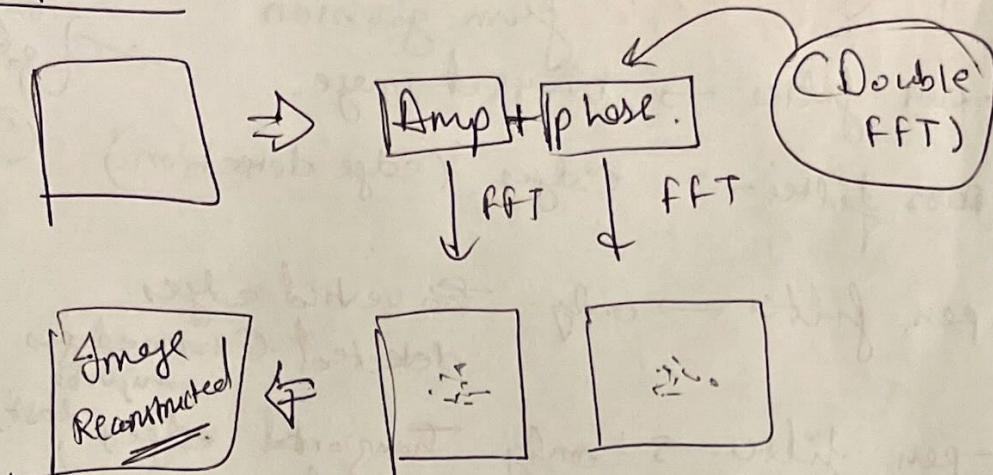


$\rightarrow f \rightarrow \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$

Laplacian filters

In FT space, this is High-pass

Compression \Rightarrow



(compressed by a lot!)

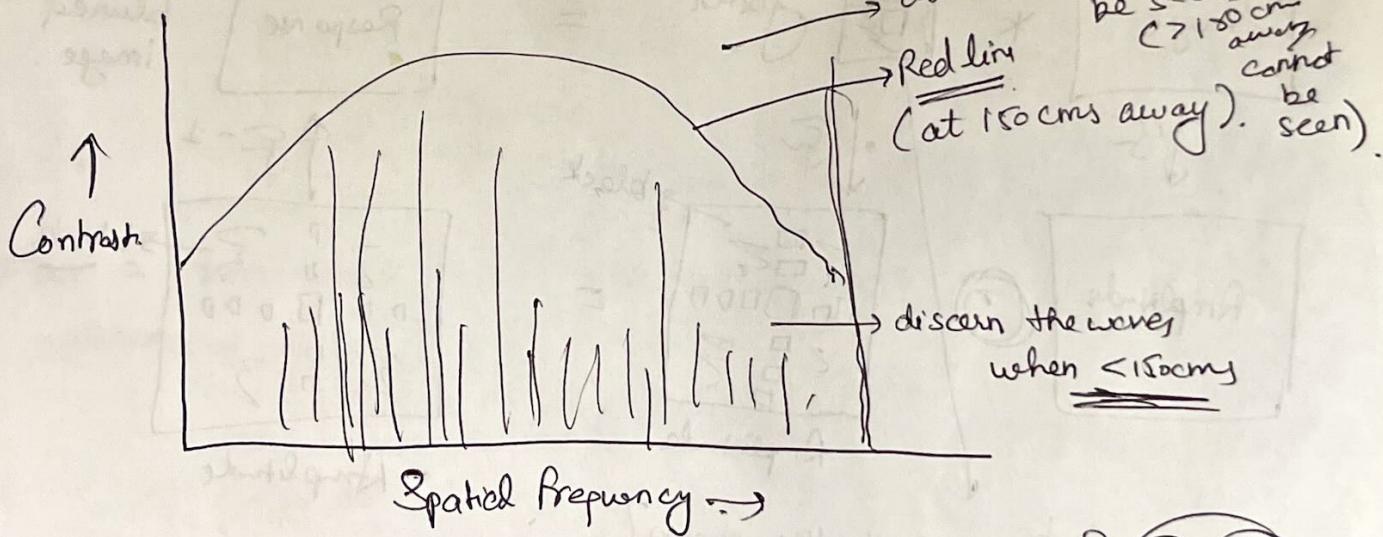
Hybrid Images \Rightarrow (for creating illusions)

FT of image 1

FT of image 2

$> 150 \text{ cms away}$

Low frequencies + High frequencies \Rightarrow Hybrid Image (Fails image within crooked)



Spatial Frequency \Rightarrow

① When far away, can see the low frequencies.

② When close, can see the high frequencies.

Result

Ans.

Machine Learning

(10)

→ we need translation invariance. (Same operation everywhere).
& scale invariance.

Regression basics:-

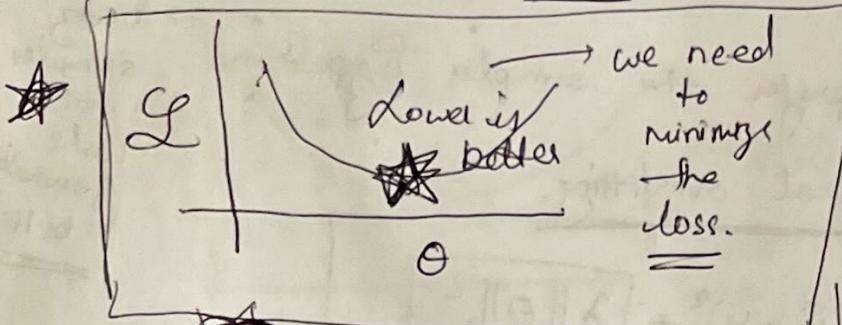
Model: $\hat{y} = f(x; \theta)$; Input image: x , Parameters θ .

Loss function: $L(\hat{y}, y)$

pred.
values

target labels

Prediction



$$\hat{y} = f(x; \theta)$$

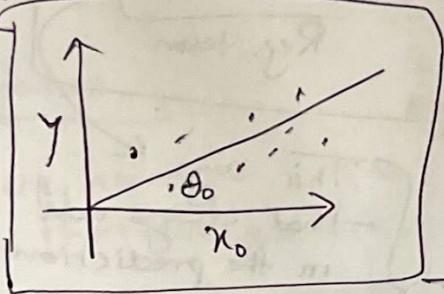
$$= \theta^T x$$

$$= \theta_0 x_0 + \theta_1 x_1 + \dots$$

$$L(\hat{y}, y) = (\hat{y} - y)^2$$

$$\hat{y} = \theta_0 x_0$$

Parameter



{ If we choose parameter badly, loss will be larger. }

⇒ for a non-linear relationship, you can't find a θ with linear relationship where loss is minimized,

Now, we use polynomial models.

2nd order polynomial:-

$$f(x; \theta) = \theta_0 x_0 + \theta_1 x_1 + [\theta_2 x_0^2 + \theta_3 x_1^2] + \theta_4 x_1 x_2 + \theta_5$$

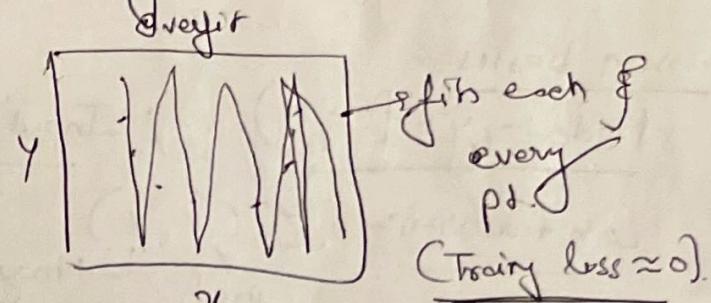
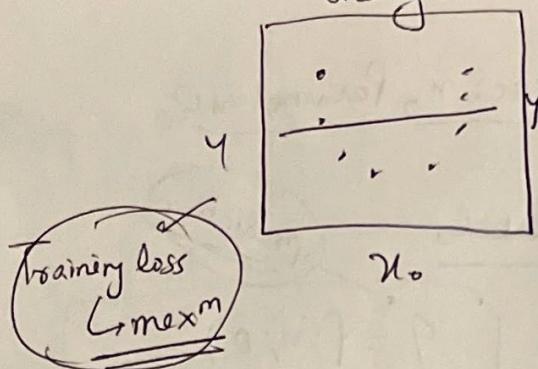
$$= \theta^T \phi(x) \quad (\text{we change the feature space})$$

Transformation function (transform the space)

{ S.t we can do linear relationship }
in that feature space.

(*) Nth order polynomial possible.

If we go beyond a point where ~~N~~ N^{th} order polynomial overfits the data (all predictions = labels). When tested, would perform poorly.



→ Occam's Razor:— We prefer the simpler hypothesis. bcz simple soft will generalize better.

Regularization To prevent overfitting.

$$\mathcal{L}(g, y) = \text{Data term} + \frac{\lambda \|\theta\|_2^2}{\text{Reg. term.}}$$

Common Loss func \rightarrow

Sq. error:—

$$\mathcal{L}(x, y) = \|x - y\|_2^2$$

This says that always add something in the prediction so, that loss doesn't go zero because of the data term.

Hinge loss:—

$$\mathcal{L}(x, y) = \max(0, 1 - x \cdot y)$$

minimize this

Cosine Similarity:—

$$\mathcal{L}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

maximize this

Learning

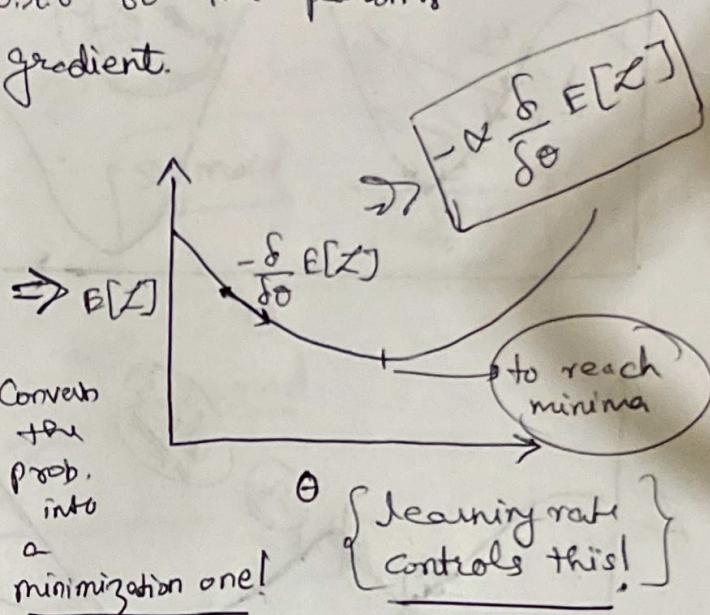
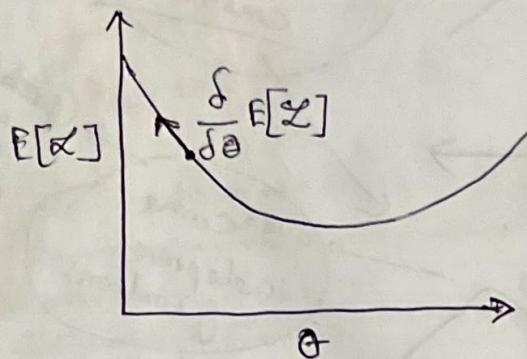
$$\min_{\theta} E_{(x_i, y_i) \in D} [\mathcal{L}(f(x_i; \theta), y_i)]$$

What are the params that minimize the value of the loss func in expectation over the entire dataset?

We want to minimize the diff. b/w prediction and the actual label, as we want to minimize the expectation of loss.

Gradient Descent:-

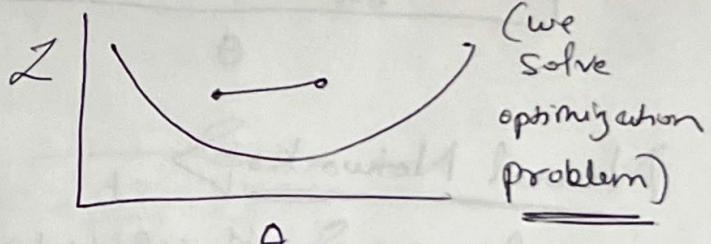
- ~~Steps~~
1. Make an initial guess for params.
 2. Calc. gradient of loss w.r.t to the params
 3. Step parallel to the gradient.



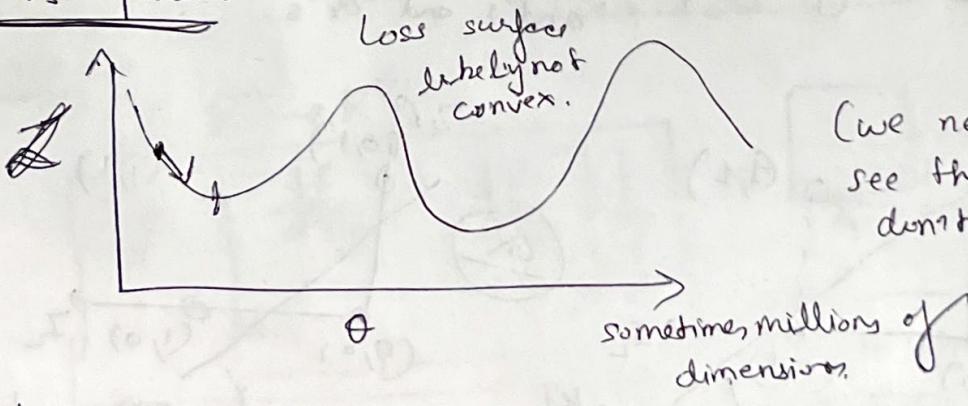
Convex loss surfaces:-

local = global minimum

If the 2 points are above, then everything in between is also above.



Local vs. Global:-



(we need to see that we don't get stuck on local maxima/minima).

- ④ We need momentum here to see if we aren't stuck at a local maxima/minima (captive)

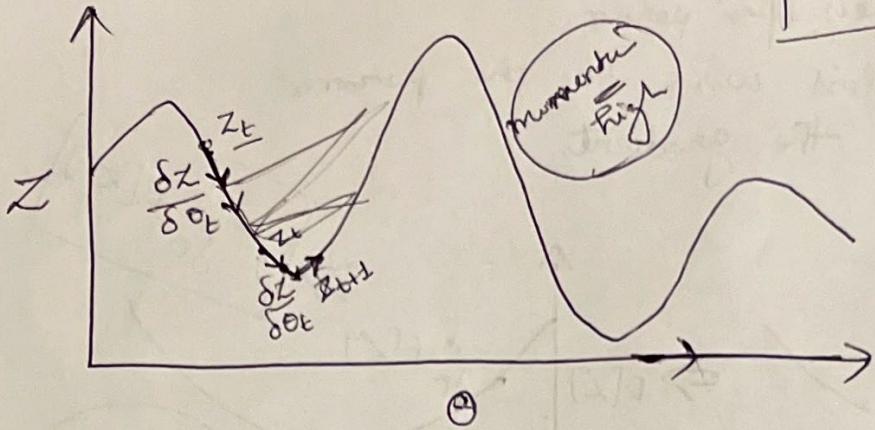
Gradient descent

$$\Theta_{t+1} = \Theta_t + \alpha \frac{\delta L}{\delta \Theta_t}$$

Grad. descent with momentum

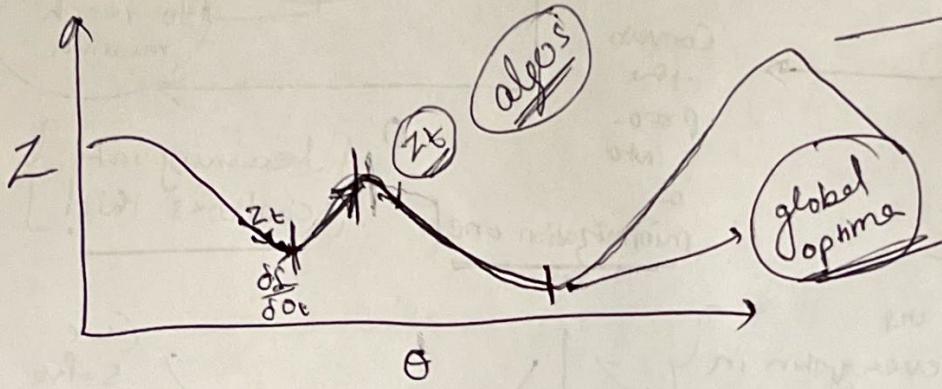
$$Z_{t+1} = \beta Z_t + \frac{\delta L}{\delta \Theta_t}$$

$$\Theta_{t+1} = \Theta_t + \alpha Z_{t+1}$$



search space constraint

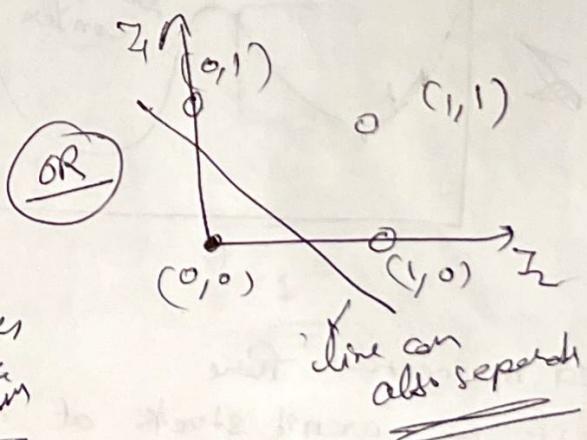
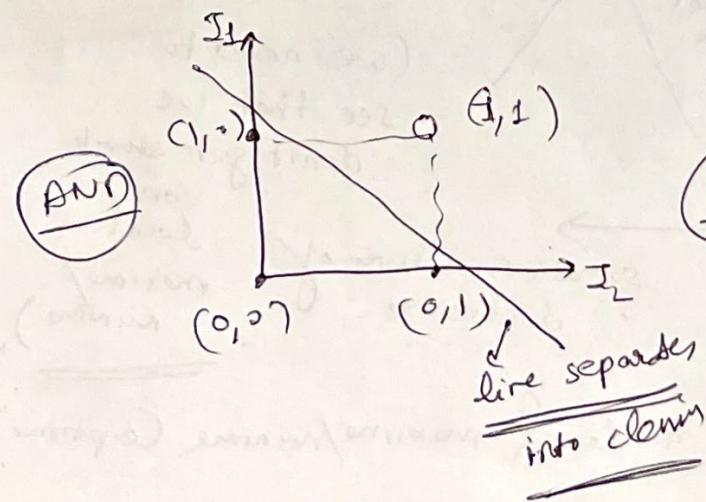
(problem w/m gradient)

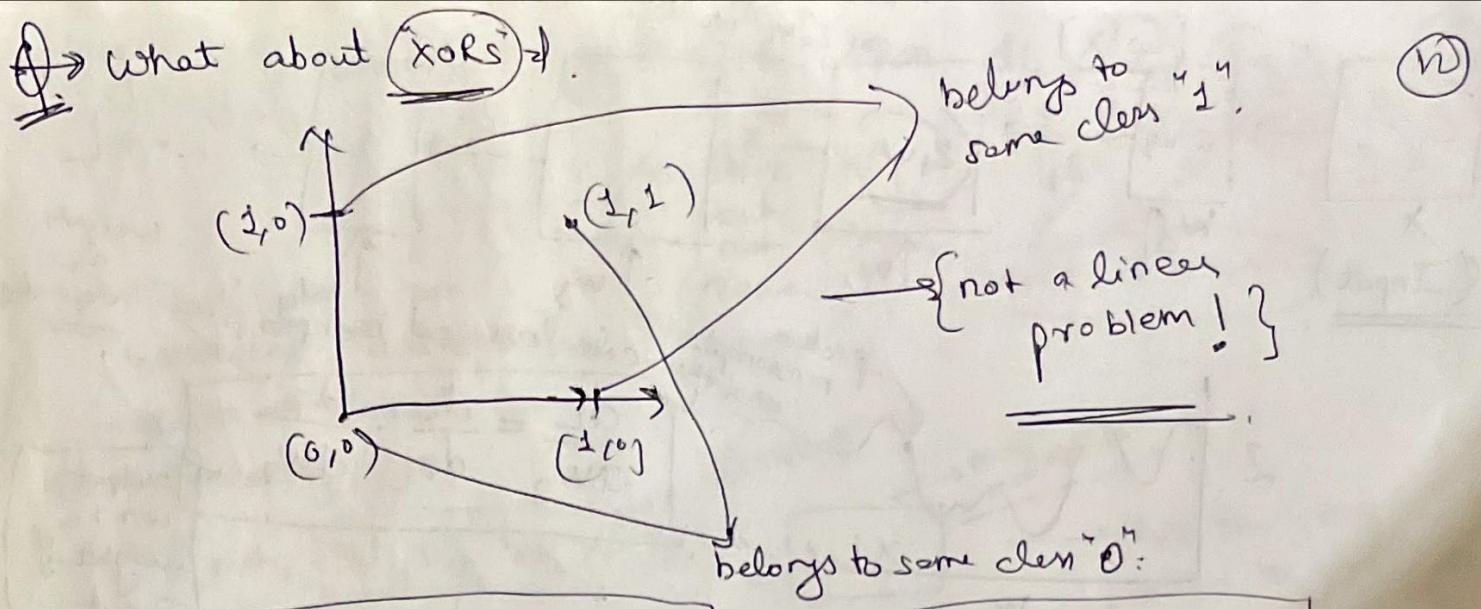


overcome adaptive gradient
+ momentum
Adem

Neural Networks

Perceptron \Rightarrow If I/P features are binary, can model logical "and's" and "or's".



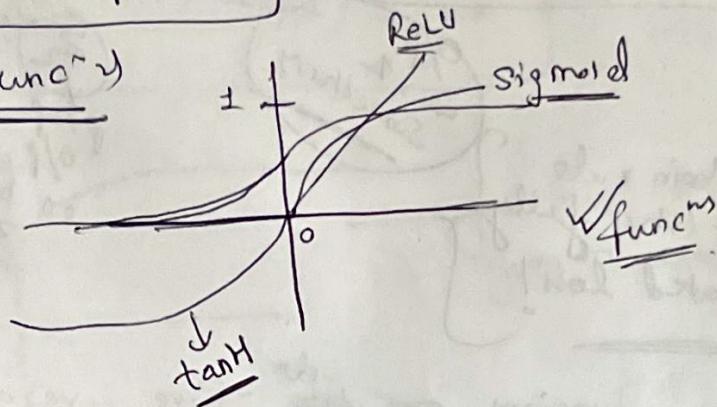


$$x_{i+1} = \phi(w_i x_i + b_i)$$

↓
Activation funcⁿ

Activation funcⁿ

$$\text{Sigmoid} := \frac{e^{-x}}{1+e^{-x}} = \frac{1}{1+e^x}$$



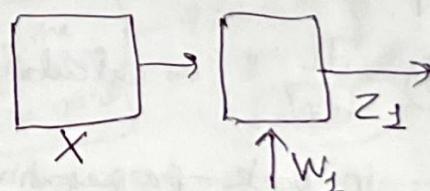
Chain Rule:-

$$y = f(x), z = g(y) \Rightarrow \left[\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx} \right]$$

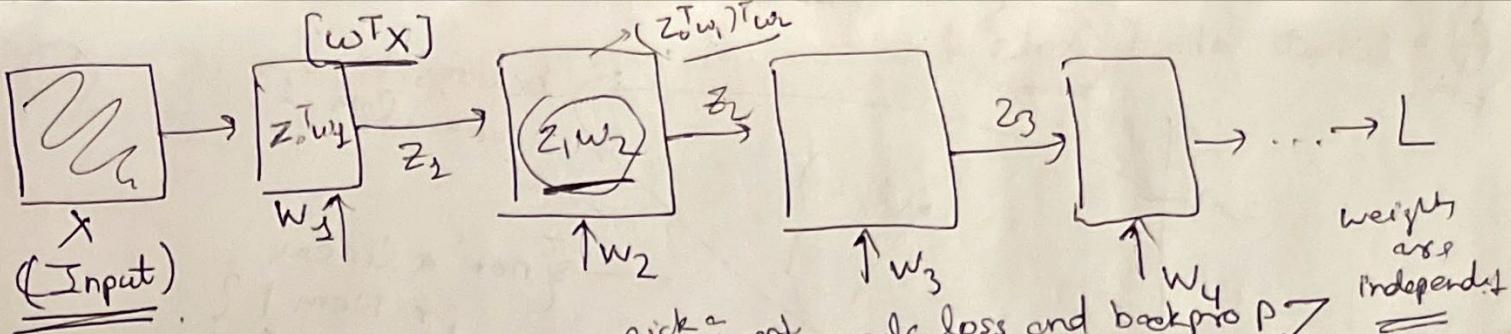
Jacobian:-

$$f: R^n \rightarrow R^m$$

$$J = \begin{bmatrix} \frac{df}{dx_1}, \dots, \frac{df}{dx_n} \end{bmatrix}_{n \times m}^T$$



$$\underline{z_1 = w_1 x}$$



pick a random pt and then calc. loss and backprop P ↗

$$\frac{dL}{dW_2} = \frac{dL}{dz_2} \cdot \frac{dz_2}{dW_2}$$

weights are independent

output w.r.t weight & input

$$\frac{dL}{dz_1} = \frac{dL}{dz_2} \cdot \frac{dz_2}{dz_1}$$

~~Imp~~
can apply chain rule now to get the fully backpropagated loss!

O/P w.r.t something

this was bcs.

z_1, z_2 aren't dependent on each other

z_2
prop O/P
on both
 z_1, z_2 ans

→ Most activation functions ~~do~~ have -ve as well as +ve as input and when the inputs are normalized, it would mirror those activation functions, that's why standardization seems so important and also because of exploding gradient problem.

④ $\frac{dL}{dW_{n-1}}$ is calculated before $\frac{dL}{dW_{n-2}}$ so, that it can be reused in back-propagation. (reuse the computations already done from the end to the start).

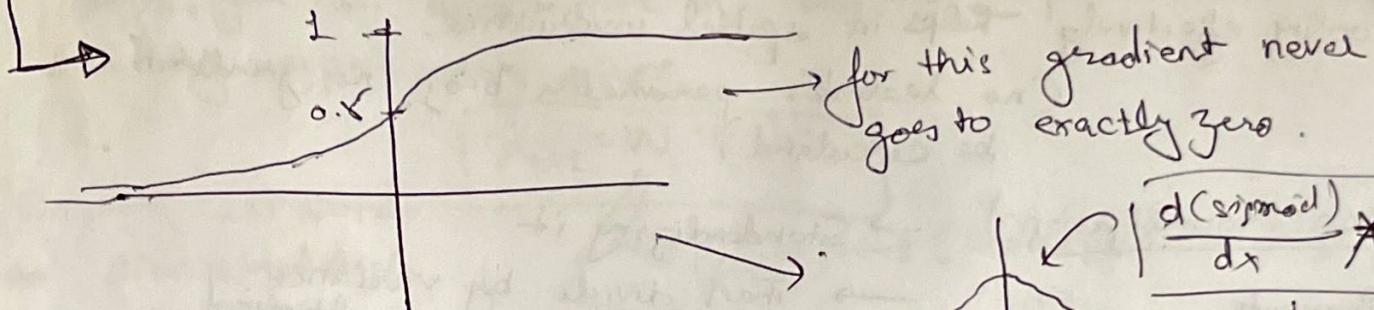
→ Vanishing gradient → if $\frac{dL}{dW_{n-1}}$ goes to zero, all $\frac{dL}{dW_{n-2}} \dots \frac{dL}{dW_1}$ goes to zero. (backprop problem).

\Rightarrow if $\frac{dL}{dW_n}$ shouldn't go to zero, then $\frac{dL}{dz_1} \cdot \frac{d z_1}{d z_2} \cdot \frac{d z_2}{d W_n}$ (13)

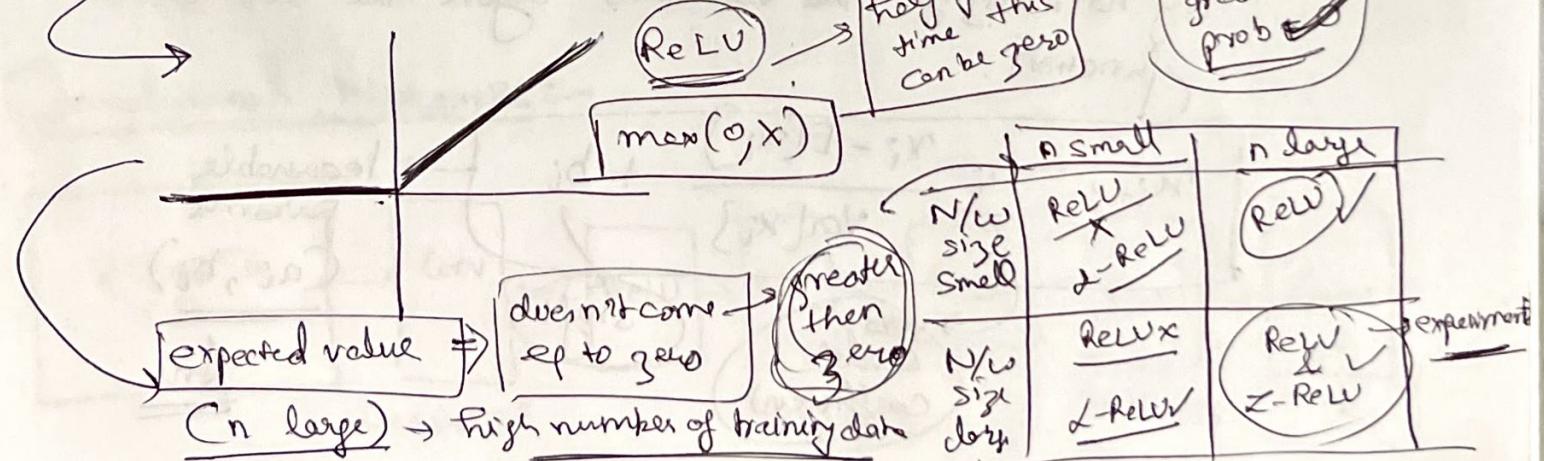
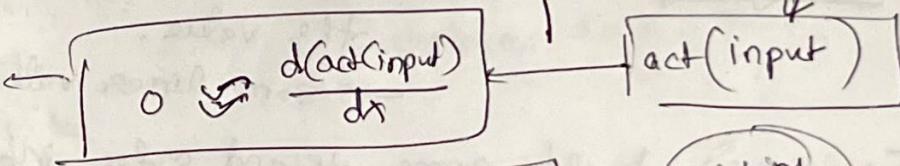
That's why we introduced activation function to introduce non-linearity.

→ if not, $[W_4 \ W_3 \ W_2 \ W_1 \ x] \xrightarrow{\text{linear op}} [Wx]$ → this has the potential to give lots of zeros.

$\dots \text{act}(W_3 \text{act}(W_2 \text{act}(W_1 x))) \dots \rightarrow$ non-linearity introduced.



effectively zero = if this goes really small



expected value \Rightarrow doesn't come up to zero
(n large) \rightarrow high number of training data

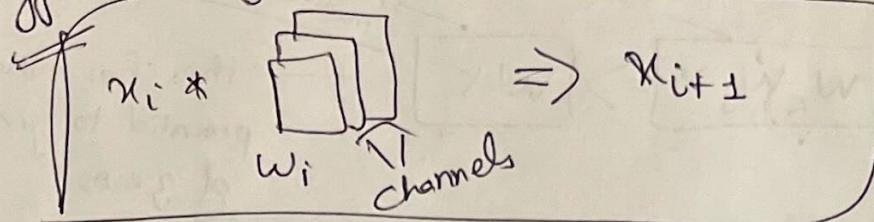
	n small	n large
N/w size Smell	ReLU	Z-ReLU
N/w size large	ReLUx	Z-ReLu
experiment	ReLU & Z-ReLu	

$$x_i * w_i = x_{i+1} \rightarrow \text{for Convolutional layer.}$$

↓
Conv. filter

$$x_{i+1} = w_i * x_i + b_i$$

→ for each layer, there might be some n filters. that detect different features in the image.



* Max Pooling → [Non-linear convolution] → only outputs the max value in the window spanned by the filter in a point in time.

→ pooling effectively helps in spatial invariance.
(no learnable parameters bcoz no gradients can be calculated).

Batch Normalization → Standardizing it
(Prevents overfitting)

- Then divide by variance
- Scalar coeff → ~~shift multiplied by added.~~
- some linear value for a batch.

Intuition → build some second-order information into the new by normalizing the variables before the activation function.

$$x_{i+1} = \frac{x_i - E[x_i]}{\sqrt{\text{Var}[x_i]}} + b_i$$

learnable params
(a_i, b_i)

Ans

Scalar coefficient

shift

Activation functions \Rightarrow

$$\text{Softmax} \quad x_{i+1}^j = \frac{\exp(x_i^j)}{\sum_k \exp(x_i^k)}$$

(multi-class)

(14) probability for an op over probabilities for all the ops. (the full vector).

Dropout \Rightarrow

$$x_{i+1}^j = \begin{cases} x_{i+1}^j & \text{with prob. } p \\ 0 & \text{otherwise.} \end{cases}$$

(to prevent overfitting).

(To provide robustness to the N/w nodes).

Width vs. Depth \Rightarrow

With sufficiently wide networks and just one layer, you can approximate any continuous function with arbitrary approximation error.

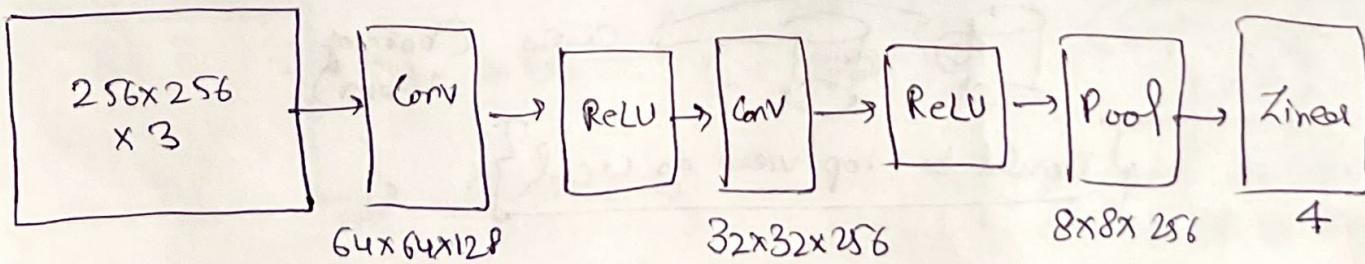
$$x_{i+1} = w_i x_i + b_i$$

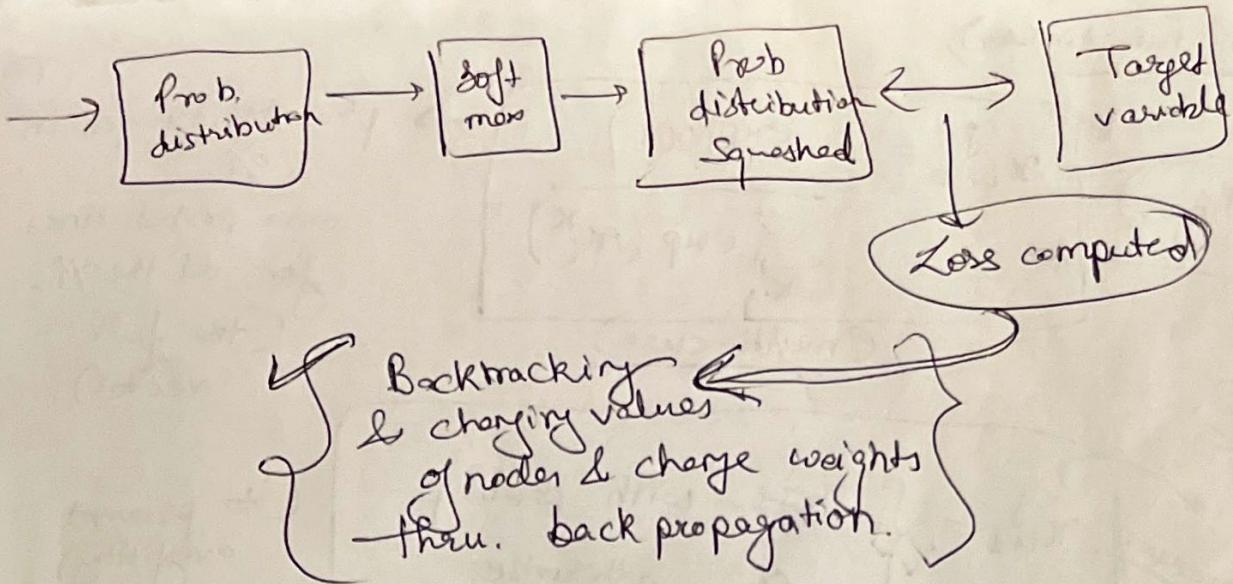
(wide) \rightarrow (corresponds to size of w).

(deep) \hookrightarrow corresponds to how many matrices do you have.

* Depth corresponds to dimensionality/complexity of the features while width corresponds to type of features that need to be extracted).

Convolutional Networks \Rightarrow





Cross Entropy Loss:-

$$L(y, \hat{y}) = H(y, \hat{y}) = \underbrace{- \sum_i y_i \log \hat{y}_i}_{\substack{\text{entropy} \\ \text{prediction}}}$$

Entropy
Average number of bits needed to encode y if the coding scheme from \hat{y} is used instead.
(the lower number of bits you need, the closer should be the distributions for y and \hat{y}).

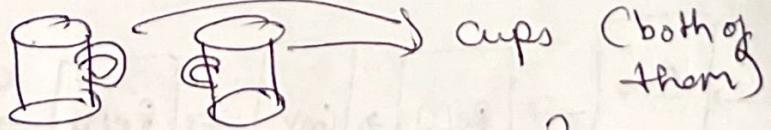
So, the loss would be less when distributions are similar.

Easy to calculate gradients.

Object Detection

Canonical Perspective! — "The best and most easily recognized view of an object".

This hinders learning if dataset is biased.



Could be top view as well

- Object Recognition is so hard because of different viewpoints.
- Changes in illumination.
- Change in scale.
- Background clutter.
- Within-class: lots of variations of an object.

18

Problems for object recognition.

Visual object categories ranges from 10,000 - 30,000 }

Lighting variations + view + other env. changes }

So many possibilities!

Crowdsourcing → paying people to do mini-tasks.

View of categories ↗

- A category is formed by some defining properties.
- Anything that matches all/ enough of the properties is part of the category.
- Anything else doesn't belong to the category.

Color categories

If 2 categories → white & black (Stage 1)

3rd added + Red (Stage 2)

4th added + Yellow/Green (Stage 3)

5th added + Yellow/Green (Stage 4)

6th added + Blue (Stage 5)

7th added + Brown (Stage 6)

8th Other colors

Perception can be done without any language)

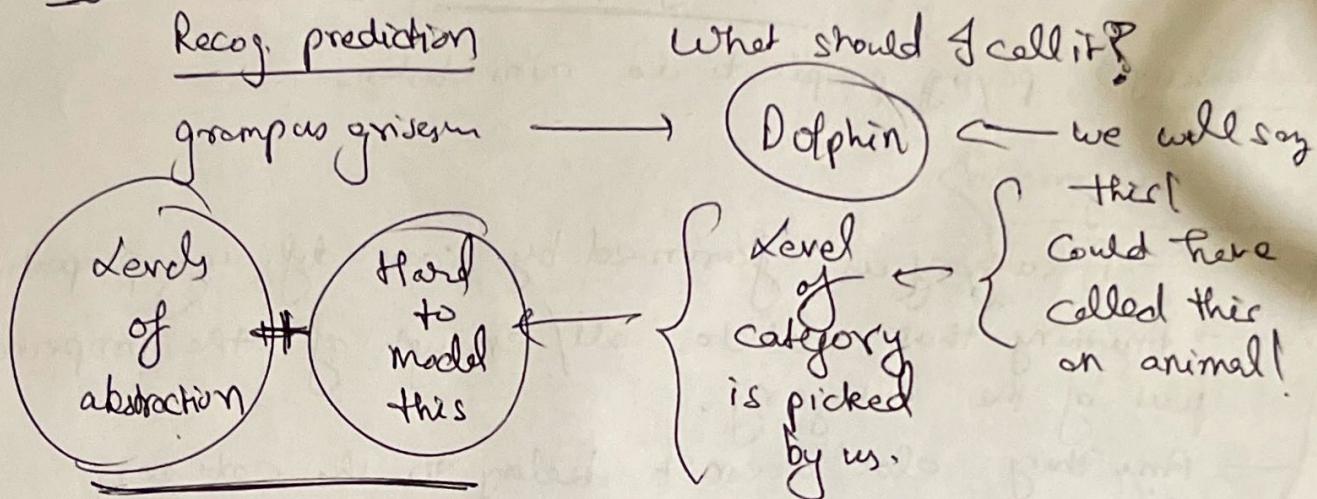
(Sepⁿ from NLP)

so on and so forth.
(Any two desc. could be diff.)

Affordance \rightarrow An object is defined by what action it affords!

(E.g. dog in front of computer) \rightarrow ? (object could be perceived differently by the machine).
{ Could help/not with training }
{ a model w.r.t perception }

Entry-level categories \rightarrow



④ Perception is always related/grounded in language!

Exemplar SVMs \rightarrow Learn a separate linear SVM for each instance (exemplar) in the dataset.

\rightarrow Each exemplar-SVM is trained with a single +ve instance.

\rightarrow Each exemplar-SVM is more defined by "what it is not" vs. "what it is similar to"?

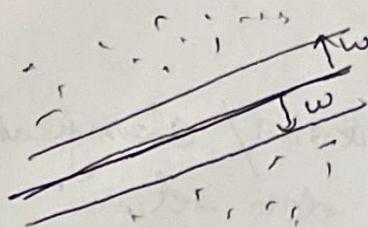
Large-scale training:

- \rightarrow Each exemplar performs its own hard -ve mining.
- \rightarrow Solve many convex learning problems.
- \rightarrow Parallel training on cluster.
(to parallelize workflow)

Support Vector Machine (Linear SVM):-

(16)

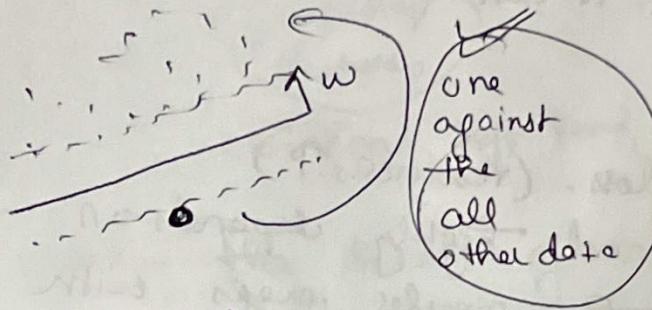
$$\min_w \frac{1}{2} \|w\|_2^2 + \sum_i \max(0, 1 - y_i w^T x_i)$$



→ when a data pt is perturbed,
the classification shouldn't
change!

(that's what SVM takes
care of).

for exemplar:-



→ used to classify
that single data
point.

Exemplar:- (in images)

The detector makes little windows to go over the whole image and detects different features in the image, considers diff features and/or geometry inside that. Then, that is translated to features in one of the images in the data set to match an appearance.

But this doesn't necessarily map really well and depends on the data set.

→ should have very large training data set.

Exemplar



This acts like a look-up table. It might also find correlations b/w diff. objects. (similarity b/w things).

for ex:- If exemplar detect a bike with some features of a person, it can say that a person is sitting on top of it.

→ This can also be used for segmentation/ classification of object from/in the images in the data set.

→ This look up table won't be exact but perturb it a little bit. (so, that it can match to images with little changes).

Problems:-

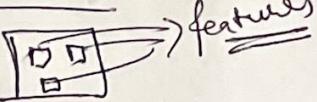
(1) Computation power loss. (resources ↑)
(2) It's too specific and highly depends on the training dataset, if it has similar images with similar geometry and view of the object (depends a lot on view of perception).

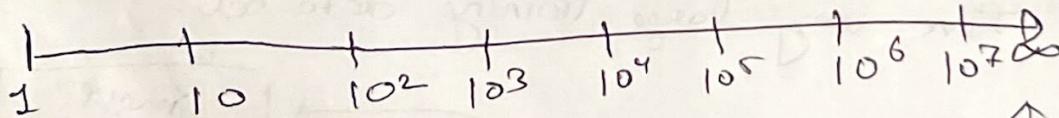
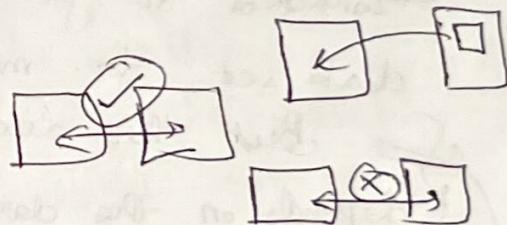
↑ that's why we need a lot of data!

↑
(This was an interpolation problem).

Now, we will go to the Extrapolation problem.

Two extremes:-

Finding features  features



Extrapolation ↑
problem

[Generalization,
Diagnostic features]

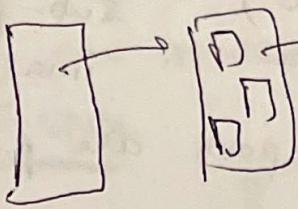
Interpolation problem

[Correspondence,
finding the
differences].

→ Deformable part models:-

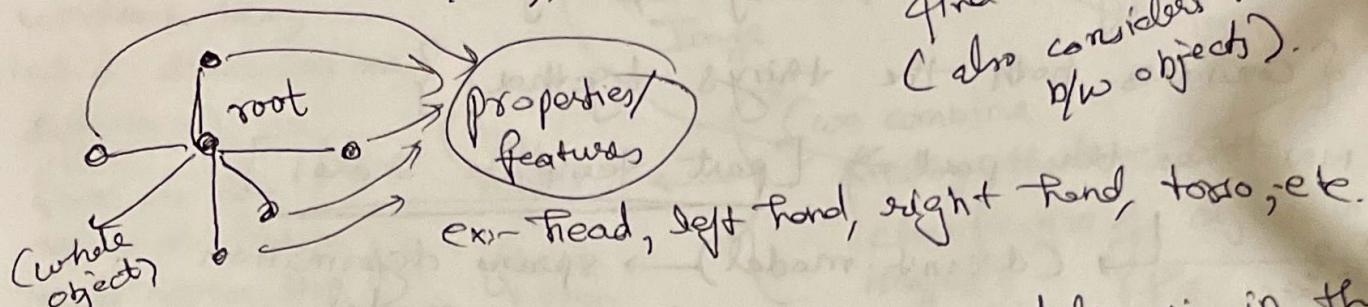
(18)

(Model encodes local appearance + pairwise geometry).



It finds features and encodes them together so if it is found together in any form, acc. to those properties, it will find that object (also considers relationships b/w objects).

Inference :- max score (x, z) .



{ Star model :- the location of the root filter is in the anchor point. }

{ Given the root location, all part locations are independent. }

(That's how properties can be mapped).

↳ this, help in classification of the whole object.

(we are trying to find relative relationships b/w sub-parts of an object).

→ This can work on different view-points as well!
Even different configuration & might apply to illumination settings as well).

Scoring function:-

$$\text{score}(x, z) = \sum_i w_i \phi(x_i, z_i) + \sum_{i,j} w_{ij} \Psi(z_i, z_j)$$

$x \rightarrow \text{image}$

$z_i \rightarrow \{x_i, y_i\}$

$z = \{z_1, z_2, \dots\}$

Part template Scores

Spring deformation model.

Score is linear in local templates w_i and spring parameters w_{ij} .

$$\text{Score}(x, z) = w \cdot \phi(x, z)$$

$\{ z \rightarrow \text{position of the sub-parts in the image}\}$

{ 1st part \rightarrow where / what all sub-parts are present in the image
and part \rightarrow relative relationships between the sub-parts of the whole object / image.

{ considers both the things together! }

~~CNN has the part \Rightarrow [part template scores]~~

why not do this?

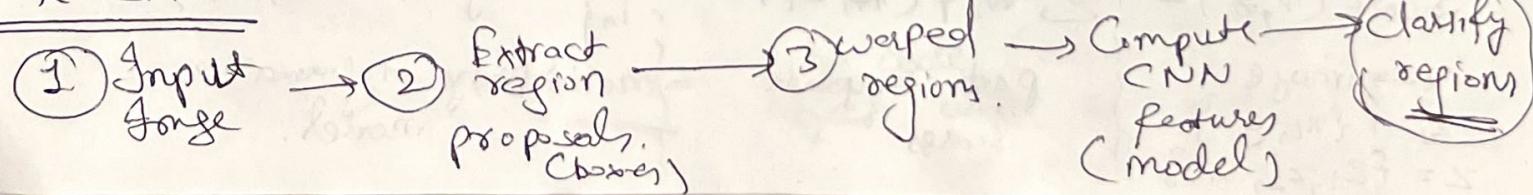
\hookrightarrow (doesn't model) \rightarrow spring deformation model.

* Side note :- we can transform images / objects into its FFTs and then store those beforehand, basically FFTs convert space & time so, that might cover both part template scores & spring deformation part. Then, training a DNN on these wave data with corresponding labels can help in classification. [or might even help in segmentation].

Localizing objects:-

- Scanning window approach & image pyramids
- select many windows over the complete image.
(candidate bounding boxes)
- \hookrightarrow run classification on each window.

R-CNN:-



→ Learn for each of the windows and combine results. (18)

↳ How do we know that learning doesn't directly mirror the training data).

⇒ If we make the structure end-to-end like in fast R-CNN, and faster R-CNN.



(max pooling is done on these windows bcz initial dimensions are different)

{ We are able to resize the feature vector rather than the image, that's why it's faster, so,

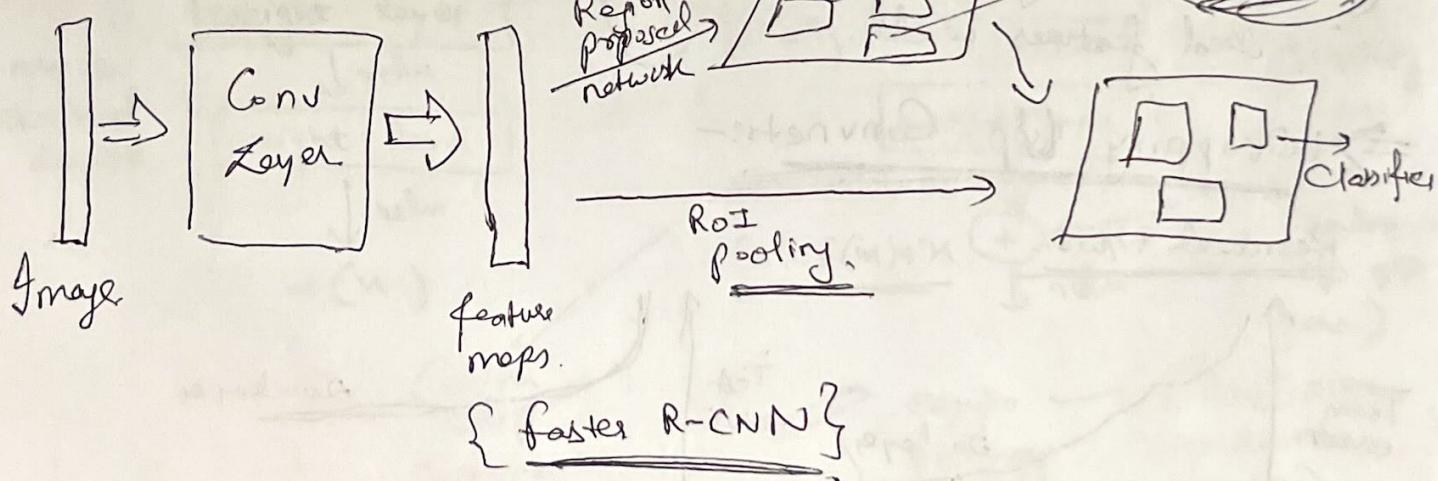
that we can reuse the computation }.

{ This makes it faster than R-CNN }.

{ We combine finding the boxes with classification of objects in the images }.

* this also saves resources.

{ To make it follow similar dimension, so, that we see feature of similar size / dimensions }.



Segmentation :-

① Semantic segmentation:-

Color code (One-hot) an object with some code and background as another.

Instance segmentation: color code different instances of an object with diff. grades of

→ the similar color and

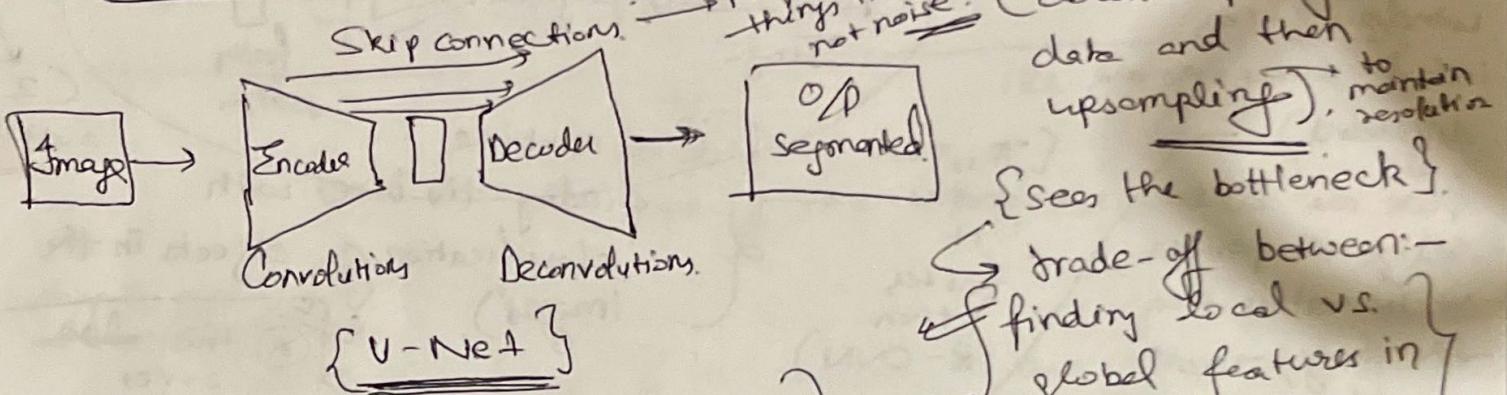
similarly whole background with one color.

Challenge:-
{unbounded number of }
O/P instances.

(like in semantic
one).

* fully Convolutional N/w. Call conv + pooling layers)

Encoder-decoder architectures:

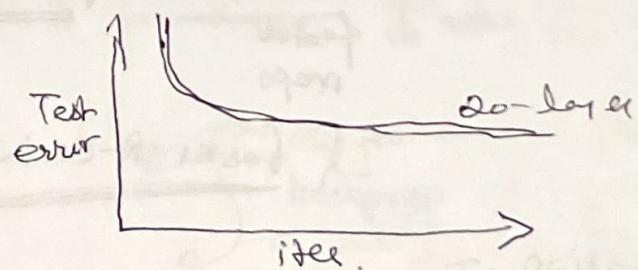
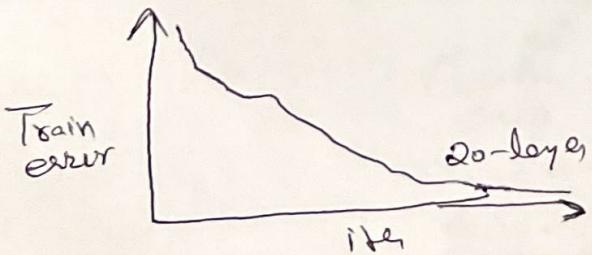


* If you never resize, then there is no way to get long range connections. (Global features)

* If you go higher in resolution, you are going to consider local features well.

⇒ Wrapping Up Convnets:-

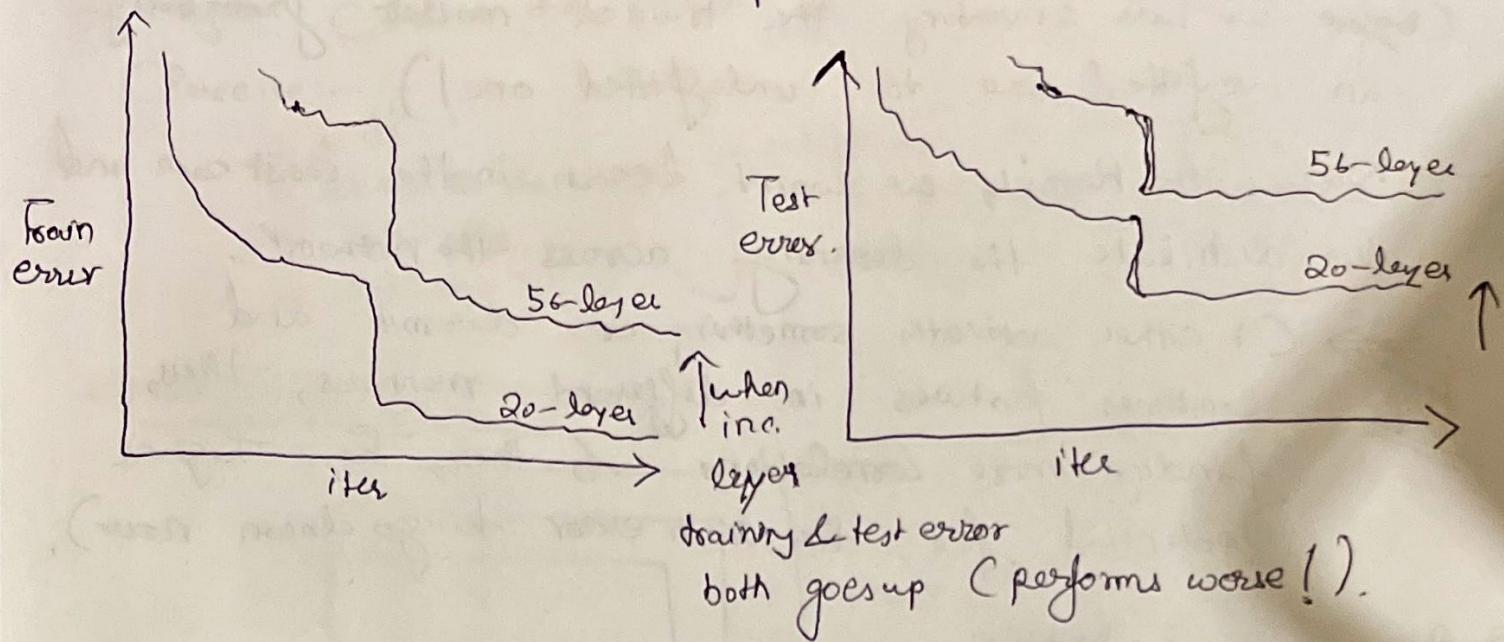
Residual N/w's:



* What should happen if I train a deeper N/w?

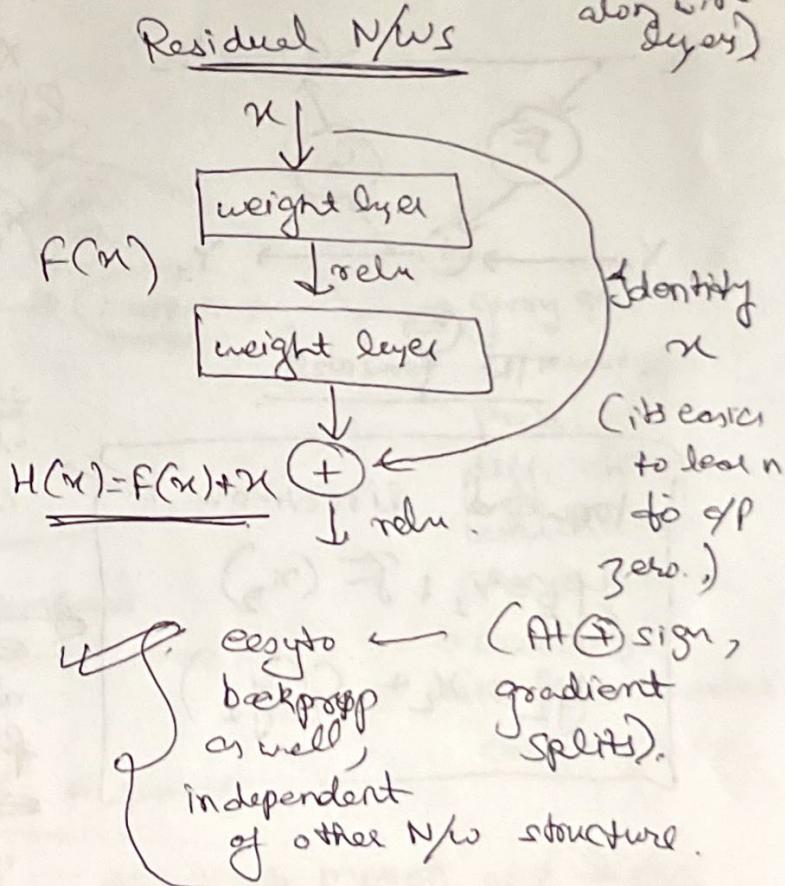
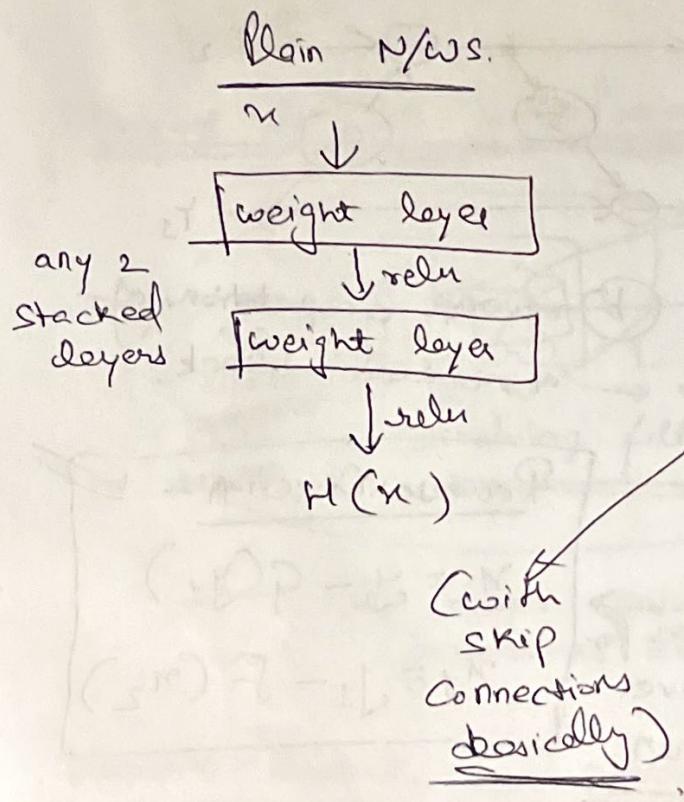
→ Train error :- As this is much lower than test error, this model overfits. We can also see that test error plateaus out.

(19)



* Why it goes up?

(you should learn to do nothing in case it would go up because this shouldn't happen). (we need to indicate along w/ layers)



→ Now, with ResNets, when we increase the number of layers, the train error would atleast go down now.

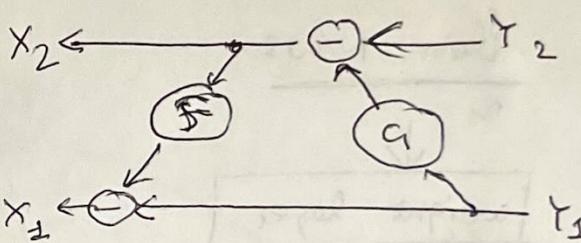
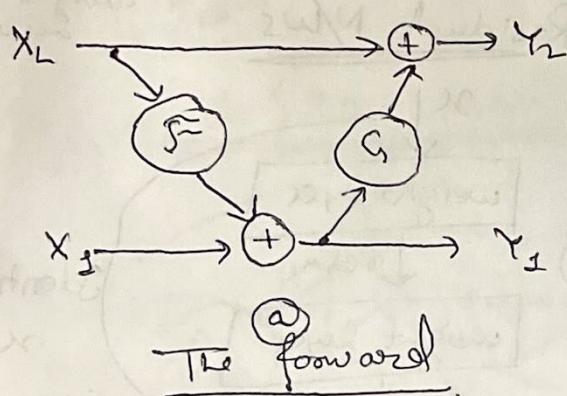
(Before we were converting the trained model from being an overfitted one to underfitted one!).

→ Now, with Resnet, we won't learn in the worst case and also distribute the learning across the network.

~~if~~ (It either activates something or doesn't and combines features in different manners, thus, finding more correlations to them, has higher potential for the train error to go down now).

②) Reversible ResNets

The backprop goes from $O(n)$ to $O(1)$ if n were nodes in the path from start to end.



b) reverse computation of a residual block

forward direction:-

$$y_1 = x_1 + F(x_2)$$

$$y_2 = x_2 + G(y_1)$$

epn for
above
flows,

Reverse Direction:-

$$y_2 = y_1 - G(y_1)$$

$$x_2 = y_1 - F(x_2)$$

Video Recognition:-

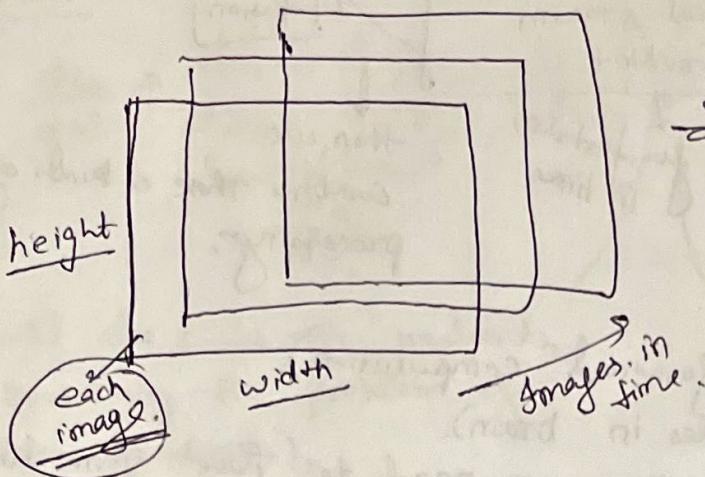
→ How to derive insights from just the video or sequence of images? (That's the problem for the machines!)

(Perceive:-

1. goals
2. objects in the video.
3. sub-tasks.
4. emotions).

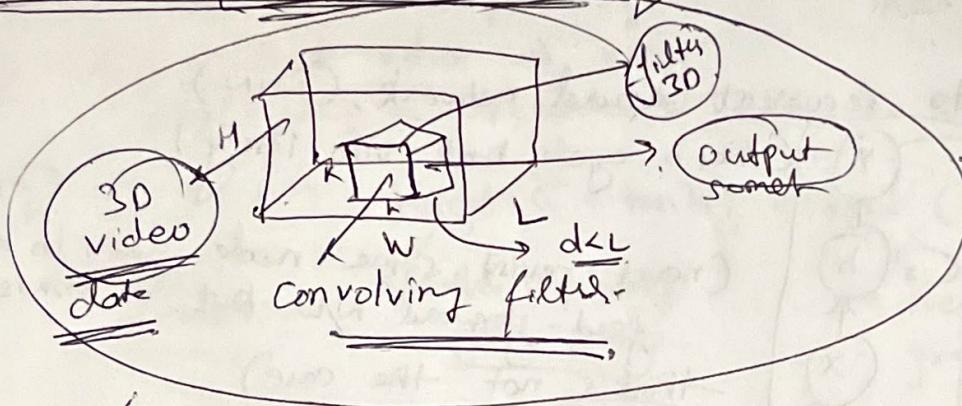
} difficult to recognize this in temporal form from a video.

⇒ Representing video:-



→ frames played fast enough creates an illusion of motion!

Now, we do {3D convolutions} \Rightarrow



→ giving %Rs in terms of prob. of diff. classes which can change in time according to the frames under consideration.

→ This learns spatio-temporal features as 3D convolutions.

Motivation:- Separate visual pathways in nature.

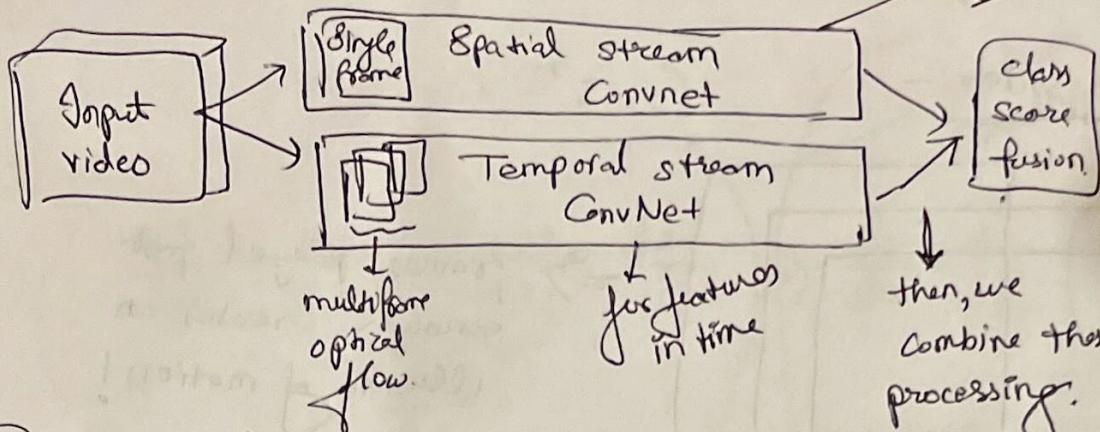
→ Dorsal stream:- "where/flow" → recognizes motion and located objects.

→ Ventral ("what") stream performs object recognition.
 (Interconnection b/w these streams can help process them together).

④ CV is good at ventral stream part but not for dorsal stream part.

* for action recognition,

we could build a 2-stream N/W, for spatial features



then, we combine those 2 kinds of processing.

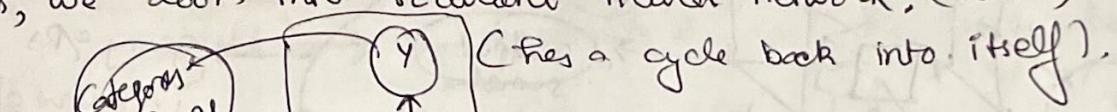
Recurrent processing:-

(has to be feed-forward computation).

(can't have cycles in brain).

but to incorporate time, we need to have some kind of looping.

Thus, we look into recurrent neural network, (RNN)

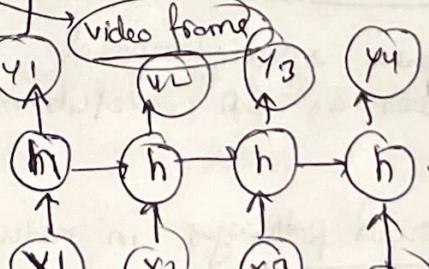


(ties a cycle back into itself).

(never revisit same node again in a feed-forward N/W but in RNNs, that's not the case).

④ this can represent any computational function!

data changes, opⁿ is the same.



$$f_{hi} = f(W_{Rn}^T x_i + W_{RR}^T h_i)$$

$$y_i = f(W_y^T h_i)$$

(opened up in time)

→ It is basically a memory being propagated through time! (21)

→ with RNNs, sequences are of variable length : they do some computation for no. of steps (possibly). If for videos, we need to model the sequence of frames. (temporal).

Backpropagation through time:-

We also do backprop for the hidden layers in time.

Forward pass:-

$$z_i = f(w_x^T x_i + w^T z_{i-1})$$

Gradients:-

$$\frac{dL}{dw} = \frac{dL}{dz_{i+1}} \cdot \underbrace{\frac{dz_{i+1}}{dz_i}}_{\frac{dL}{dz_i}} \cdot \frac{dz_i}{dw}$$

gradients are multiplied (multiplied with itself) (for all z_i 's).

$$\frac{dL}{dw} = \frac{dL}{dz_T} \left(\prod_{j=i}^{T-1} \frac{dz_{j+1}}{dz_j} \right) \cdot \frac{dz_i}{dw}$$

very prominent in nature,
if gradients get very small/
very big.

problem of vanishing/
exploding gradient comes
into picture.

→ we could also scale the gradients,
(for countering the problem of
exploding gradients),

But still problem of vanishing
gradients could come into picture!

If gradient value is very
small) → as recurrence can be very large.

ZSTM \Rightarrow Input (σ), forget ($\sigma + \tanh h$), output ($\sigma \times \tanh h$) gates.

(Simple activation)

Resolves exploding grad. problem

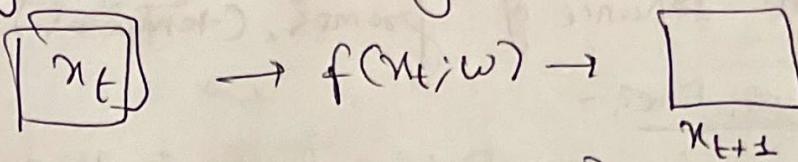
Resolves vanishing grad. problem

→ Speed of action alters perceptual judgement, so, if we slow down frames in a video and train a CNN on it, the perception can change. (We could also model this slow motion).

→ We always see the video according to the action (perception).

→ we don't pay attention to other features.

future generation → Try to predict what's gonna happen next depending on the actions before.



Why can't we do this?

Minimize Euclidean distance →

$$\min_w \sum_i \| f(x_t^i, w) - x_{t+1}^i \|_2^2$$

(in case of multiple modes (actions with some perks), it's hard to predict, averaging out actions doesn't help).

→ we might know what's gonna happen but how?

(it's very hard problem!).

What a machine could anticipate?

for example:

— A child stand outside a candy store, reaching for the door. What does the child predict will happen once they go inside?

→ The child does not predict the color of lollipops inside!
→ Rather, the child predict their own sensations of what they will see, hear and taste.

this
is
what a
M/C would
predict

→ this is what we had want to predict

(we might need multiple predictions based on perception).

→ (multiple futures possible).

← (but one of them should match!)

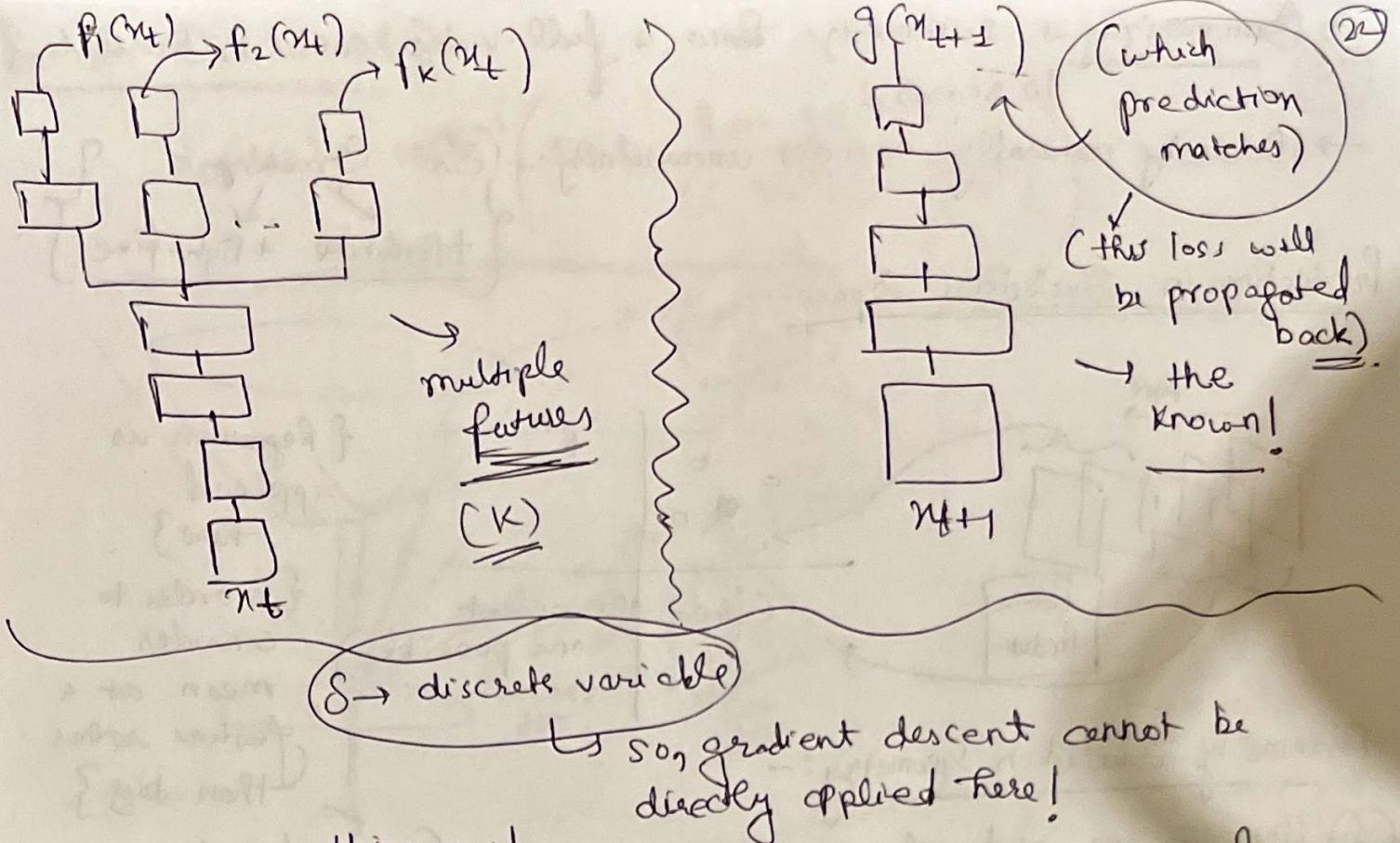
latent variable

If we
did know
this, then
we could
incorporate
that in
loss
function

$$\min_{f, g} \sum_i \sum_k (\delta_k^i) \| f_k(x_t^i) - g(x_{t+1}^i) \|_2^2$$

$$\text{s.t. } \delta^i \in \{0, 1\}^K \text{ and } \|\delta^i\|_1 = 1 \quad \forall i$$

↑ allow only one ↑ residual from Cl latent)



we use something new!

Expectation-Maximization

- Guess the latent variable
- Try to solve the loss and back-propagate.
- Repeat the above steps.

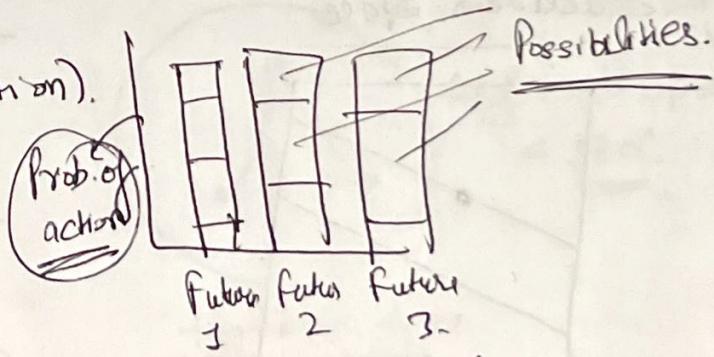
(holding 4/S
at a
time)

- 1. E-Step :— fill in missing variables and estimate latent variables
- 2. M-Step :— fit the model (with back-propagation in this case).
- 3. Repeat.

↓ we compute probability of action).

(we need to decide time frame too).

→ we could also turn this into a feedback solution spec

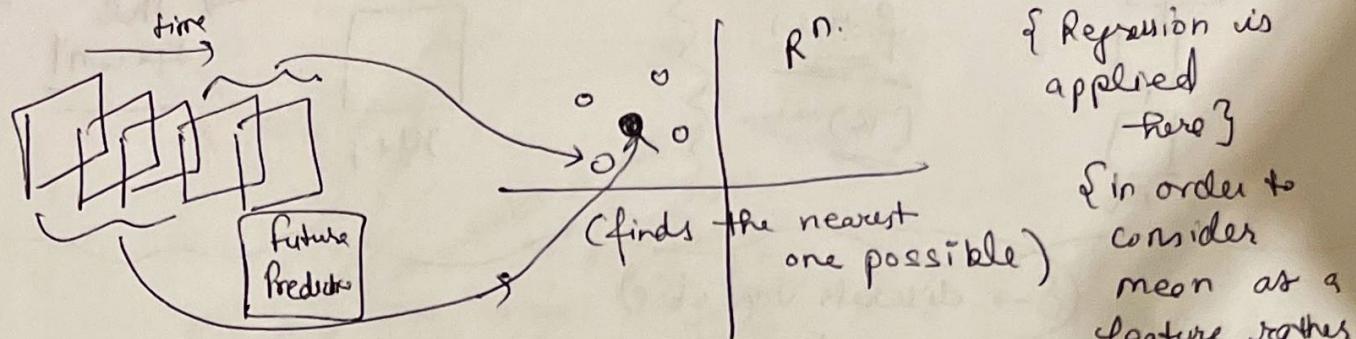


spec

→ Summarizing → summarize from a full-video could be translated to scores.

→ hierarchy naturally encodes uncertainty:- {Ex:- Greeting
+ Handshake + High-five}

Predicting in Euclidean Space:-

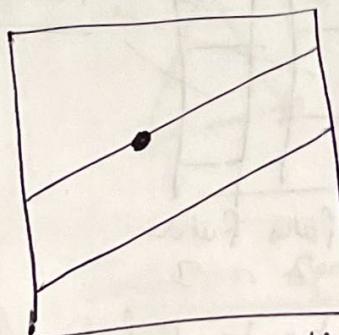


Axioms of Euclidean geometry:-

- ① There is one and only one line segment between any two given points.
- ② Any line segment can be extended continuously to a line.
- ③ There is one & only 1 circle with any given center and radius.
- ④ All right \angle s are congruent to one another.
- ⑤ Given any straight line and a point ^{not} on it, there exists one and only one straight line which passes through that point & never intersects the 1st line. (parallel lines).

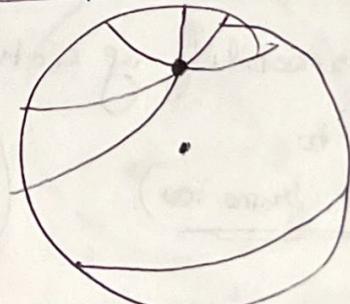
{infinitely many} → for hyperbolic geometry

Euclidean Space



{Only one parallel line}

Hyperbolic Space

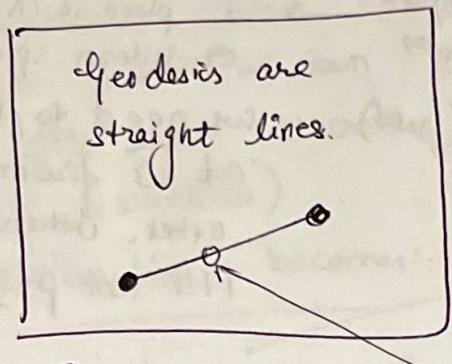
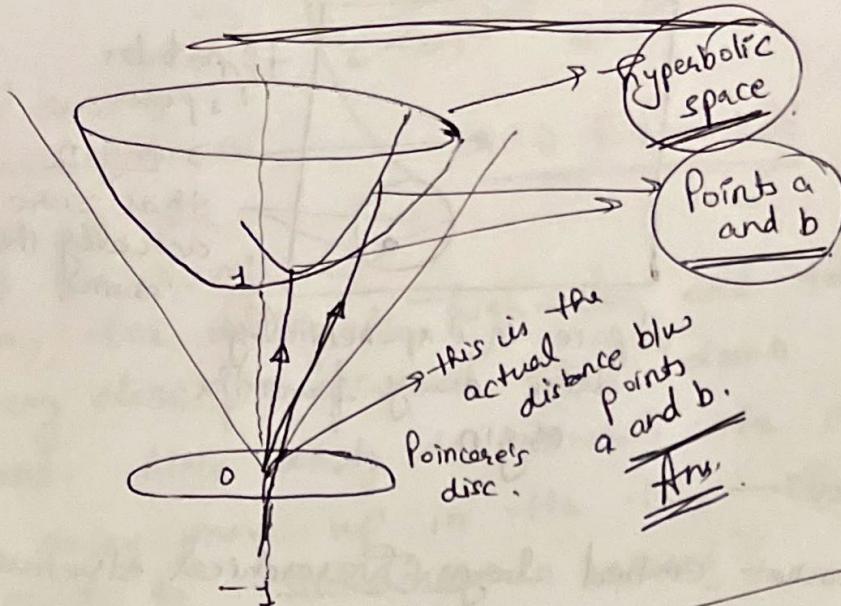


{Any nix parallel lines.}

Poincaré Disk Model \Rightarrow

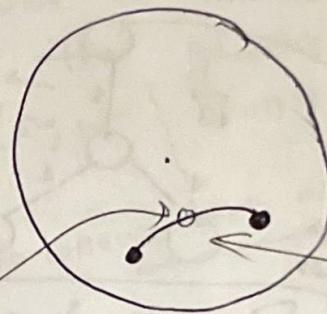
(23)

$$d_D(a, b) = \cosh^{-1} \left(1 + 2 \cdot \frac{\|a - b\|^2}{(1 - \|a\|^2)(1 - \|b\|^2)} \right)$$



Euclidean space

Mean points
(half-way)

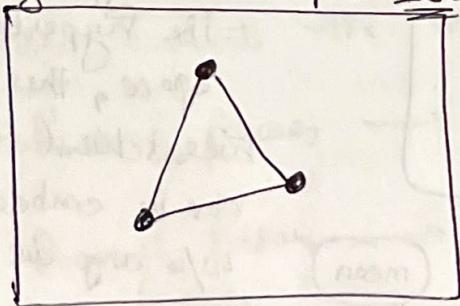


Hyperbolic space

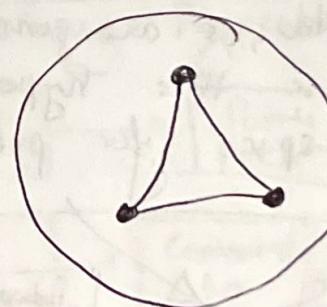
Geodesics don't "look" straight.

Any mean pt in this space is actually closest point to the origin)

Angles must add up to $\leq 180^\circ$.



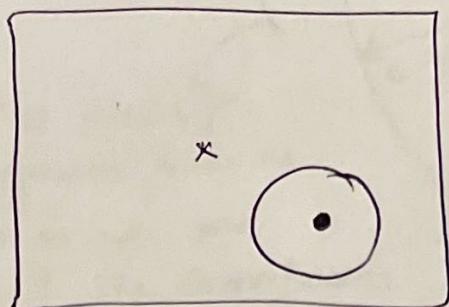
Euclidean space



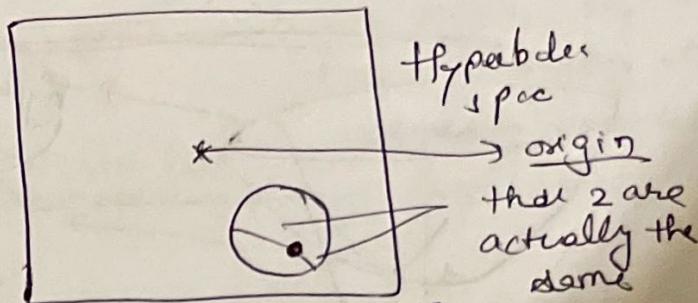
Hyperbolic space

Angles add upto $< 180^\circ$.

→ All angles are $< 90^\circ$ in a hyperbolic space ∵ the shapes such as squares & rectangles aren't formed in the hyperbolic space.



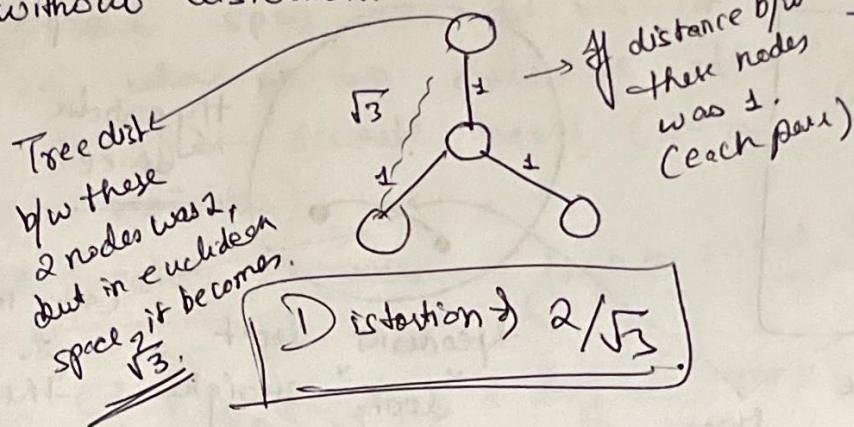
Euclidean space
Standard circle



Space is exponentially dense away from the origin.

⇒ Distortion of the space:-

→ Euclidean space cannot embed large hierarchical structures without distortion.

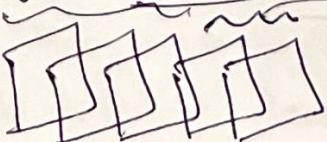


To place it in Euclidean space, we need to place it at $\sqrt{3}$ from each other. Otherwise, it's not possible.

Hyperbolic space naturally embeds trees (because of the inherent density / distance in the hyperbolic space, this hierarchical structure can be embedded w/o any distortion).

If now, we are gonna use this hyperbolic space for predictions

Hedging the bet



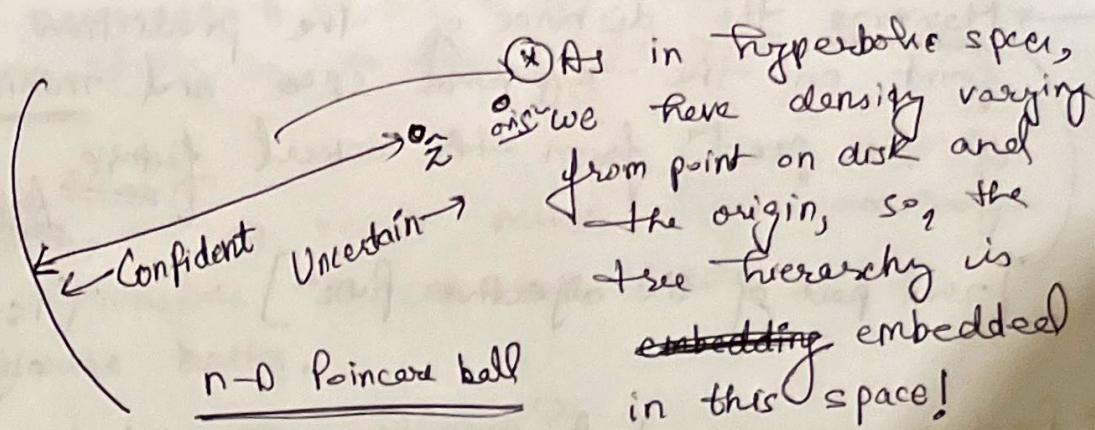
future prediction

mean

• (origin)
n-D Poincaré ball

the parent of all pts.

24



→ the density difference tells us that as we are around the circumference boundary, the density is high there and to predict a point out of very close pts, we need to be damn sure, thus, more confident. Since density dec towards the origin, pred^m becomes easy as we move up in the tree hierarchy. (similar to hierarchies in nature!).

→ hyperbolic space can be viewed as a continuous form of the tree where we can calculate the gradients easily now. It's not discrete now! (so, we won't need the latent variables).

So, loss function becomes:- $\xrightarrow{\text{func to convert that distance}}$

$$\min_i \sum_i \left[d_h^2(\hat{z}_i, z_i) + \log \sum_j \exp(-d_h^2(\hat{z}_i, z_j)) \right]$$

Obj func

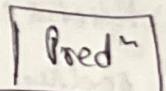
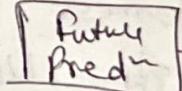
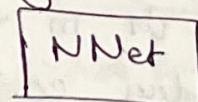
future prediction

actual future.

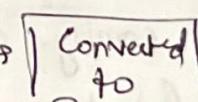
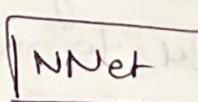
Neural nets converts the images in the hyperbolic space

(Both the models are trained on the same objective func.).

feature frames



Actual future frames



Compared to compute the loss

{The hierarchical structure of the tree is learned over time}.

→ Maximize the distance of the prediction from any random point on the hyperbolic space and minimize the distance of the predⁿ from the actual future.

[2nd part of the objective funcⁿ]

[1st part of the objective funcⁿ].

⑧ Features are encoded in the "hyperbolic space as well".
↳ the scenes in the images are stored in the hyperbolic tree structure as well (and thus, as tree hierarchies which are learnt by the NNet).

The z's are actually vector representations of the frames/images and can be treated as an embedding for the training.
(We sum over all the samples in our tree to find the loss from the objective function).

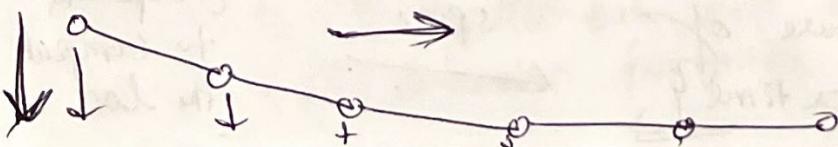
(Basically we predict the mean of all the points corresponding to the frames in the video)

→ Here, we design an objective function!

→ We need to create a hyperbolic classifier to train the N/w for the use-case.

Hierarchy provides us the notion of uncertainty.

↳ because based on the input frames, we can forecast a particular category at a particular level in the hierarchical tree structure formed which is in turn based on the form of ball radius. As this radius goes low, that means that our prediction is coming closer to the ground truth in the hierarchical tree structure.



point
ball radius

Note:- As we would never reach as where the leaf nodes in the hyperbolic space should be, so, the loss might never become zero.

(25)

→ In hyperbolic space, we can reach much higher accuracy in less number of dimensions as we can model the hierarchical structure better.

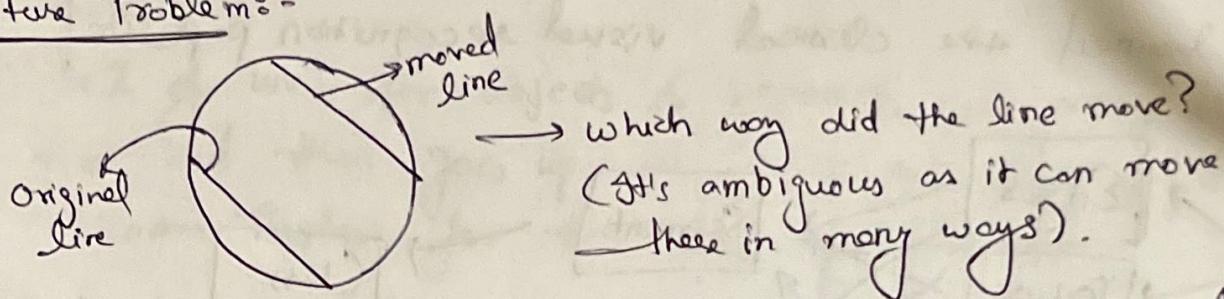
→ How to replicate the motion stream?

Optical flow:-

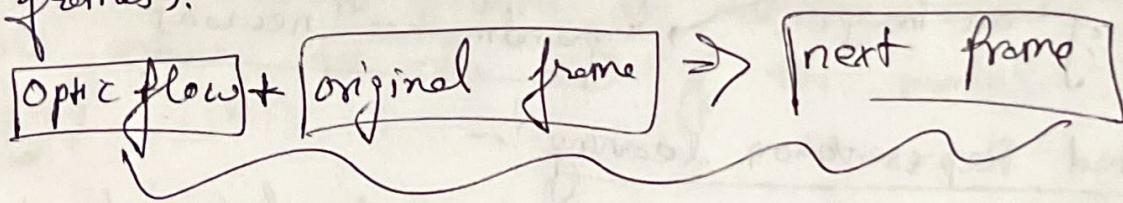
flow field → assign a flow vector to each pixel. → to suggest the flow in a scene (set of frames).

Visualize :- flow magnitude as saturation, orientation as hue.

Aperture Problem:-



To find the optic flow, motion magnification can be used. Then, the optic flow map of the image/video can be used and trained with to make predictions of the flow. (Or learn to extract such frames).



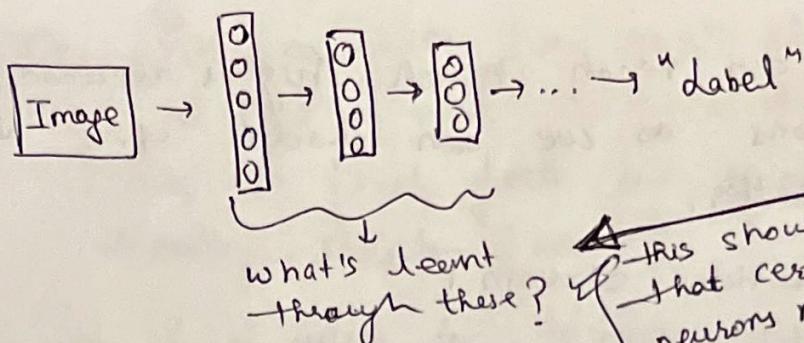
(This kind of prediction can be done quite easily)

Possible Qns:-

1. filter Convolutions mapping
2. hyperbolic space.
3. a little bit about backprop.
4. Activation functions.

Interpretability :-

→ what do deep nets internally learn?



"Grandmother" neurons:-

→ Always one neuron is activated.

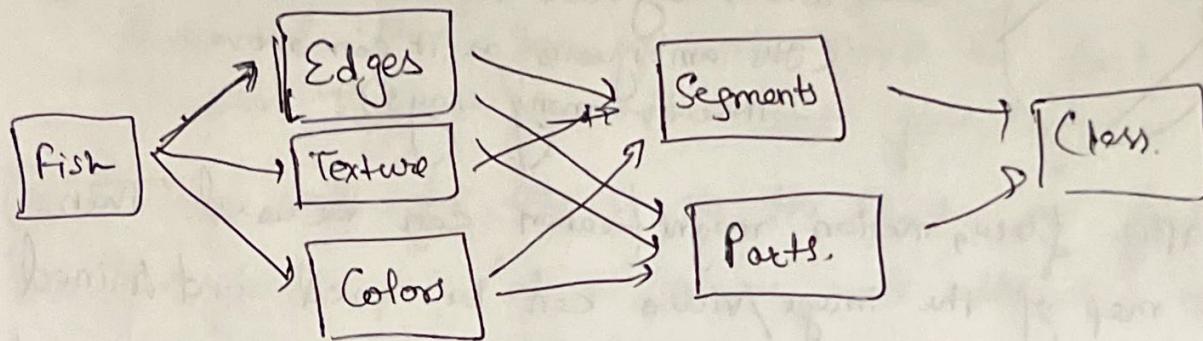
(even one neuron might be activated for a particular class).

↳ This shows that certain neurons might have some info. about particular label/type of activity!

⇒ Generally, some set of neurons are associated with a class

rather than just one, this makes memory more resilient! Such might be done in neural networks as well.

⇒ CNNs learned the classical visual recognition pipeline!



→ Different neurons are activated differently, to detect particular features of an image. → "grandmother" neurons.

Distributed Representation learning :-

→ In this we distribute the learning for each of the neurons on different systems as each neuron is activated for some particular type of feature.

{Image patches activate some particular set of neurons most strongly}.

Object detectors Emergence in Deep scene CNNs

(26)

→ As we go deeper in the model:-

- ⑩ The percent of units for simple elements & colors is the most at first and it goes down constantly.
(bcz at first few layers are responsible for detecting simple features).
- ⑩ The % of units for texture materials is high at first and then goes down slowly.
- ⑩ The % of units for regions/surfaces is the least at first and then goes up by a lot (at a higher rate).
- ⑩ The % of units for objects & segments is very less at first and then goes up at quite a fast pace and becomes highest for the last few layers.
- ⑩ The % of units for scenes is the least at first and then goes up (at a slower rate though) in the last few layers.

Thus, model becomes smarter at detecting coarse details while moving to the end of the model (going deeper).

Investigating a representation via similarity analysis:-

for all the image/subparts modeled in the model, we form a representation and dissimilarity matrix to find how dissimilar these images are! Activation vectors are formed from the

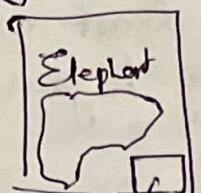
images, that is used to find the dissimilarity b/w those 2 vectors & thus the images. (L1 norm).

$$\|\mathbf{h}_i - \mathbf{h}_j\|$$

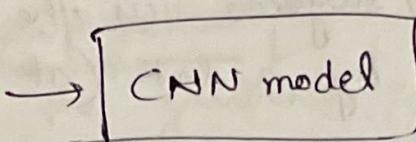
→ Neural activation vector

We can form the matrix for all the images/patterns learnt by a model and generate some kind of heat map to see what all features are being extracted & how many!

Saliency by Occlusion:-



we mask
some part
of the image



$$p(\text{elephant}) = 0.98$$

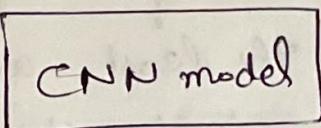
When we mask some part of the image & find op prob. for that class drops, then that masked part was significant!



masking
some part
of the image

thus score

head of elephant
is marked



$$p(\text{elephant}) = 0.78$$

By occluding some parts of the image, we can find which part of the image helps the classification!

After doing the above and moving mask over the full image, we can generate a heat map and see how the change in score happens; higher the change in score, more imp. that aspect of the image and higher the heat & thus these heat maps gives us the imp. part of the image for that class.]

(27)

Saliency maps Reveal Bias

By occluding parts of the image by increasing the mesh size, we could find the point till when we mask, the image is classified correctly. For example:- Husky is classified as one till the whole Husky is masked with some minor part left and snow. This happens because mostly with snow, Huskies are associated. Thus reveals the bias in the model on the surrounding of the obj. under consideration. ~~fmp~~

Guided back propagation

- we compute the gradient of hidden unit w.r.t the input pixels, this gives us the change w.r.t the input. Basically, if an input pixel is significant to classify the image correctly it gonna have a +ve gradient while backprop. If not significant, then going to have -ve gradient.
- Also, we only propagate the +ve gradient backwards & truncate -ve gradients to zero.
- Intuition) find which regions cause the object, not the regions that do not cause it!
- ⇒ for maximally activated patches, we find op of the neurons s.t. only the most significant parts are highlighted after using guided backprop.

This helps us in interpretability by a lot!

Per Gradient Ascent

This optimization problem is non-convex and not guaranteed

to respect natural image priors.

This way we tell the model to give us the best representation of a class and it provides us with something ~~with~~ from which we can see the boxes associated and also the imp. features for that class (or atleast what the model thinks is!).

$$\max_{\mathbf{u}} \mathcal{L}(\mathbf{u}) - \lambda \|\mathbf{u}\|_2^2 \quad \xrightarrow{\text{w.r.t input}}$$

⇒ Pick a unit, perturb the activation to go up rather than go down in the obj. function.

e.g. we saturate a particular type of neuron for the images to find most imp. features for that class.
(We want to turn on rather than turn off some of the neurons that we need for a class).

* This is done on a trained model. In fact, all the interpretability is done on trained models.

Flows— (Related to motion flow)

How you track an object in time?

{Dynamics in video} → how to harness w/o supervision?

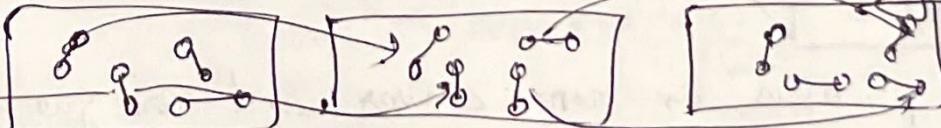
Common factor) what connected pieces are going to move together in a video?

Correspondence ⇒ how an object/person is the same one across diff frames in the video.

We need to find these!

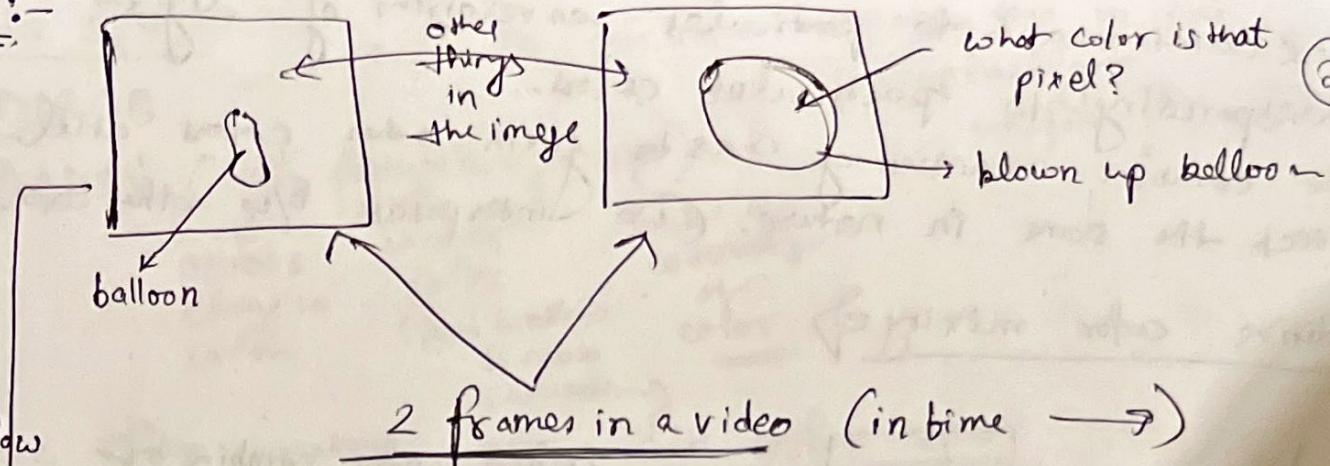
Object tracking) → problem of recognition and consider above 2 things.

Common feature



Correspondence

Ex:-



How
does
my
c
tell
this?
(about
color)

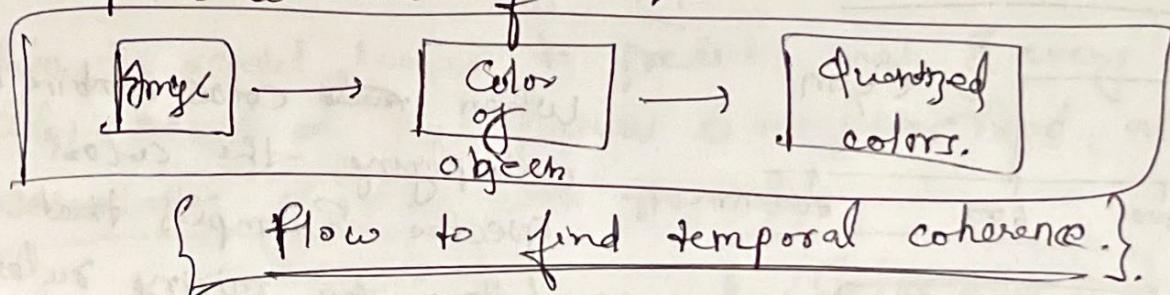
① It needs the track that the objects are the same in time (correspondence).

② How does the object move across time (in terms of frames) (common fate).

We can learn to track an object in a frame across the full video and maybe color that object in the gray-scale video (moving racing car).

With some exceptions, we can find some coherence in colors for diff. objects in the frames of the video.

For achieving the objective of tracking, we need to find the temporal coherence of color.

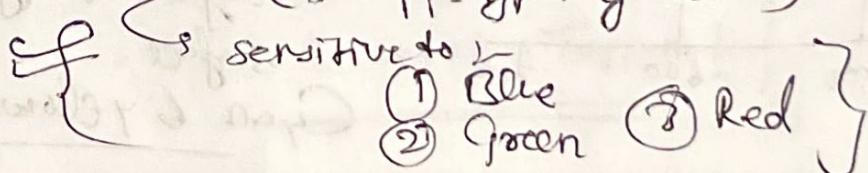


About Human eye:-

Retina has:-

- ① Rods → responsible for determining intensity of light, doesn't care about colors.
- ② Cones → sensitive to colors,

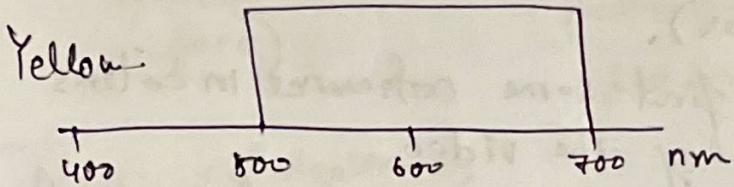
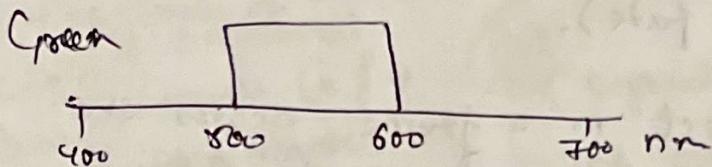
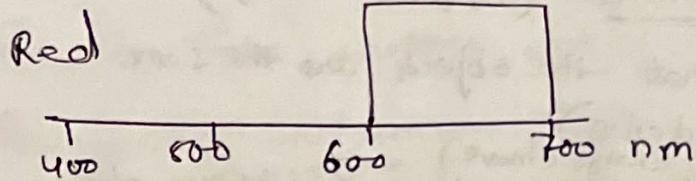
(3 diff. types of cones)



These cones react to particular wavelengths of light corresponding to particular colors.

→ we combine reaction of cones to form other colors and detect the same in nature. (we interpolate b/w other colors).

Additive color mixing

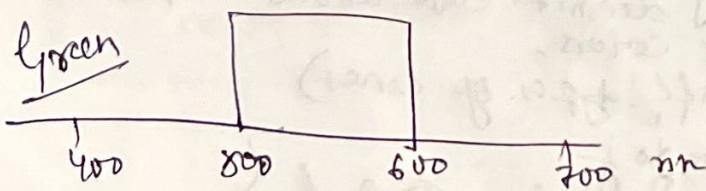
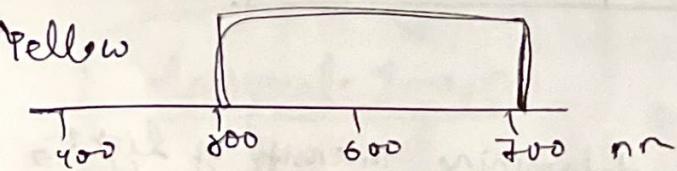
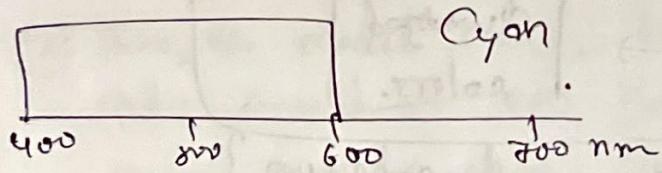


When colors combine by adding the color spectra.

Ex:- color displays that follow this mixing rule:-
tiny display dots on a monitor screen, multiple projectors aimed at a screen.

R + G make Yellow!

Subtractive color mixing



When ~~the~~ colors combine by multiplying the color spectra. Examples that follow this mixing rule:- most photographic films, paint, colored optical filters, crayons & light reflecting off a diffuse surface.

Cyan & yellow → Green

Color spaces

(in euclidean space this is not perceptually uniform!) (29)

① RGB

② HSV (Hue, Saturation, Value)

↓
across
change in
Colors

co-ordin
color
(from white
to max
saturation)

↓
value of
intensity
of
color.

So, we have another one:- Lab color space

Perceptually uniform:- Euclidean distance matches human perception of color similarity.

L → brightness of color.
(intensity)

a → 2 channels (x & y)

b → co-ordinates
covering a color
space)

(30)
separates the color dimensions from
the intensity dimension making it
perceptually uniform, -then is used a
lot now!)

Image colorization

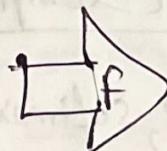
Training data $\rightarrow \{x, y\}$ where { $x \rightarrow$ grayscale image }
 $y \rightarrow$ colored image }.

* Then, the model learns to predict each & every pixel color based on some features & neighborhood and the correlation of those with the colors.

Grayscale Image

L channel

$$x \in \mathbb{R}^{H \times W \times 1}$$



Color info:-

ab channels.

$$y \in \mathbb{R}^{H \times W \times 2}$$

{ This is what happens! }

Choosing loss of representation

If we choose a loss fn such as:-

$$L(f(x), y) =$$

$$\|f(x) - y\|_2^2$$

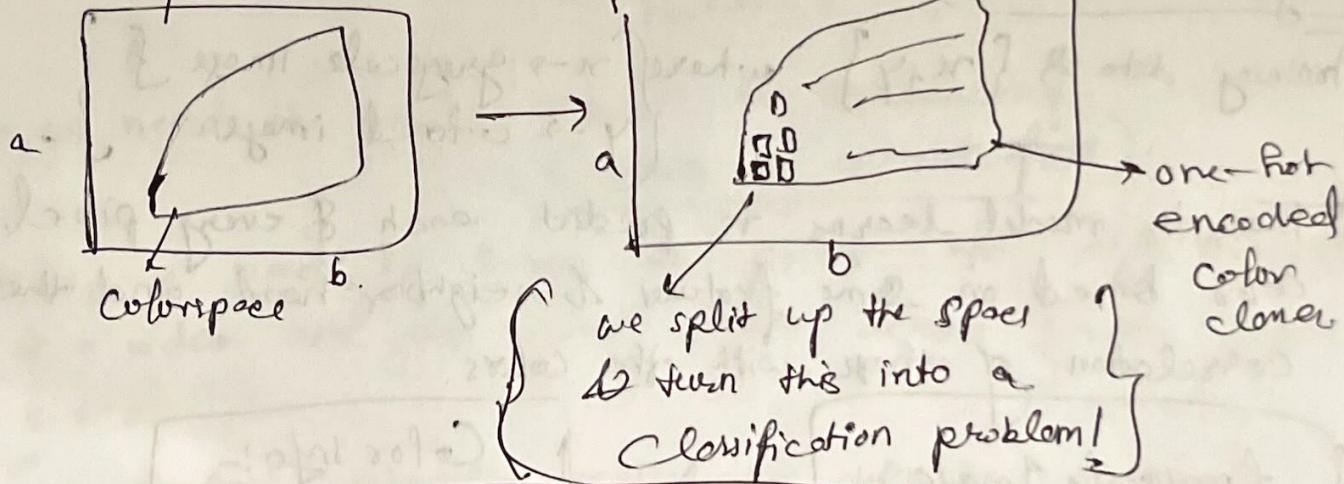
This is going to train the NW

s.t. it will predict the averaged color of an object/pixel based on the learnt associated color with an object or object + surroundings.

(taking avg doesn't help us here much!)

One-way to convert this regression problem would be to one-hot encode k discrete classes representing k diff colors in the lab color space. (ab channels). This converts the problem to a classification one and thus, moves away from calc averages which is involved in a regression problem

$$y \in \mathbb{R}^{H \times W \times 2}$$



Loss funcn

$$L(y, f_\theta(x)) = H(y, \text{softmax}(f_\theta(x)))$$

{This gives us more vivid colors sometimes!}

It might still give diff colors but learns to detect object correctly in the image.

Instructive failures

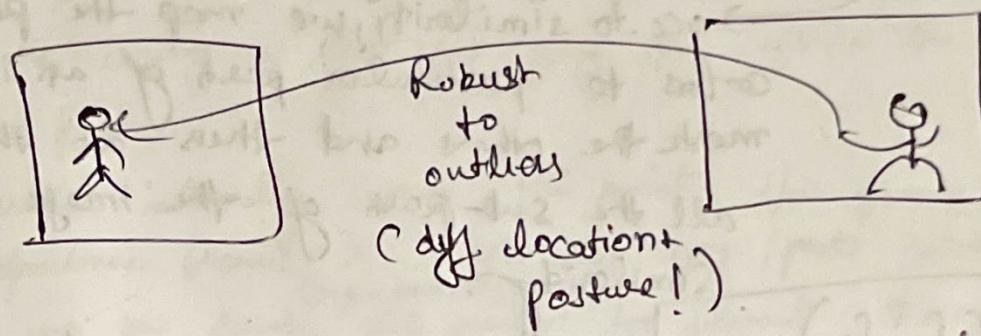
(10)

Could lead to wrong coloring as well for ex: if we had lots of images in our training data with dogs with tongues out, then it could color a dog with no tongue out as still ~~the~~ the tongue color. (corony, right?)

Then, we go back to the problem of object tracking where we need to know to copy colors & detect object correctly first. (for this, we need to learn to do correspondence).

→ This should be robust to outliers as well!

(Learning a policy how we should copy/paste colors).



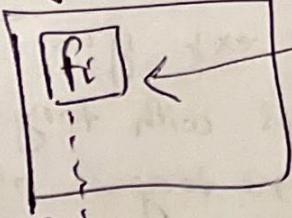
Also, we need it to be robust to occlusions
(driven by common fate)

→ So, we need to learn to point to diff. things in a video and for some object across frames in the video.

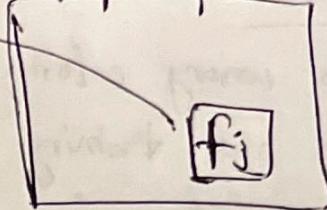
(we construct embedding vectors for different features in different frames and construct a similarity matrix
(similarities A_{ij}): b/w patches of images.

(Also, we would convert these to color space and then mapping the same heat map for these).

Reference frame

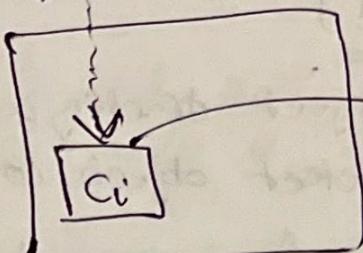


Input frame

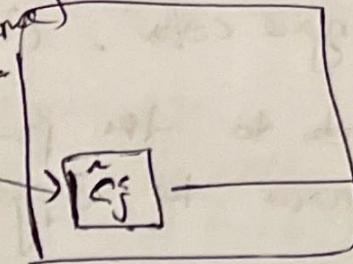


A_{ij}

&
Need to
find
the
correspondence



Reference colors



Target colors.

(Predict
color for
each pixel)
we form
a target
map
(color)

$$\hat{c}_j = \sum_i A_{ij} c_i$$

(Weighted sum of all
colors in the input (ref) image and
weighted by similarity matrix).

acc. to similarity, we map the particular
colors to particular parts of an image &
mask the others and then do this for
all the sub-parts of the image acc to
similarity.

$$A_{ij} = \frac{\exp(f_i^T f_j)}{\sum_k \exp(f_k^T f_i)}$$

Similar to softmax!

probs.

Ans

objective
func

$$\min_f L \left(\hat{c}_j, \sum_i A_{ij} c_i \right)$$

(Distance func should be able Ans).
to copy past colors, which is modeled by minimizing
the loss func above) \rightarrow (Applied to
video coloring)

(31)

CGI : Computer generated Imagery

Lumière brothers : Inventors of motion picture, 1895.

When they invented it, they thought it to be completely useless.

→ also invented first practical color camera, 1903.

Then sold it to a magician called Georges Méliès.

↳ discovered special effects

He accidentally discovered double exposure and used it in his magic tricks.

In the colorization problem, rather than us trying to explicitly train the machine, we're tricking the machine to solve it.

We are creating a task that going to require the machine to solve a problem.

The model internally learns this pointer (w) frames to learn correspondence because it needs to copy/paste colors.

Now, we can reuse ^(repurpose the model) this pointer for different tasks.

Instead of propagating colors, we are now propagating the segmentation mask of an object. → to track objects!

We can also reuse the same equation used before :-

$$\hat{c}_j = \sum_i A_{ij} c_i \text{ where } A_{ij} = \frac{\exp(f_i^T f_j)}{\sum_k \exp(f_k^T f_j)}$$

The meaning of c changes now; it no longer denotes color, but now denotes segmentation mask.

It can also work on complex issues such as occlusion since the model is trained on massive amounts of data.

Neural networks, in general, specialize in one particular task but this can generalize ~~to~~ the different types of tracking tasks that are not related to colorization.

We can also do pose tracking.

* Visualizing Embeddings (PCA/dimensionality reduction)

Project embedding to 3 dimension and visualize as RGB.

Objects ~~with~~ that are the same color are close in the embedding space. The neural nets are not told about the objects but they are able to ~~the~~ color an object with a different color — learns to distinguish / do grouping on its own. It makes mistakes.

calculate
euclidean
distance

Embedding are very ↑ in dimension & might be confusing most of the times

Self-supervised Learning

How do you solve all diff. problems with a model?

(It's a really difficult task).

- Rather than supervising the M/Cs, we want them to supervise themselves & learn recursively through the process.

Prediction hypothesis:- Our mind when sees something, when redrawn again, we would extrapolate some information not always seen in reality.

Learning without examples: This introduces noise such that we can learn without giving it what everything is!

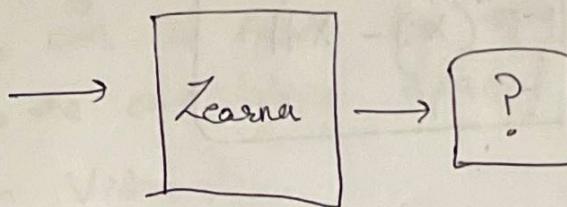
(Includes unsupervised learning and reinforcement learning).

Data

$\{x_1\}$

$\{x_2\}$

$\{x_3\}$



Unsupervised Representation learning:-

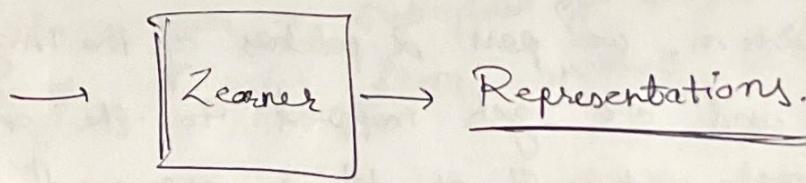
Learn features / representations of the world that help us learn patterns that can be applied elsewhere.

Data

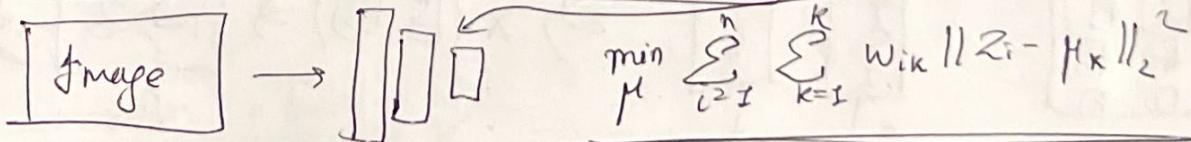
$\{x_1\}$

$\{x_2\}$

$\{x_3\}$

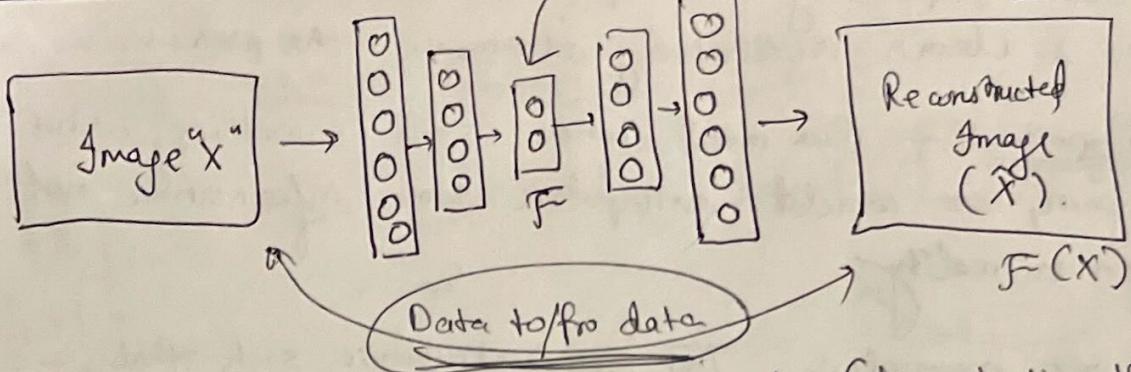


Unsupervised Learning by Clustering:- compressed image code (vector \approx)



$$\text{s.t. } \sum_{k=1}^K w_{ik} = 1 \Rightarrow w_{ik} \geq 0 \quad \underline{\text{k-Means Objective}}$$

"Auto encoder" \Rightarrow



\Rightarrow Essentially, it learns to memorize this! (by extracting the most important features, reducing the dimensionality by removing the correlated features).

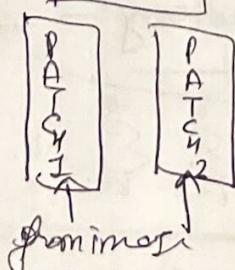
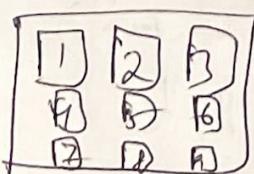
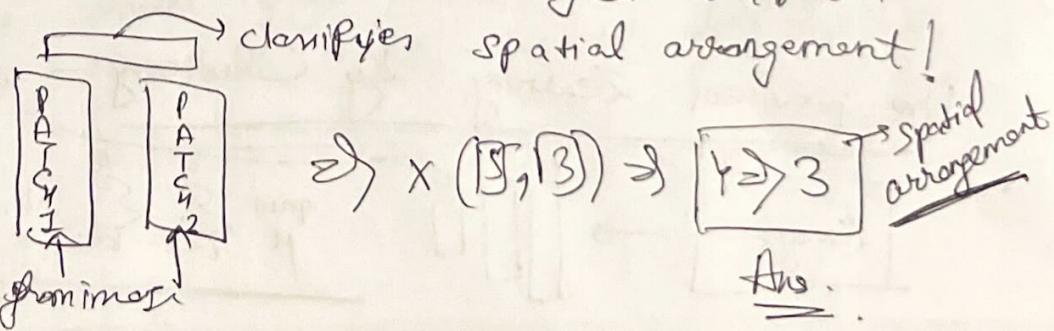
{ related to diff features of an object/image }

for which the expected class can be minimized.

for self-supervised, we extrapolate:-

\Rightarrow Split the data in 2 parts and predict one from another!
(This makes the model learn correlations in the image data).

for context prediction, we pass 2 patches of the image to the neural network and one gets mapped to the character as a classification task of where it should be arranged w.r.t. one another.



⑩ Even with random initializations of the patches, the self-supervised model learns the correlations b/w diff. patches/object in different images on its own. (33)

33

Learning by Counting

Learning by Counting \rightarrow
Learn recursive structures to images \rightarrow
Learning recursive

$$f(\boxed{\text{I}}) + f(\boxed{\text{L}}) + f(\boxed{\text{D}}) + f(\boxed{\text{U}}) = f(\boxed{\text{F}})$$

↓
 2 imp.
 features
 such as eyes

↓
 2 legs
 as
 features

↓
 was
 broken into
 4 quadrants

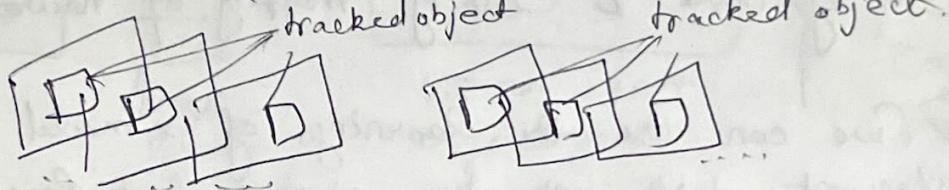
full image

Sum of total no. of eyes in the final image

eyes in the final image should be the same as the no. of eyes present in all the embeddings that's added up. This implicit count of features can be learnt by the model. (as close as possible to the embedding space!)

Learning from Video →

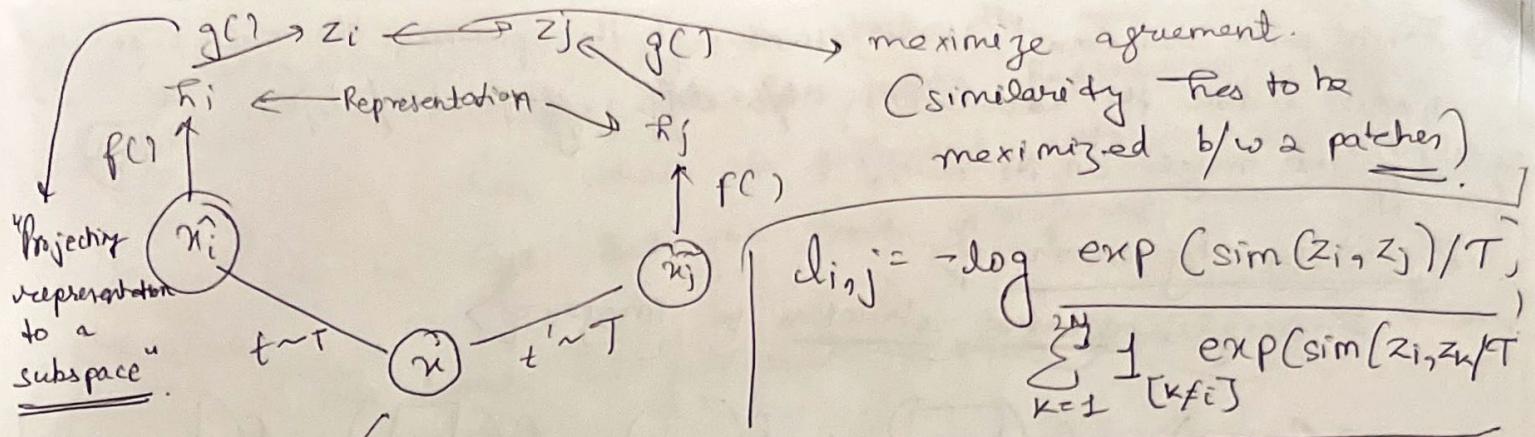
Unsupervised tracking in video



→ We rank the images or sub-part of images according to the distance of one patch from another in the embedding feature space.

Distance b/w similar images will be less, when compared with a random image, distance will be large.

Contrastive Learning



{ Diff. views are created out of an image }

to do contrastive learning, finding contrast of similarity b/w diff. views.

Sim \rightarrow Similarity metric

Ans

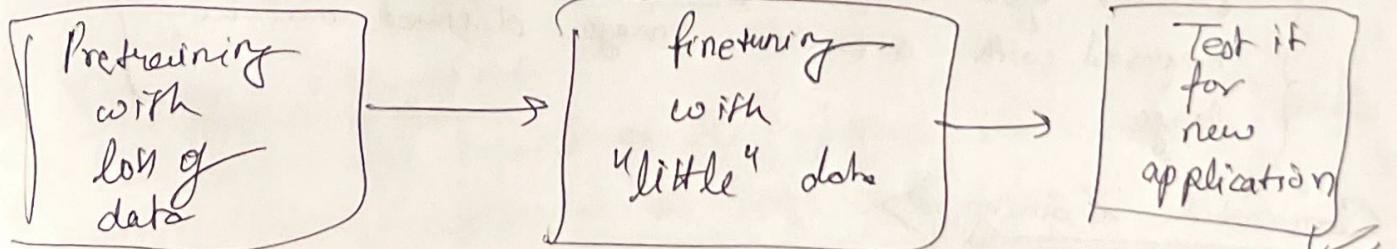
z_i, z_j are from some image

z_k \rightarrow view from a random image

clarify which 2 images belong to same image, whole metric that image as far as possible from a random image

Contrastive Learning needs a lot more parameters to reach the accuracy for a supervised model solving the same task.
(As need lot of data & lot of mapping of neurons in NN).

finetuning \rightarrow (we can use the learnings of a model to use it for a different task just by changing last few layers by retraining the last few layers (might have to change the loss funcⁿ) on the images for the specific application for which finetuning is being done).



* for contrastive learning, z vector stores the minimum required information for minimizing the loss funcⁿ. Representation layer has a lot of data, when that is projected into a subspace, a lot of info. is lost but the most significant ones are preserved! 34

(This learns the features in a much better manner than auto-encoders).

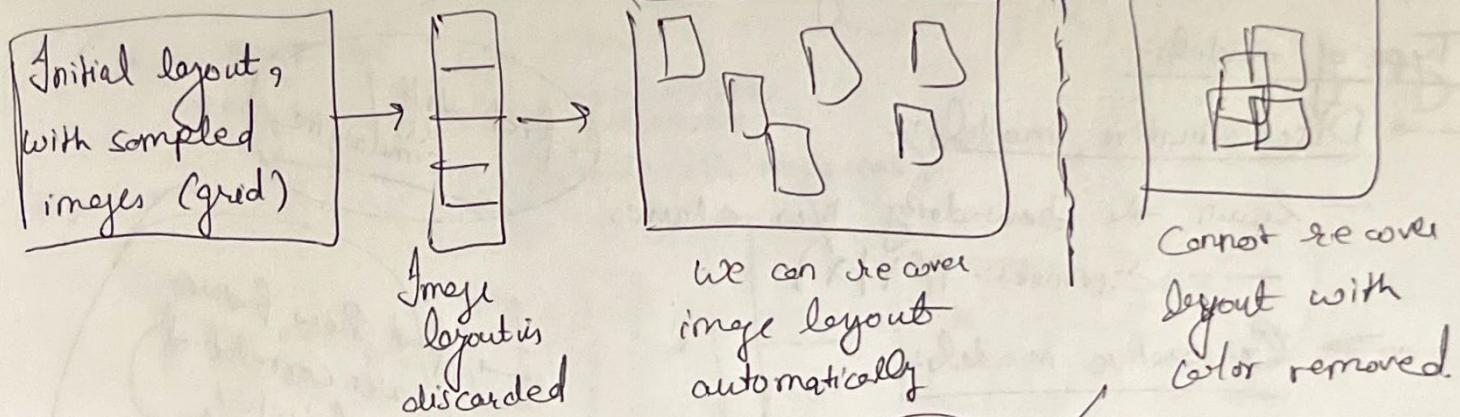
* Ways to generate the views of an image:-

- ① Crop & resize
- ② Flipped
- ③ Color distorted
- ④ Color distorted (jitter)
- ⑤ Rotate
- ⑥ Cutout
- ⑦ Gaussian noise
- ⑧ Gaussian blur
- ⑨ Sobel filtering: Cause of generalizing as well!

Shortcuts: - { NNets are exceptional at finding patterns! }

Trained a model to predict the absolute (instead of relative) posⁿ of a patch.

→ It works until color is removed!



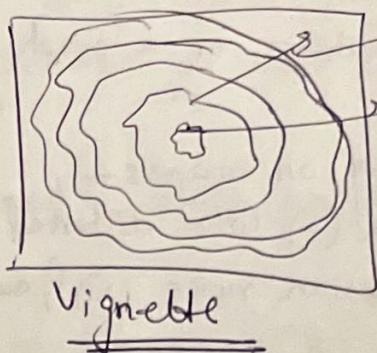
→ "There is some extra info. encoded in the color channels that's very imp for this task"

Chromatic Aberration-

This is inherent because of the optical lens used, there are some inherent patterns in the image and also vignettes.

- ④ There are implicitly encoded in the color channels of an image which might be helping in the use case before.
- ~~Ans~~ → If the chromatic aberration is not centered in the image, then the image might be cropped, we can find whether an image was cropped.

In a photo taken with a camera:



→ a darker gradient along the edges of a photo.

→ Due to camera lenses.

Synthesis:— (Also a supervised learning approach!)

→ Now we can change images by synthesizing diff. features/info. from diff. frames/images to form a counterfeit image.

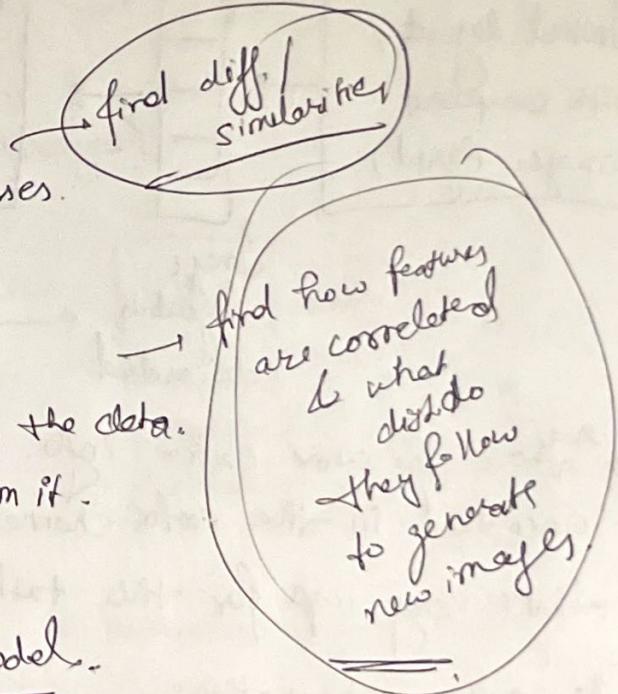
(Basis)

So, how do we build these generative models!

Type of models:

→ Discriminative models:

- Learn the boundary b/w classes.
- Estimates: $P(Y|X)$.



→ Generative models:

- Learn the distribution of the data.
- Often you can sample from it.
- Estimates: $P(Y, X)$.

BigGAN

→ example of the above model.

We were learning discriminative models till now.

Now comes models through which we do not just learn boundary but the distribution. (35)

A lot of models of synthesis are generative models!

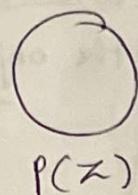
→ we can generate images from natural language (sentences entered).
(even learn reflection of objects).

DALE-2

But this can be used to create fake images now, it becomes an easy task.

⇒ Deep generative models are transformations of distributions ⇒

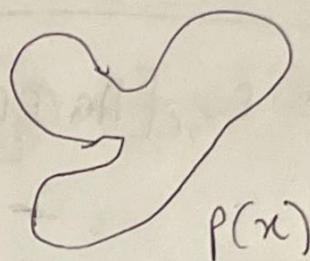
Prior distribution



(done thru. a
net)

G_z
estimate
how to
transform
the
distribution.

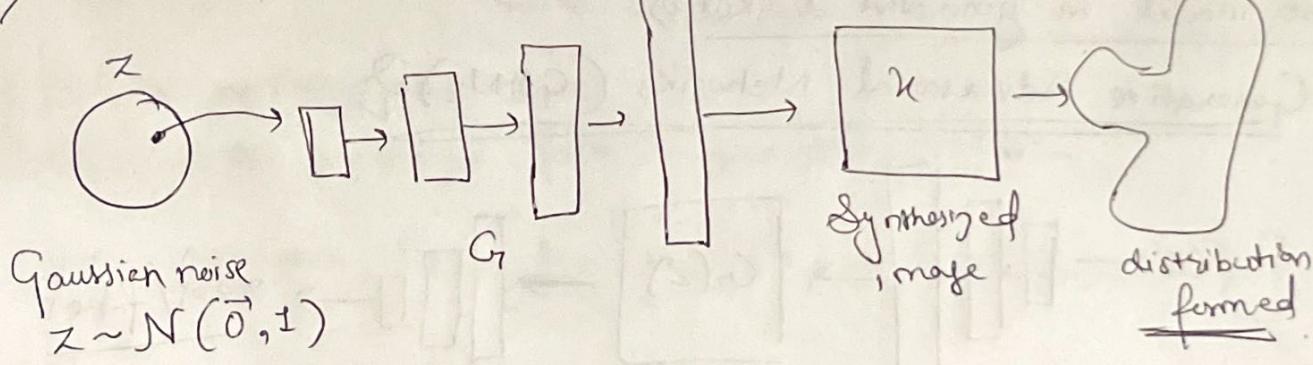
Target distribution



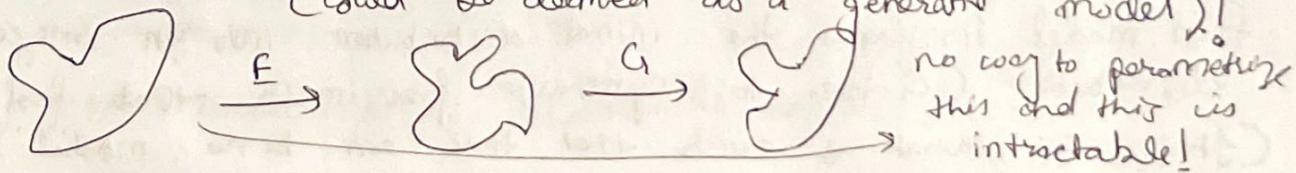
(can be
very
complex)

(learn the mappings)

⇒



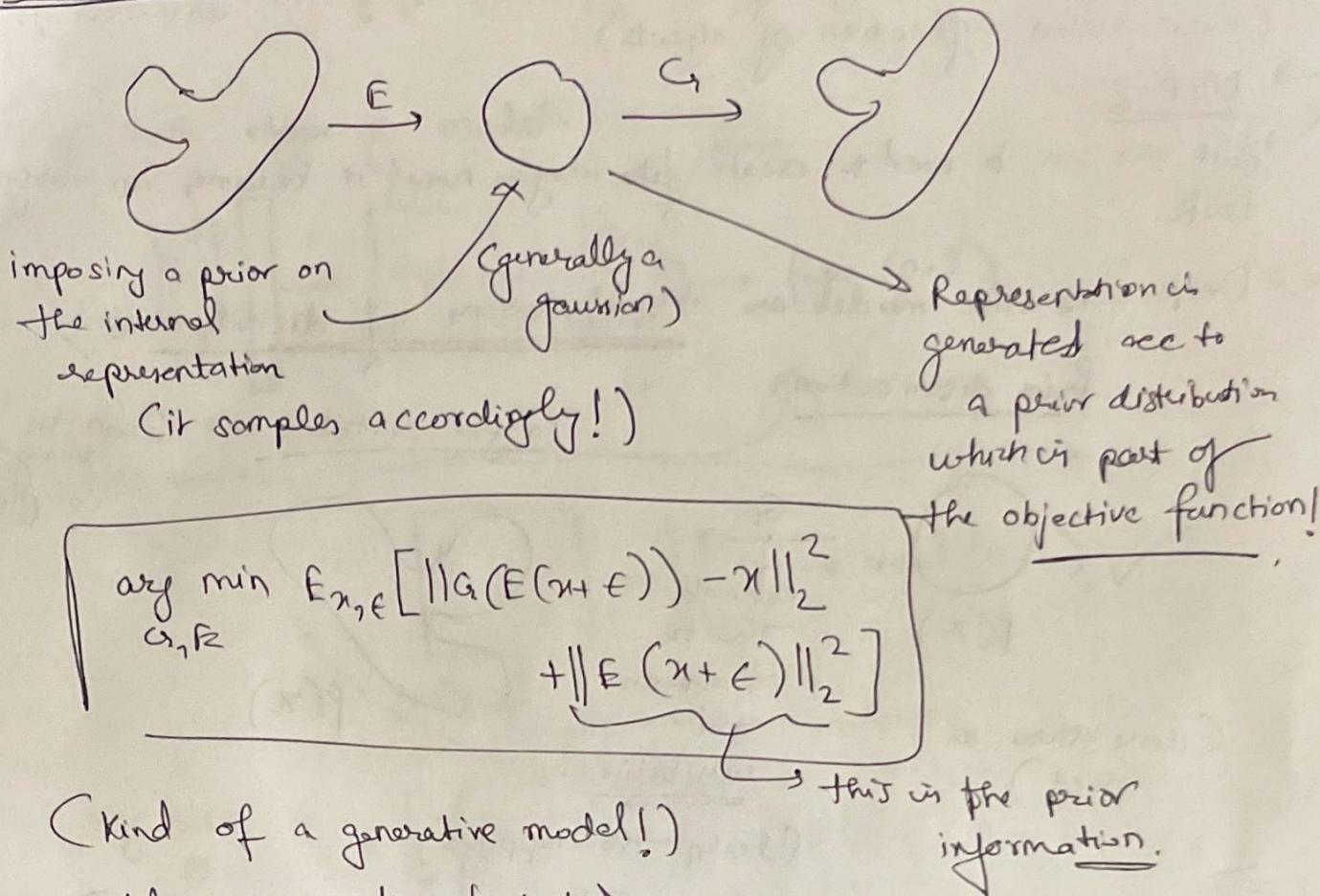
Classical autoencoder ⇒ we transform it in between from which the original distribution can be regenerated.
(Could be deemed as a generative model)!



That's why not exactly a generative one!

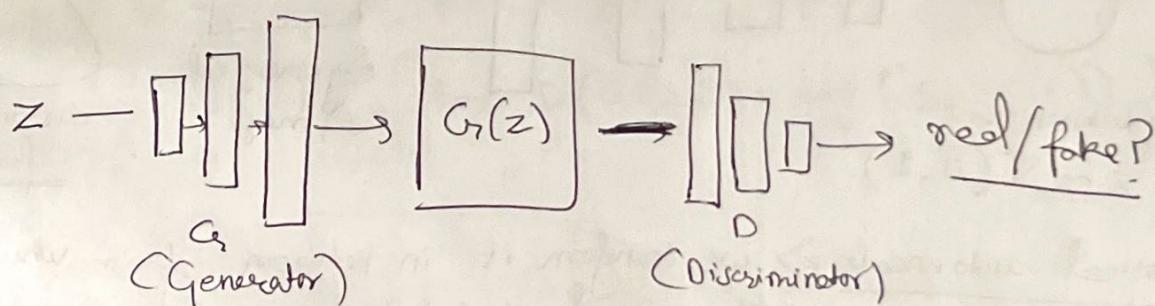
$$\arg \min_{G, E} \mathbb{E}_x \left[\| G(E(x)) - u \|_2^2 \right]$$

Variational Autoencoder



Best model in generative dataset

Generative Adversarial Networks (GANs) ↗



we have 2 diff. models competing against each other.

first model transforms the initial distribution into an image distribution (it tries to synthesize fake images that fool [

(It fails to generate as such that there can be no model that

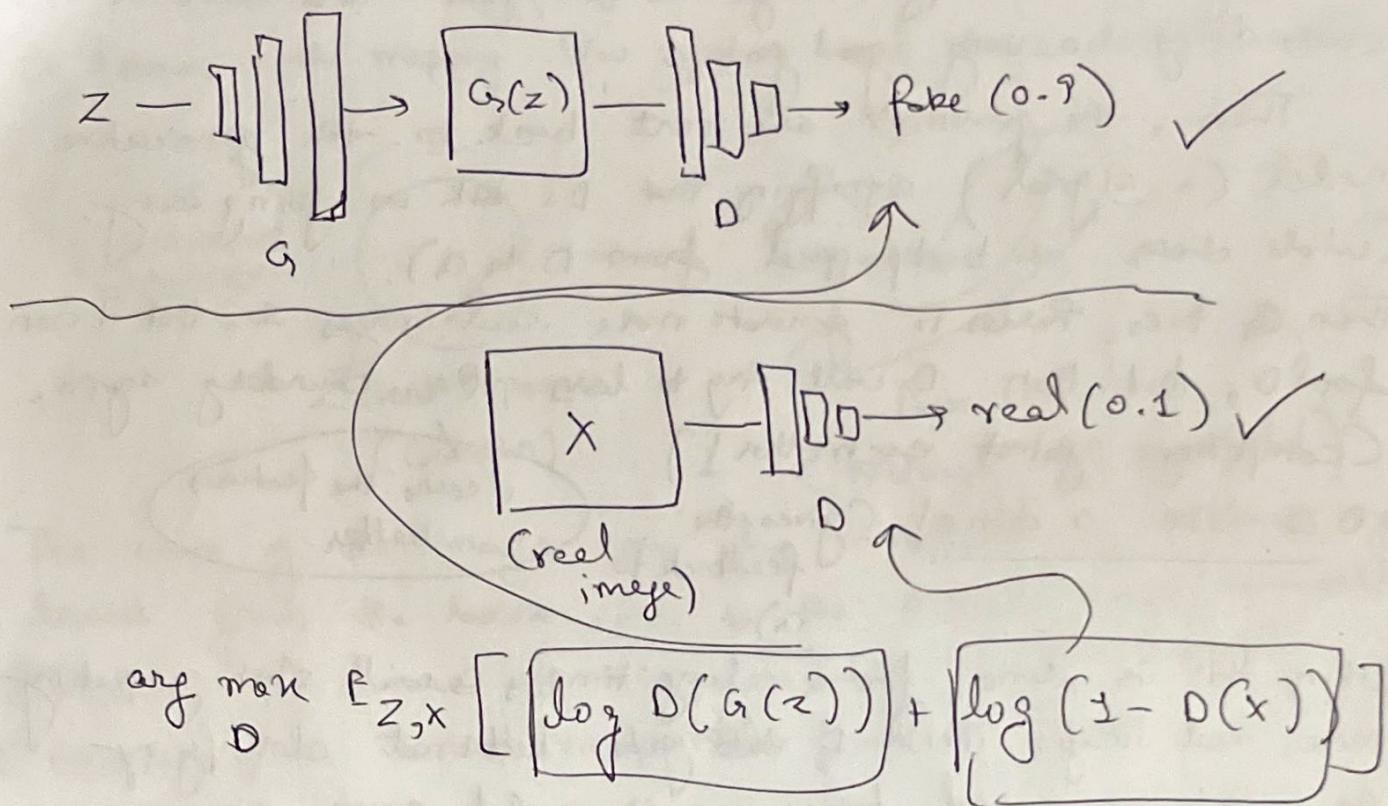
can tell us that $G(z)$ come from a generator or is an actual image) \rightarrow (generative prob.)

(36)

Then, D tries to identify the fakes! (Classification prob)

G is trying to maximize the objective of generating real fakes while D is trying to identify those. G is adversary of D.

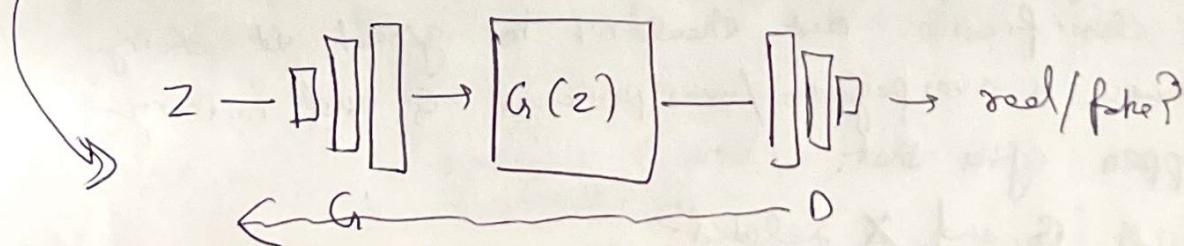
We map both type of models in one.



G tries to minimize the same!

G tries to synthesize fake images that fool D:-

$$\arg \min_G E_{z,x} [\log D(G(z)) + \log (1 - D(x))]$$



Now, when used together. Essentially we are :-

G tries to synthesize fake images that fool the best D -

$$\arg \min_G \max_D \mathbb{E}_{z,x} \left[\log D(G(z)) + \log (1 - D(x)) \right]$$

Objective funcⁿ

Initially, noise is passed to G to generate something, D has an easy job to compare / discriminate real / fake. D learns the boundary to classify image as real/fake and once it starts doing a very good job,

Then, the gradients are sent back to the generative model (a signal) signifying that D 's task is going easy. (whole chain is backpropagated from D to G).

Then, G tries harder to generate more real images so that it can fool D , but then D will try to learn the boundary again.

(Competing against each other!)

D is like a critic! (gives the feedback to G).

\nearrow
Leans the features better.

When this is done for a long time, G will start generating more real images while D also gets better at classifying. It might go around forever as it could never converge. we need to balance G and D such that its descent enough. (That's why hard to train)

④ If D does really good then no gradients will be propagated back to G , effectively stopping the learning, so, D should be good at classification but shouldn't be great at doing that. otherwise, it overperforms / overpowers G and nothing might happen after that.

Are gradients of G and D related?

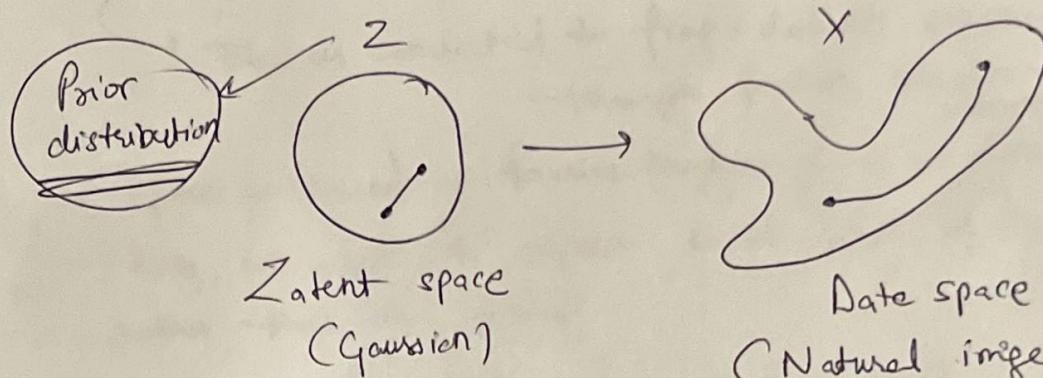
→ They are but not directly! This is primarily because G

never gets to see what it is actually generating, it just knows that what it generated was real/fake (prob. of it being real/fake) and the loss is propagated back from D to G. (37)

* G only learns through the gradients of D!
(Indirect signal).

Works so well ~~because~~ because D tells what boundaries for classes are and that gets propagated to G. G then considers that learning. ($x \rightarrow$ Training image)

It learns the mapping b/w gaussian and image distribution!



The views of the image are also manifested in the mapping learnt from the latent space to the actual image representations.
↳ Angle of the camera \rightarrow this correlation is also learnt
↳ Disentangles the camera viewpoint from the features in an image

(In space, one dimension is direction of view point) $\xrightarrow{\text{correlation found!}}$
(finding new basis in gaussian space)

Different things can also be learnt:

- size of the object
 - color of the object
 - viewpt. of the object
- } StyleGANs

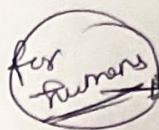
BiGGAN \rightarrow is a massive GAN!

How do we know that GAN hasn't memorized all the images?
→ It's replicated in diff. scenarios to suggest that it's not all

memorization but it is also able to extrapolate.
Mode collapse \Rightarrow from obj. func for GAN, we know that we sample
from z , it's not bound to cover the whole space and
natural image x is very indirectly put into g . (through gradients)
 G just focuses on generating something that helps it trick D.
Sometimes, some classes might collapse altogether. Thus, bias
can come into GANs. This is a problem with GANs.

Vision and Sound

→ Interaction b/w modalities is quite important!

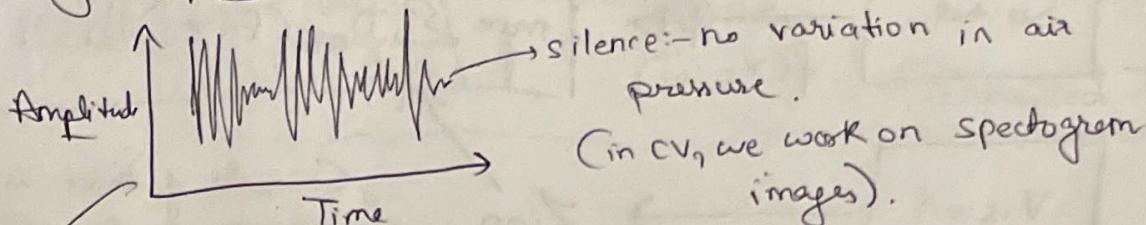


(How do we integrate those together).

We believe our eyes more than other senses.

(bcz. visual system is faster).

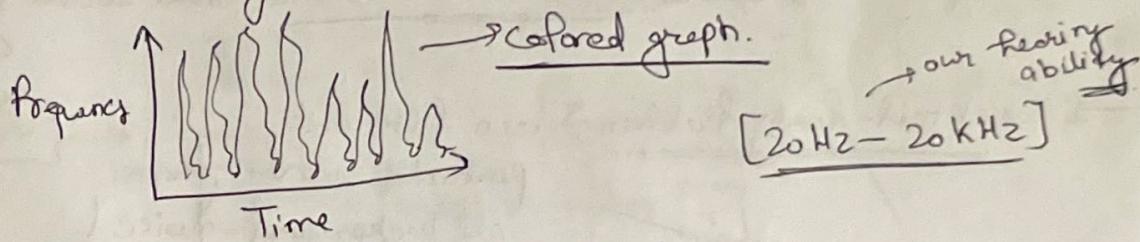
⇒ Waveform: - pressure variation over time.



This is converted to freq-domain representation through Fourier transform.

(often windowed → Fourier transform of a window).

→ bcz. we want to capture local details in the frequencies rather than the global!



⇒ for different noises, we have diff. kind of spectrograms.

{ Sometimes images can be hidden in Fourier transforms of sound }

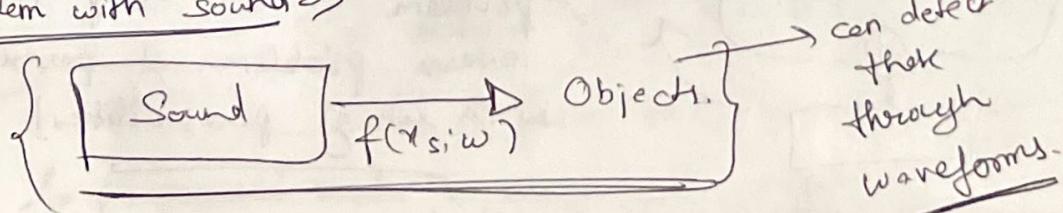
Diff. b/w sound & light ⇒

wavelength of light < wavelength of sound]

+ concept of diffraction, bcz. of this sound diffracts a lot more than light (vision)

④ Light doesn't diffract much. (very little)

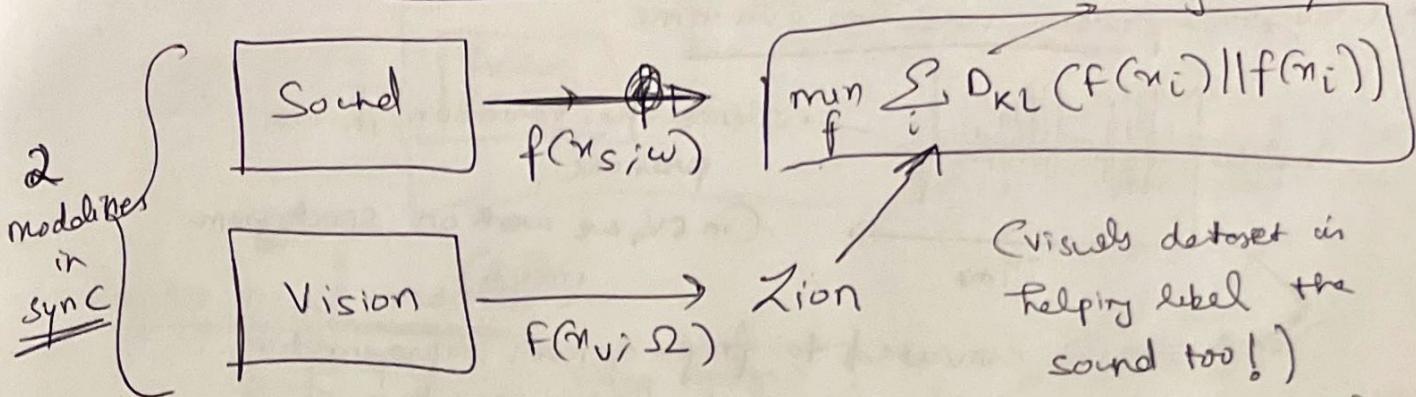
Problem with sound ⇒



Neuroplasticity \Rightarrow When a person loses sight, their visual cortex adapts for sound, smell & touch! Takes about six months.

{2 weeks for this}

There is some "Natural Synchronization" b/w those 2 modalities → divergence b/w the two!



→ we could label many unlabeled videos. (through SoundNet).
 (helps in classification)

⇒ (finding correlations b/w vision & sound).

→ Help devices ~~to~~ be environment-aware thru sound (as power consumption for sound is way less).

→ Cocktail-Party Problem \Rightarrow can filter out particular sounds from all background noise!

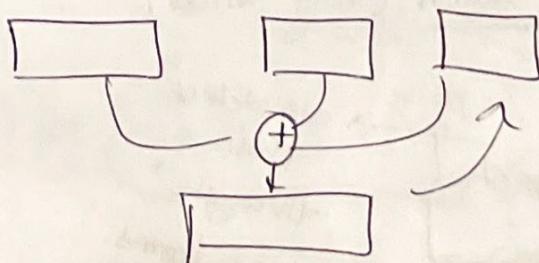
* How can H/C do this?

→ Audiovisual Grounding \Rightarrow

Video + mixture of sounds

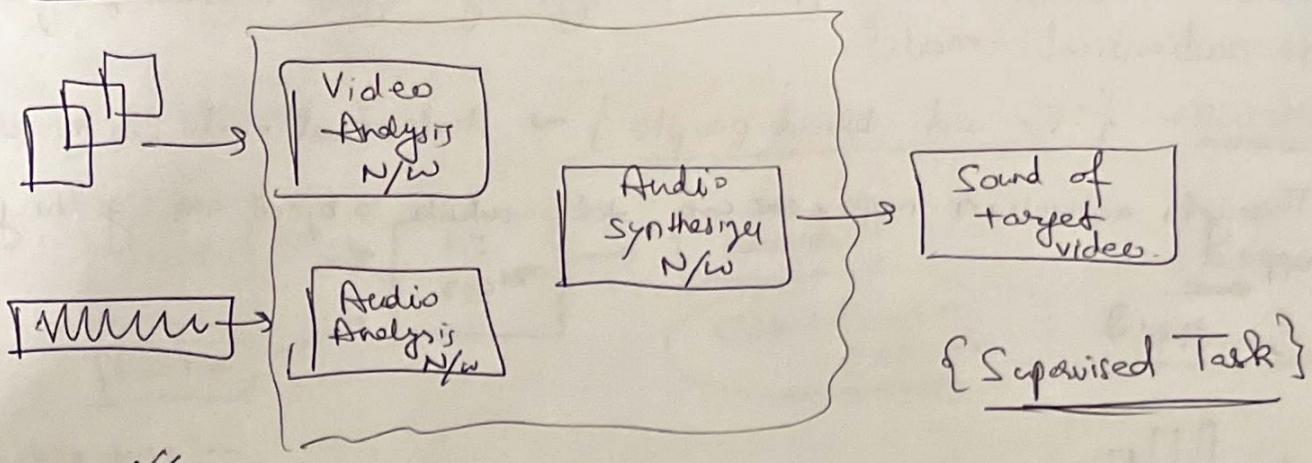
{find which path of sounds come from which objects / path in the video.}

→ How do we recover originals in a mixed waveform of waveforms of various objects / parts?

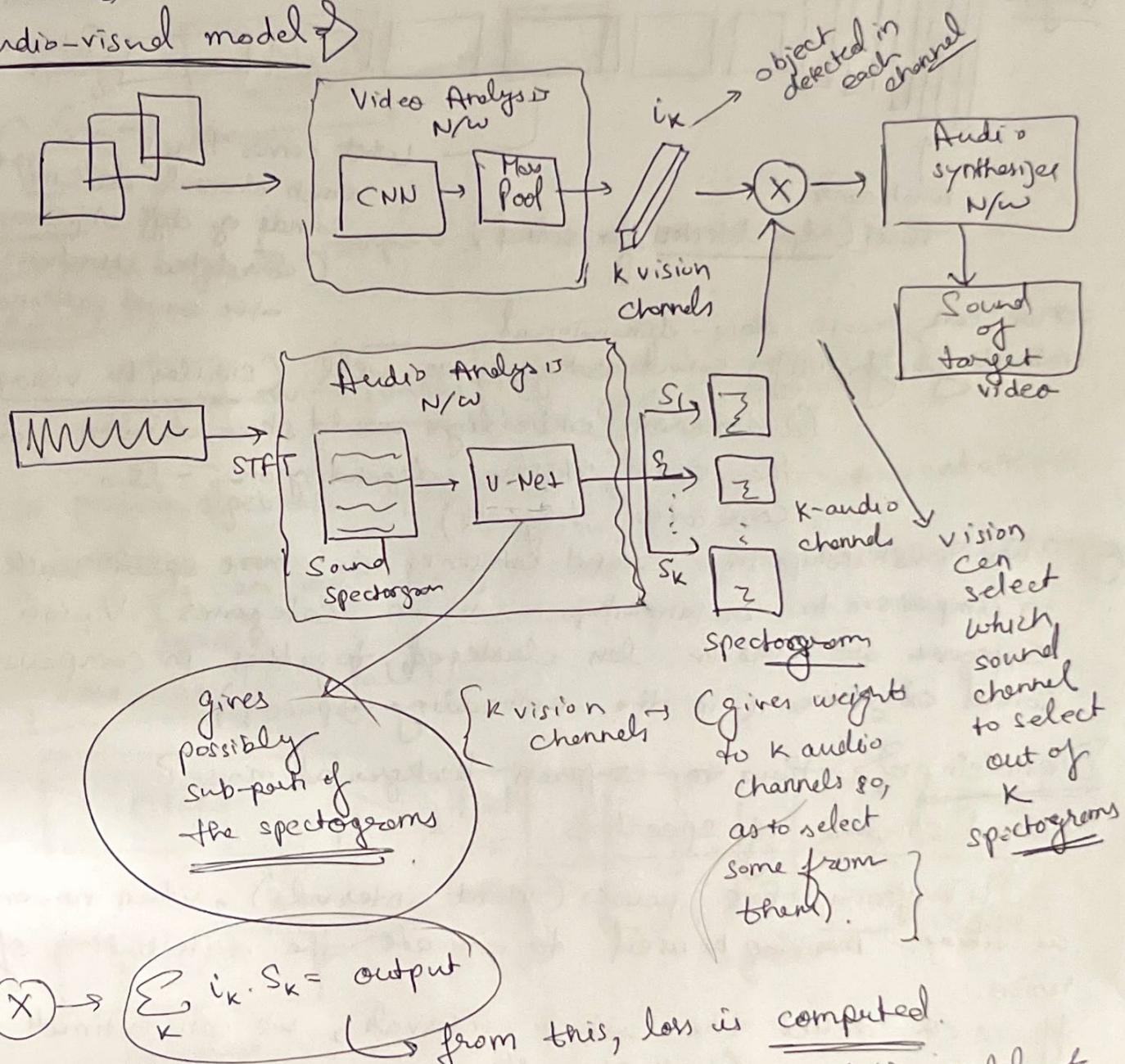


from just audio, it's ill-posed problem as basically it's an inverse problem of permutation with lots of possibilities!

→ Vision can help!



Audio-visual model

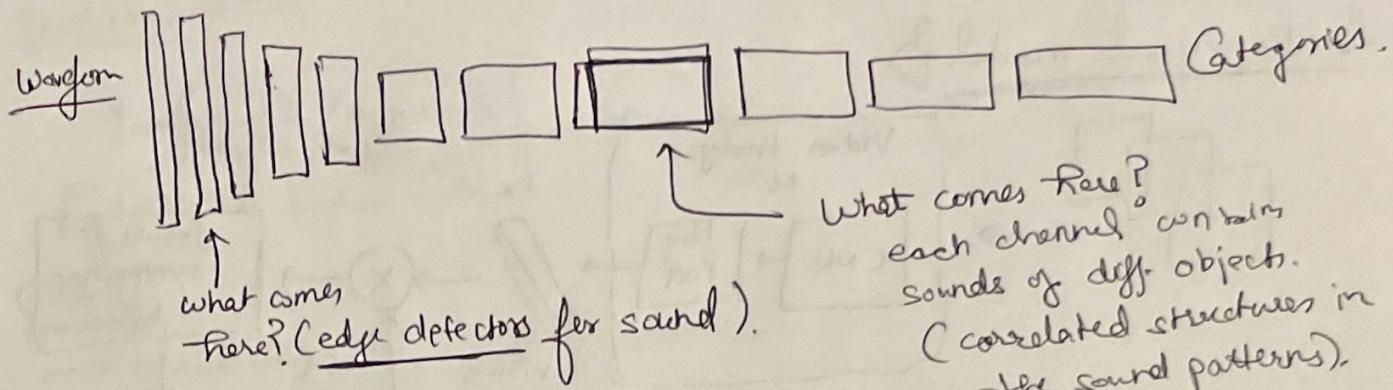


be propagated back to change the weights of all the models & weights applied to spectrograms.

* If we occlude an object from the video, then the associated sound shouldn't be mapped to anything after learning from the audiovisual model.

(B) Side note: {To aid blind people} → Audivoisual mode can be used.
→ Through activation maps, we can tell which objects are getting mapped.

SoundNet ↗



→ we can create low-dimensional embeddings for diff sound categories as well (similar to vision)

(2-dimensional embeddings would show distances b/c the sound categories depending on how correlated they are).

* These relationships in sound categories are more ~~scattered~~ scattered in comparison to relationships in vision categories. Vision categories are more or less clustered together in comparison to sound categories (in the embedding space!).

Denoising ↗ How to suppress background noise?

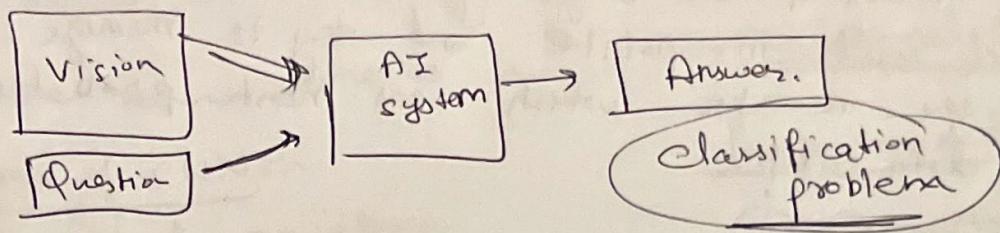
Natural structure of speech

Waveforms have pauses (silent intervals), when no-one is talking. This can be used to estimate the distribution of noise.

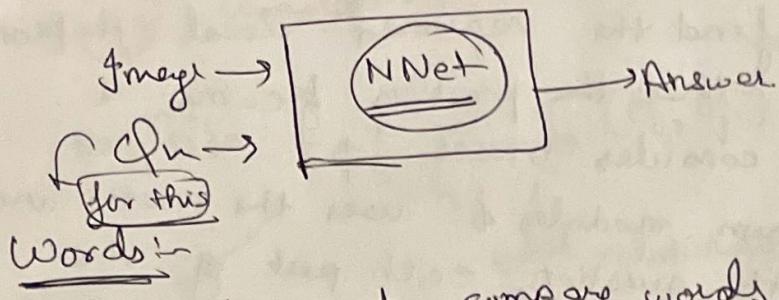
If we can detect the silent intervals, we can estimate the noise distribution & suppress it.

→ It also works for multi-person inputs.
→ It even works for sounds of silence (only detect words / inputs). 40

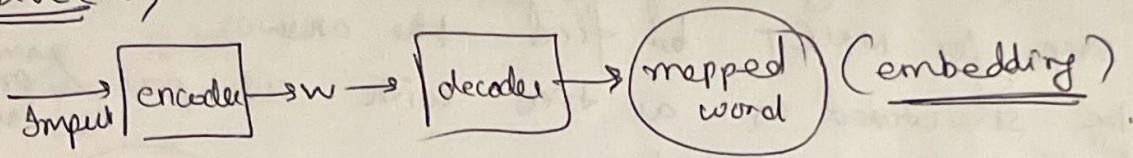
Vision & Language



Architecture :-

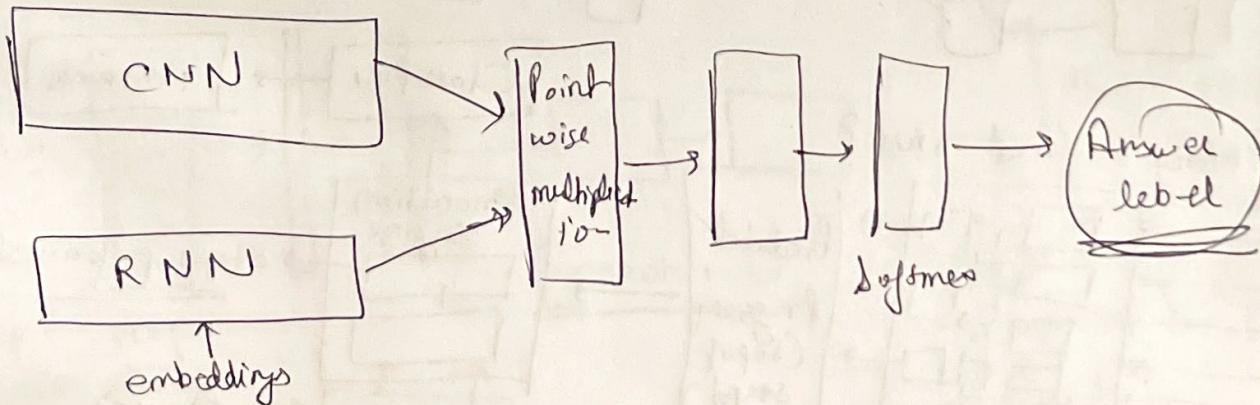


word2vec :-



→ Can perform algebraic operations on vector representation of words,
(vector addⁿ & subⁿ).
(you can even find closest neighbors in the embedding space).

So, now we form full N/w2v



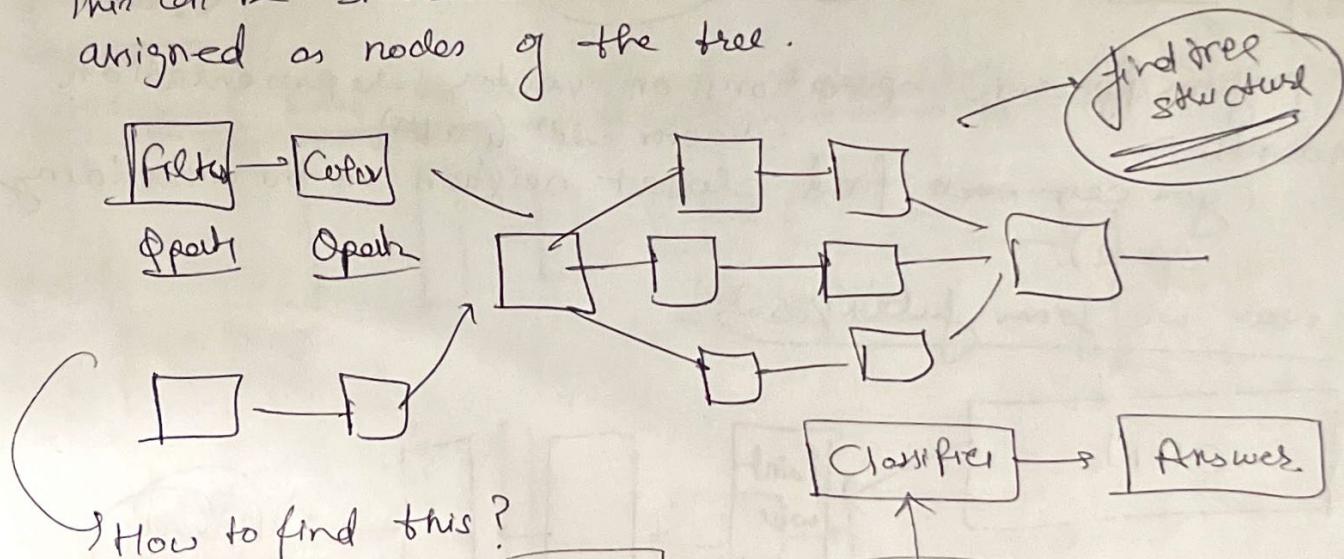
→ Sometimes the results aren't even consistent within itself. Generally, we extrapolate in nature from the knowledge that we have. But here we interpolate because the nature of this model is to find Qn-Answer pairs.

→ If there are no trains in the image given and asked for no. of trains, the model will interpolate and try to minimize the error and give the number which was most probable according to the training data set!

Compositional VQA → When you need to consider all the contextual knowledge in an image to find the required local significant features acc. to the question, then the problem becomes a hard one. Also, we need to consider biases.

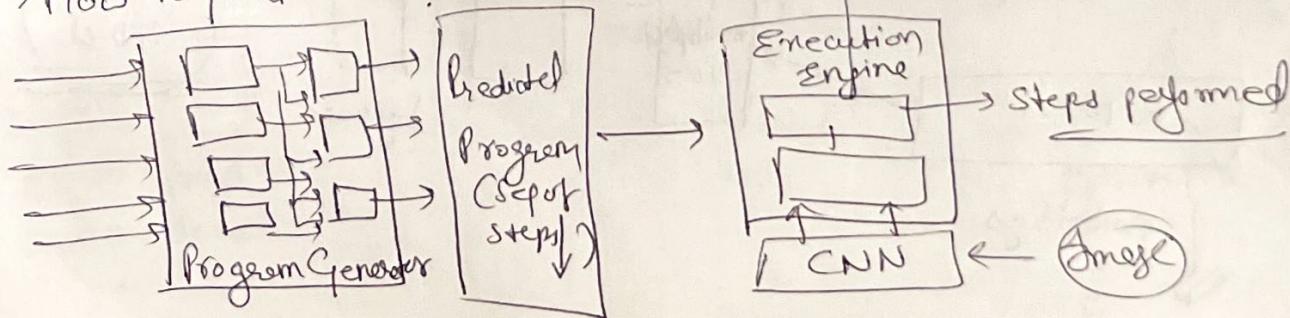
→ We have to find the program modules & uses the NNET and assigns tags to each part of the question, each part of the question will be assigned to a program/question. This, basically becomes steps for NNET to find the answer.

This can be structured as a tree with sub-parts & programs assigned as nodes of the tree.



How to find this?

EMBEDDINGS



Relational N/w :-

(40)

Based on a relational question, we need to find representations of the question to find the final answer. Because of so much structure before, one assumption (wrong) can impact a lot. Thus, generally, we won't follow a similar structure as before.

→ Is there nn same as nn? (Relational Qn).

Properties of CNNs:-

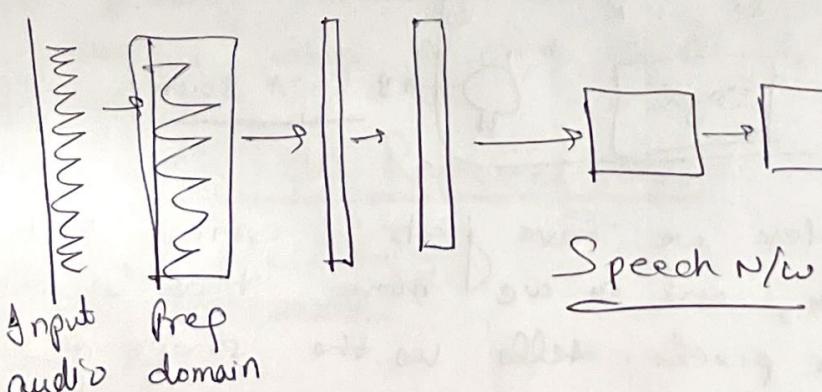
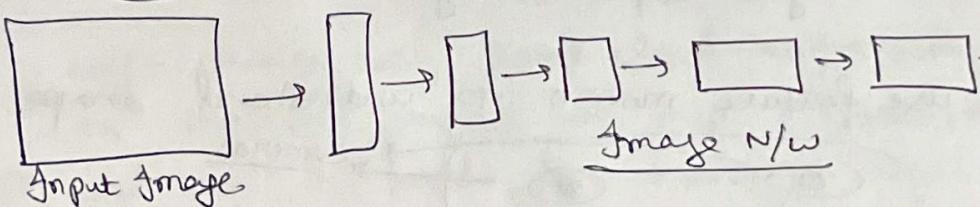
- 1) Spatial equivalence.
- 2) Translational invariance

$$RN(O) = f_{\phi} \left(\sum_{i,j} g_O(O_i, O_j) \right)$$

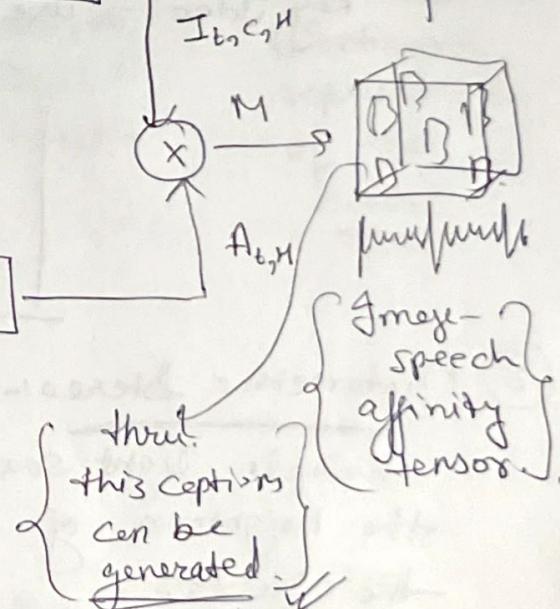
(For this also, we need to distribute the image into patches!)

~~(relationships b/w pairs of object in the image (finding that and adding up for all the possible cases)).~~

DAVENet (Deep Audio-Visual Embedding N/w) :-



Correlating areas
Computing image-caption similarity



3D - Vision \Rightarrow

\rightarrow why put ~~on~~ eyes together? \rightarrow 3D Perception (depth sense)

\Rightarrow Depth is so much more important.

Distance b/w eyes / cameras \rightarrow baselines.

Stereo vision \Rightarrow when eyes are seen together with some distance b/w the eyes.

(the far eyes are from each other, the better they are at depth perception).

Predators will have stereo vision with large distance.
Prey have eyes as far as possible to cover 360°.

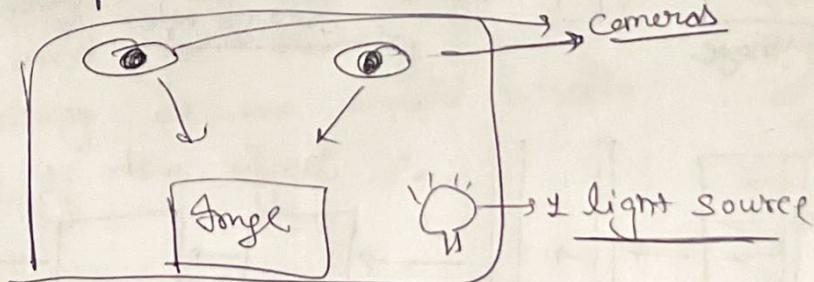
\rightarrow Depth without objects?

We do not perceive depth thru. objects but due to some low underlying detection by 2 eyes, (in an integrated way).

Ways to perceive depth \Rightarrow

① Binocular stereo - when you move camera from one place to another one and take the same image, if the pixels for the object has moved by a lot, then that object is actually nearer than the object that didn't move much across the images captured by the cameras at the two positions.

Key Idea - use feature motion to understand shape.



② Photometric Stereo - Here we have just 1 camera but multiple light sources, and as we move these around, the brightness of the pixels tells us the shape of the object.

Key Idea:- use pixel brightness to understand shape.

(42)

Why learn 3D vision?

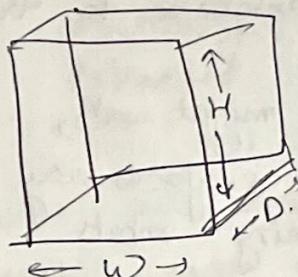
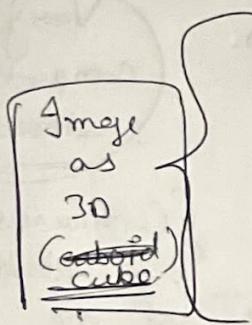
→ To calculate the depth of the pixel, we can move the view pt around see the distance traveled by the particular pixel but this makes perception quite hard, that's why we learn 3D vision.

→ Due to the light source direction, the depth perception changes & it also does this bcoz. of vanishing pts.
We need to combine priors about the env. to correctly perceive 3D data.

→ The world is not ideal diffuse.

LIDAR Sensor → works like sonar to detect objects and can also be used to create 3D models of the environment
(Used in self driving cars)
(Car is not good with LIDAR).

→ We represent 3D Images as Voxels. (VOXEL) (Volume element)
(similar to pixel).
(represented as a cube).



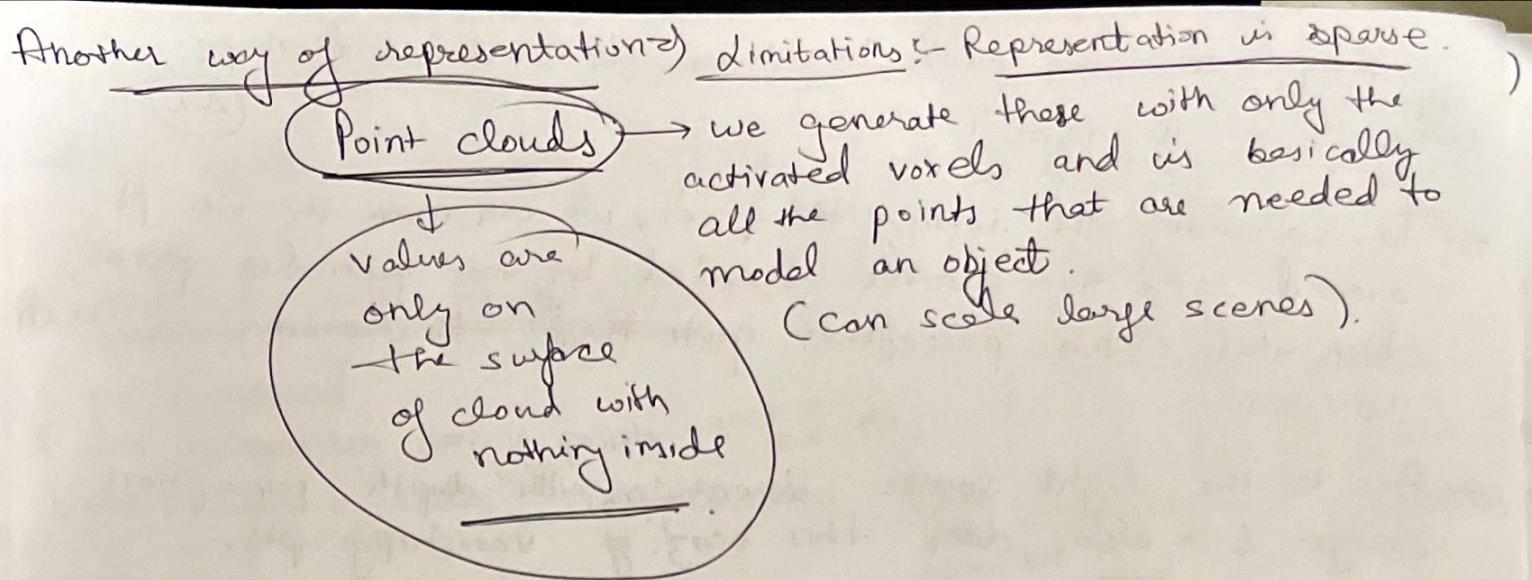
Limitations:-

a) Very memory intensive as most of the space is empty.

b) Trade-off with voxel resolution.

(depth can increase by a lot if we break the object into very small voxels).

To reduce memory, we might have to make a trade-off of the voxels to represent an object.



Method → Represent diff. parts of the object as made up of Δs.

Limitations:- Difficult to integrate with neural networks and takes expertise to label.

No camera/sensor directly senses in a mesh-form.

Implicit surfaces →

$$f(x, y, z) = 0$$

only parameters of this func are stored

(Model a 3D shape / scene with func value on the surface of an object equal to zero, otherwise not zero)

Limitations:- Computationally extensive to train.

Very compactly stored.

* For a self-driving car, point cloud might work but with implicit surfaces, we map surfaces very well and can be used for training robots while grabbing something as we would need to know how surface of the object is.

(Continuous representation of the scenes)

Learning with 3D Representations:-

With Point clouds →

PointNet →

→ classification → part segmentation → semantic segmentation.

43

→ Point cloud is basically a list of values (for all pts. of a surface)

Q Challenges with Point Clouds -

→ Unordered

→ Interaction among points → we don't know this!

→ Invariance under transformation.

Consider adding 10 to everything,
point cloud's origin is arbitrary

So, nothing changes (invariant)

but it would impact brightness/something for the image.
(we need it to change/vary mostly!)

X	Y	Z
-	-	-
-	-	-

can have
-ve or
well as
+ve
values

helps
with
CNNs, but
we don't know!

Q How to build a NNET for this?

→ what about sorting?

→ Pick a deterministic order, and feed results into a CNN?
(still not able to work with invariance under transformation!)

(Also, if there's noise, it might change the ordering drastically, which might impact relationships extracted by CNNs).

Thus, this won't work well!

Q → what about invariance?

① Treat the order as data augmentation.

② Shuffle sequence, feed into RNN & train.

Shuffle again, feed into RNN and train, etc.

③ Will this provide invariance to the order?
(no, as we can't go thru all the permutations possible.)

PointNets:- Invariants

tries to fix these problems above!

$$f(\{x_1, \dots, x_n\}) = g(h(x_1), \dots, h(x_n))$$

$\left\{ \begin{array}{l} h \rightarrow \text{simple neural network} \\ g \rightarrow \text{order-invariant func} \end{array} \right.$

Examples:- sum, mult, max

Technically, this is a universal approximator, and it could learn to create voxels by partitioning the space. But it learns something more powerful!

Scene Segmentation & Classification can be done using this.

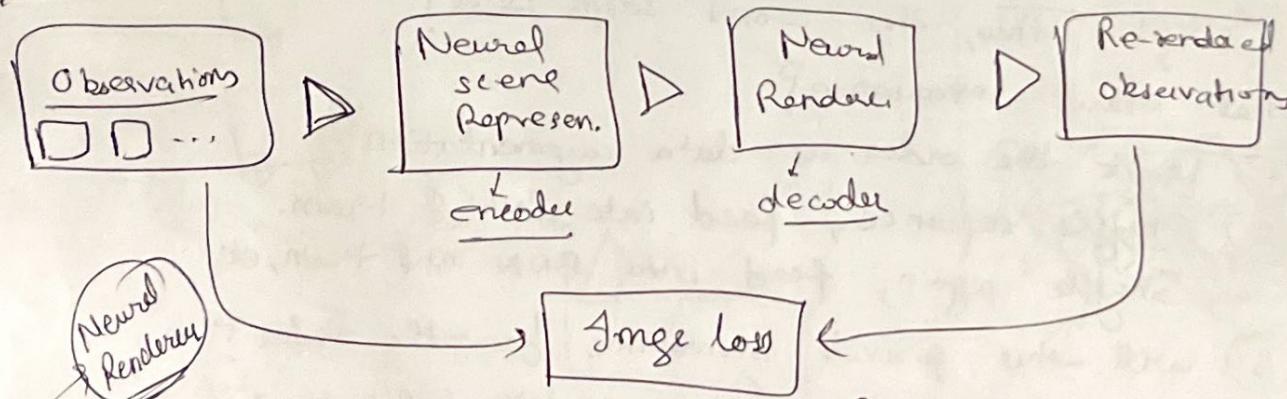
Q What are the critical points?

Which I/P points are contributing to the global feature?
 (global features)
 (from the layers) { we need to get these! }
 ↳ to see how model is learning.
 generally determines the shapes.

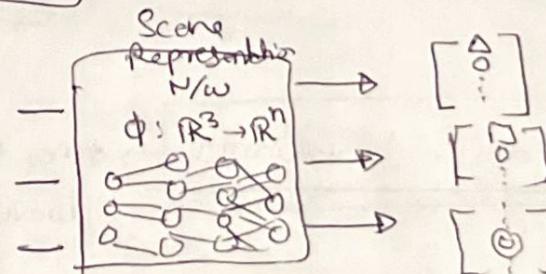
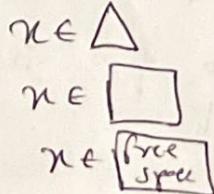
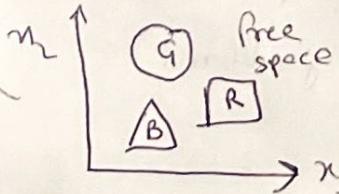
Upper-bound shape - critical pt ⇒ points that do not matter!

Implicit Surfaces

Scene Representation N/w sd) (write an auto-encoder).



Q How do we do this?



Ans

The scene is converted into parameters of the neural network.

(44)

Neural Renderer Step 1: Intersection Testing.

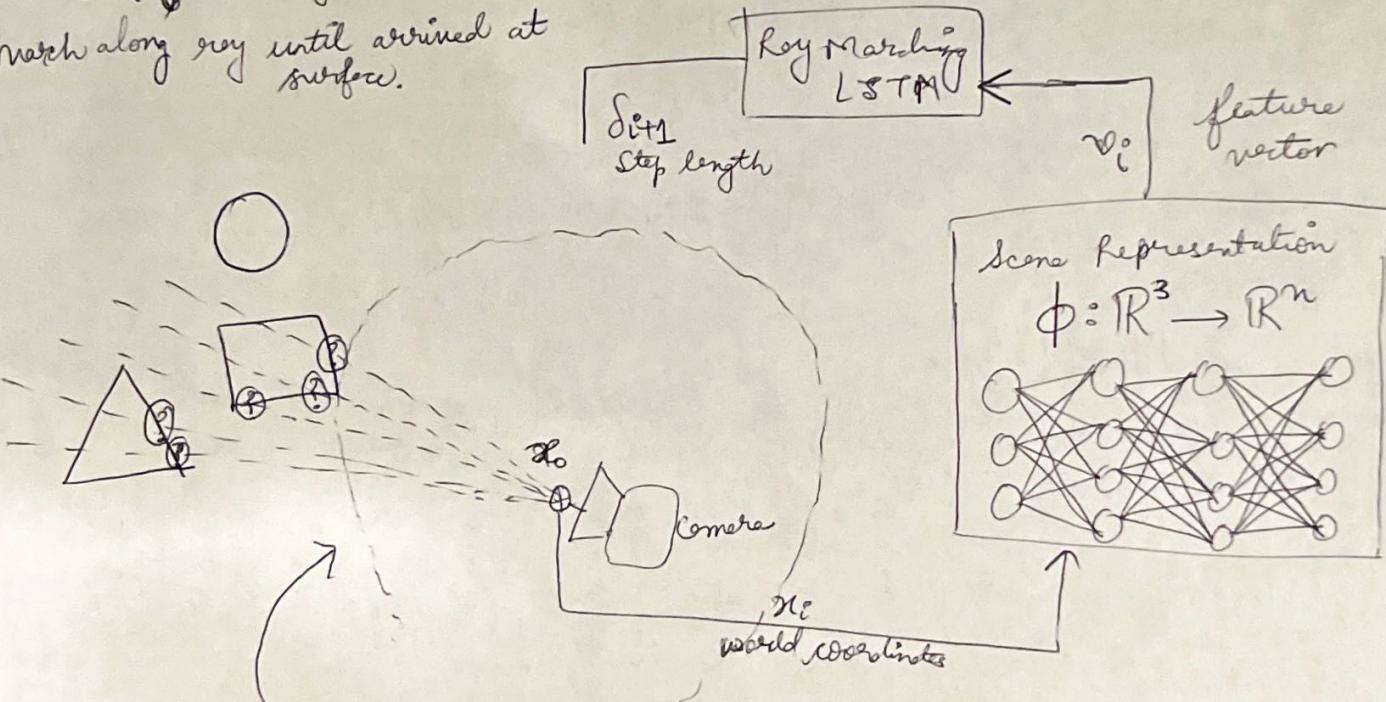
To represent a scene, we need to input the image and figure out what the weights are of this NN. Since NNs are universal approximators, there is a way to take any 2D/3D scene & compress it into the weights of a NN.

The Decoder (renderer) is going to learn to query that representation.

We can combine this renderer with classical algorithms in computer graphics. Eg of intersection of NN + graphics.

→ It needs to read out from the encoder to reconstruct the I/P.

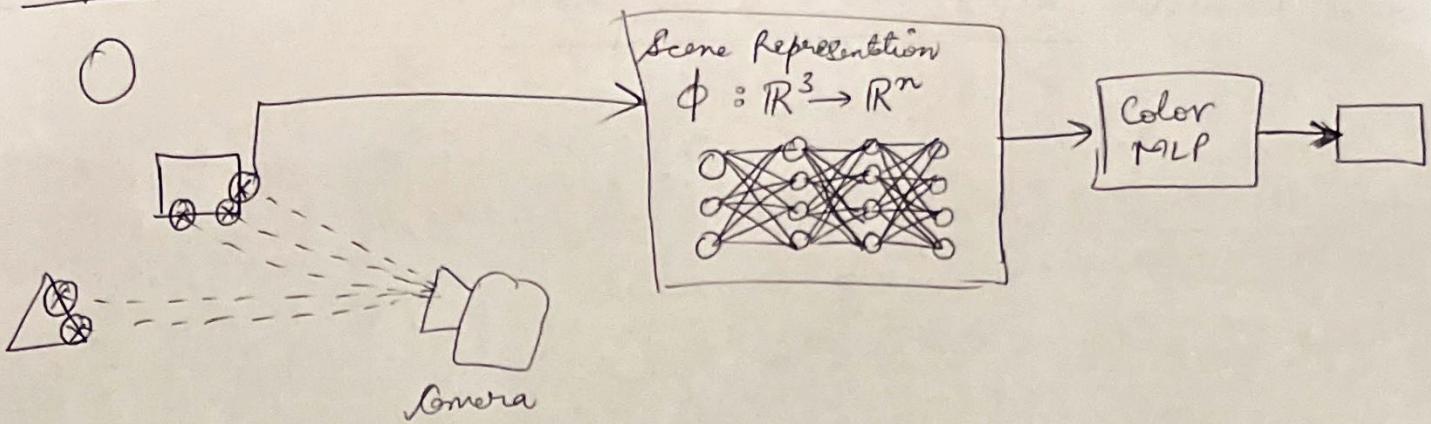
* Ray Marching used in Neural renderer Step 1: intersection testing
March along ray until arrived at surface.



feasible step length : distance to closest scene image

When you shoot your ray, march along it & hit a surface, you record that value; example: color of the object.

Step 2 : Color Generation



Can now train end-to-end with posed images only!

* You can backpropagate at each step since every step is differentiable.

Ray marching LSTM tells you WHICH DIRECTION to walk in and also stores your history of the direction you have taken. LSTM tells you how to sample along that ray. The scene representation is continuous & so are the rays.

Because the scene representation is continuous, we can easily increase the resolution of the camera (eg 1Mpix camera) by increasing the # of rays.

This ^{scene representation (encoder)} model works well since it samples a couple different views of the object.

Ray marching assumes that we are operating w/o any transparent objects and that lighting does not matter.

→ we move the camera around if we cannot see an object much/ not at all!