

lc 1

$$x_1 \rightarrow \begin{matrix} \text{P} \\ \boxed{1} \end{matrix} \rightarrow p(x_1) = 0.98$$

$$\square \rightarrow \begin{matrix} \text{P} \\ \boxed{0} \end{matrix} \rightarrow p(x_2) = 0.01$$

↳ handwritten digit don't have black pixels at LT corner.

$$\bullet \rightarrow \begin{matrix} \text{P} \\ \boxed{1} \end{matrix} \rightarrow p(x_3) = 0.09$$

$$\left[ \nabla_{x_i} p(x_i) \right] \because \nabla_{\text{ideal}}^{\text{ideal}} \text{ means lots of values.}$$

Not  $\frac{\partial p(x_i)}{\partial x_i}$  because we calculated it for all 9 different pixels

$$\text{loss} = f(w, x) = \text{MSE}(\hat{y}, \overset{\text{actual}}{y})$$

$$\frac{\partial \text{loss}}{\partial x_{1,1}}, \quad \frac{\partial \text{loss}}{\partial x_{1,2}}, \dots, \frac{\partial \text{loss}}{\partial x_{1,n}}$$

↳ change in loss w.r.t. value of pixel  $(1,1)$

↳ pred or sample

Hower

$\Delta$  - Delta mean Small change

$\frac{\Delta \text{loss}}{\Delta x_{(i,j)}}$  if we make small  
charge small enough,  
it will become delta  
derivative. So,

but in note

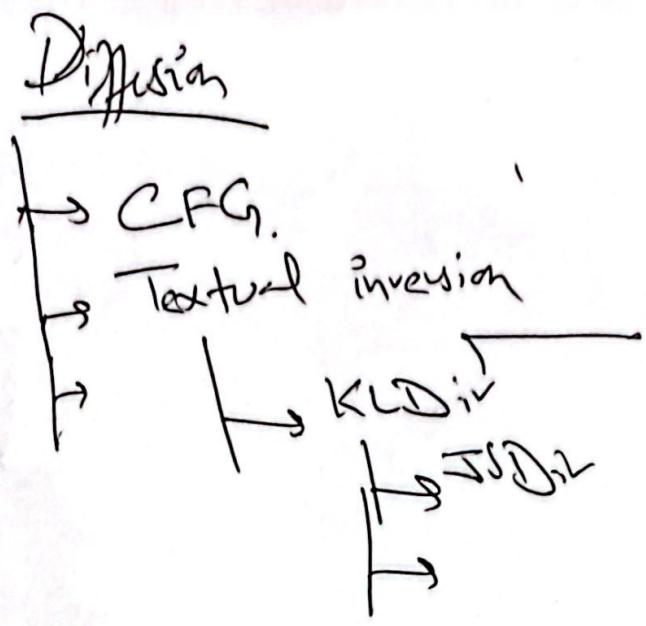
$\nabla \text{loss}$  is actually  $\nabla \text{loss}$  for too many  
 $\nabla x$  values

So, instead of changing weights of  $p_i$   
a model, we are changing inputs into the  
model. So, we are trying to make  $x_3$  more  
like digit. So, we can change pixels accordingly

to this gradient  $\nabla_{x_3} p(x_3)$ : So we

are gonna take every pixel and  
we are going to subtract its ~~no~~ gradient  
a little bit.

$$x_3 = x_3 - c \times [\nabla_p(x_3)]$$



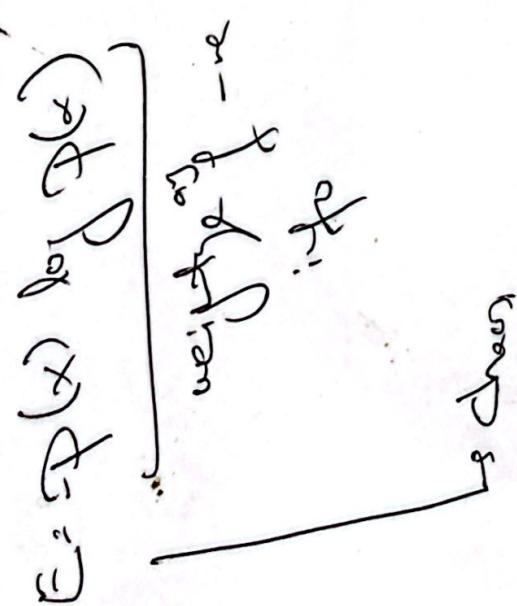
U-Net - predict noise

Encoder :  $\rightarrow$  latent

Decoder : latent  $\rightarrow$  img.

Computationally expensive

$\rightarrow$



$$H(x) = -\sum_i p(x) \log(p(x))$$

$$D_{KL}(P||Q) = \sum_i p(i) \log \frac{p(i)}{q(i)}$$

2)

As I chose right pixels, how does it change probability that this is a digit. which tells which pixels to make darker which pixels to make lighter.

So, when you do this by doing each pixel one at a time to calculate derivative, this is called a finite differencing method of

calculating derivatives and it's slow.. because we have to call the function 784 times for a  $28 \times 28$  image but it's slow & we would have to use finite differencing.

we can do this by

→ f.backward()

→ Xs.grad

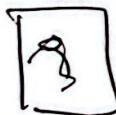
So we don't particularly need the thing that calculates probabilities. we only need the thing that tells us what pixels to change to calculate probabilities.

but where is  $\hat{f}$ ?

$f$  is a NN that tells us which pixels to change to make an image look more like handwritten digit. So, we could create some training data & use that to get info we want.

We can pass samples like

- digit
- doesn't look like digit
- more corrupted digit



{  
↳ the char which is more  
like a handwritten digit }  
which occurs less but that  
would be complicated. What  
if we could predict how  
much noise we added?

4)

So, we take input & pass through NN.  
we try to predict what noise was ( $\hat{n}$ )  
and actual noise is  $n$

so loss fn would be

$$\text{loss} = \sqrt{\frac{(\hat{n} - n)^2}{\#}} \quad \therefore \text{MSE}(n, \hat{n})$$

- So input noisy digit whose  $n$  we are controlling
- we predict noise for it.
- calculate  $\text{MSE}(n, \hat{n})$ .

→ we have  $n$  because we are mainly adding it

So, if we can predict the noise, we get exactly what we want which is

$$x_3 = x_3 - \underbrace{p(x_3)}_{\downarrow} c \cdot \overline{\sqrt{p(x_3)}}$$

$$\boxed{x_3} = \boxed{x_3} + \boxed{\epsilon} \sim N(0, 0.3)$$

Squash this and remove that from  $\boxed{x_3}$  to get  $\boxed{x_3}$

So, if we just input pure noise to the model (NN), it's gonna split up like about which part of it, it thinks is noise and it's going to leave behind the bits that look most like digits.

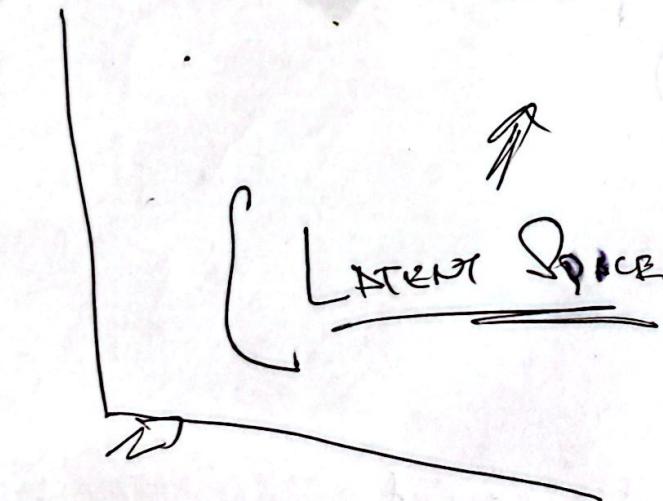
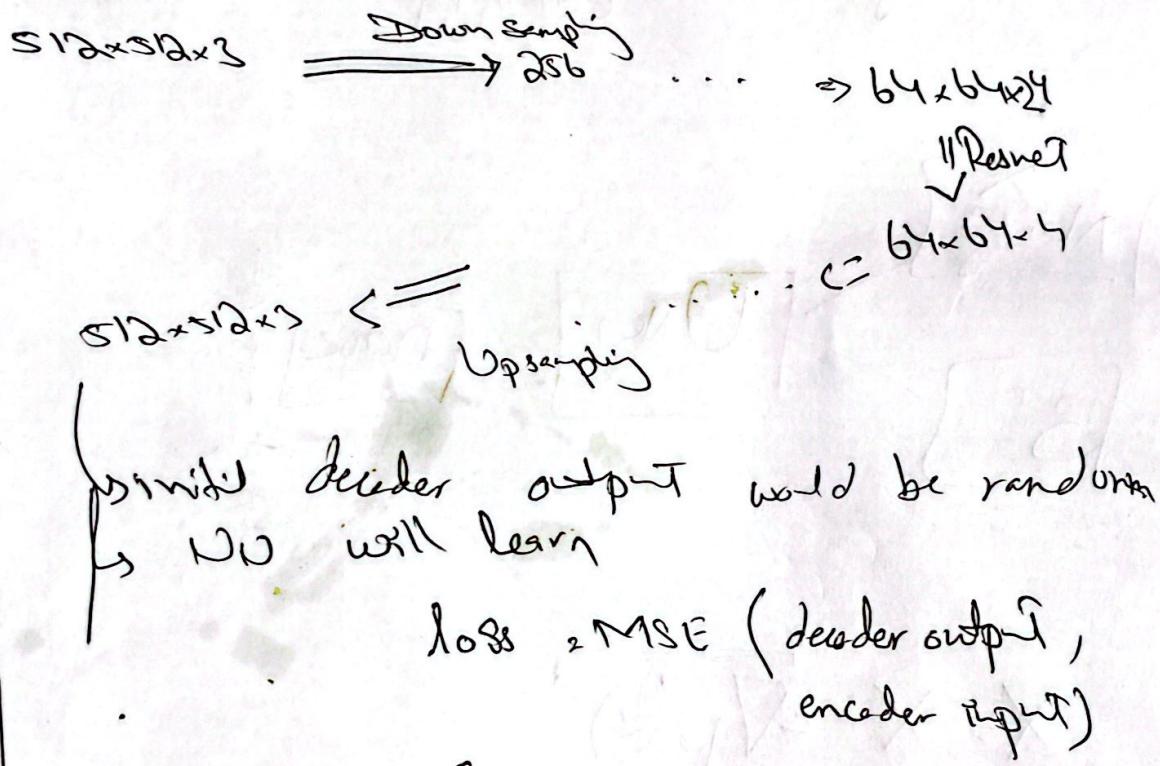
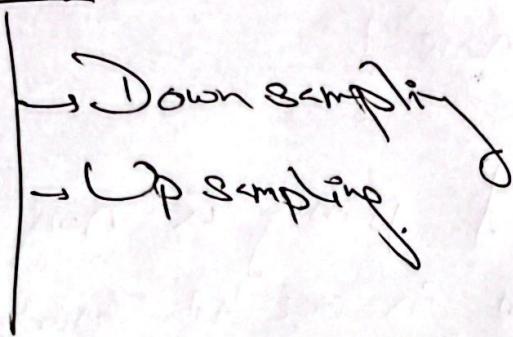
The NN is called U-Net.

So, for U-Net

- noisy img (input)
- output : noise

5)

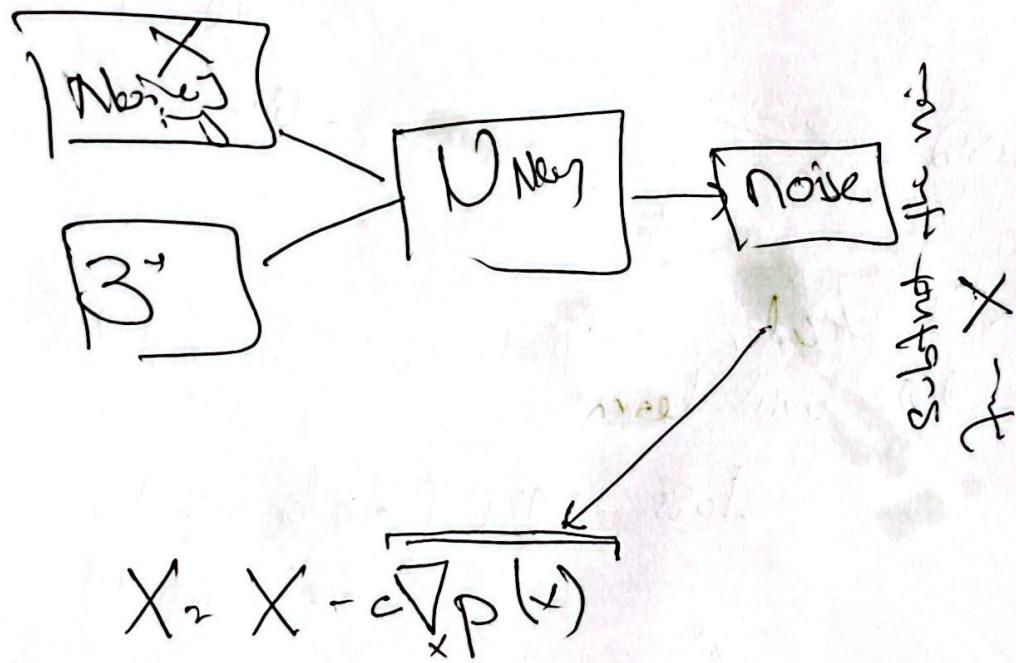
Encoder:



Dimensional Reduction  
Noise Removal

but how to tell "produce me an image  
of teddy riding a horse"?

In addition of passing in noisy  
input, we can pass in one hot encoding  
of what digit it is [HOT].



This is called Epidane.

but for "a cute teddy" we can't  
use one hot encoding

3)

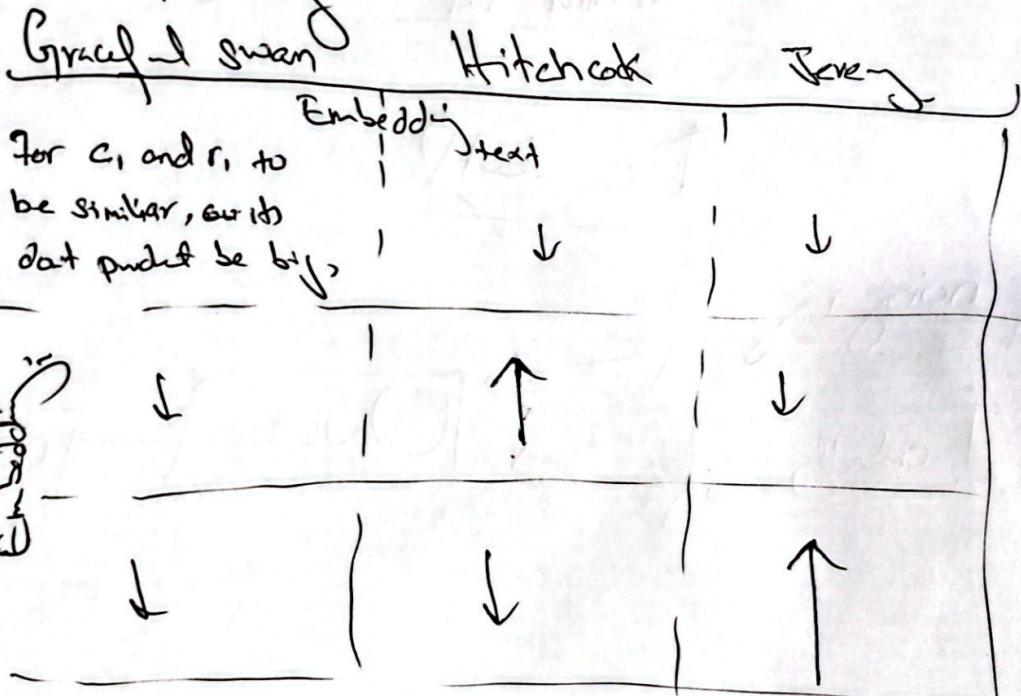
## Embeddings

" a Cut teddy "

Start create two model

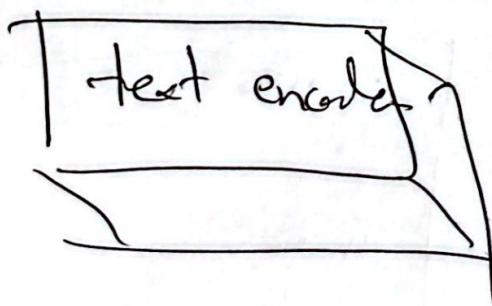
→ Text encoder

→ tiny encoder

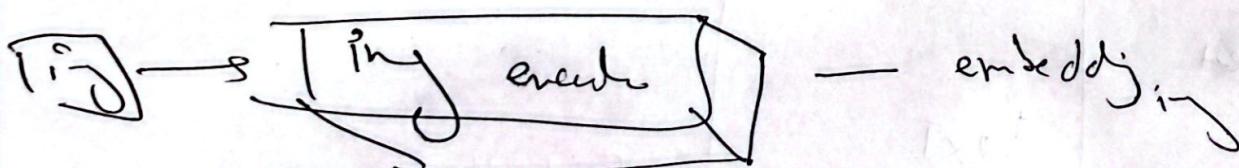


↳ All diagonal entries are selected.

" a graph  
swan "



embedding text

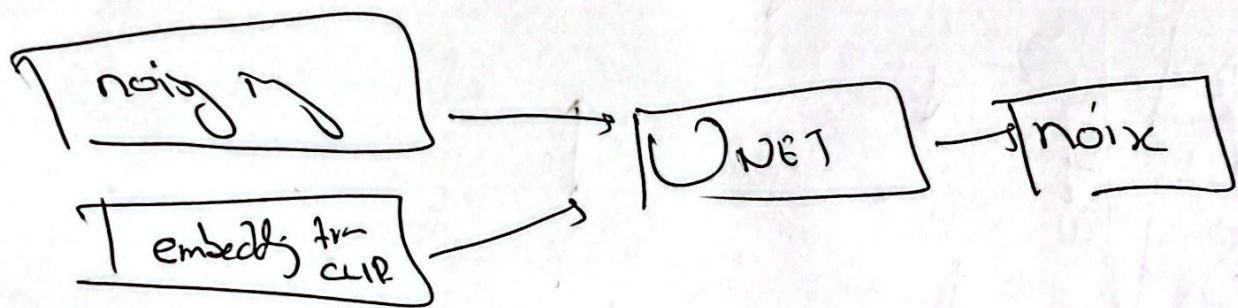


Dot product =  $\sum_i$  (element wise multiplication)

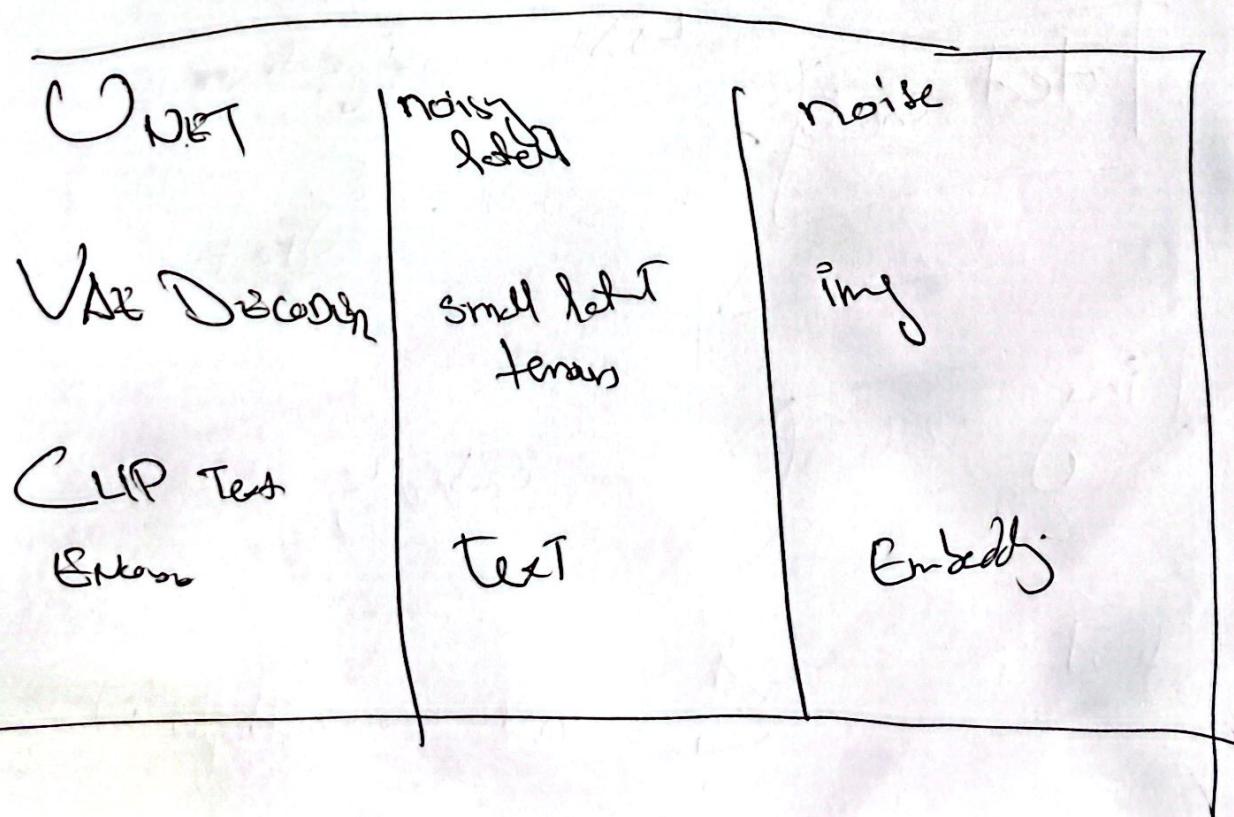
$$\{ \begin{pmatrix} \text{diag} \\ \text{esti} \end{pmatrix} \} - \{ \begin{pmatrix} \text{non-diag} \\ \text{esti} \end{pmatrix} \}$$

"Contrastive loss"

"CLIP"



8.



8)

8



$$\boxed{\text{Latent}} = \boxed{\text{Latent}} - c \cdot \underbrace{\left[ \nabla_{x_j} p(x_j) \right]}_{\text{gradients}}$$

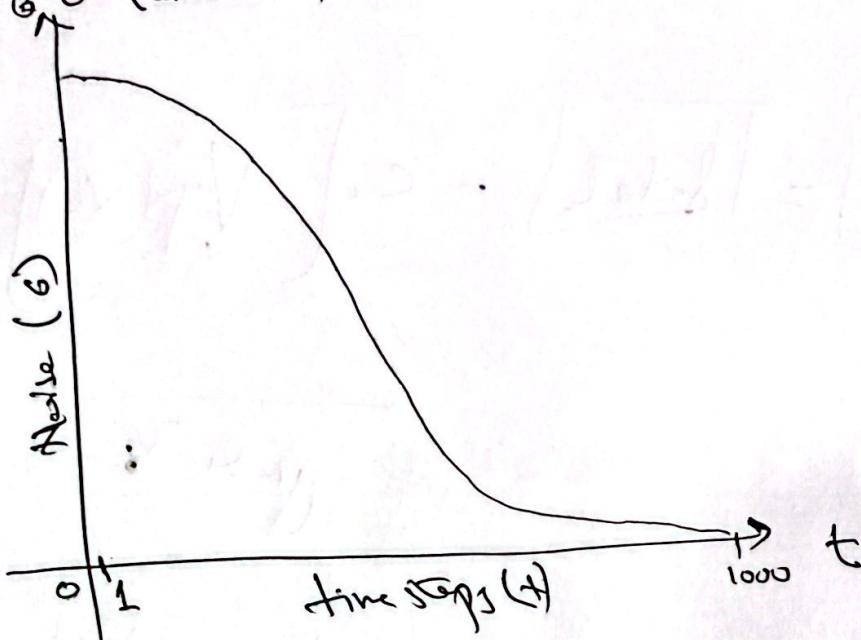
L  
D(CNN)  
C Degrees  
P Pics

→ These gradients  
are often called  
Score fn.

## Time steps

↳ we use varying amount of noise.

$\beta_t \mathcal{G}$  (amount of noise to use)



Pick a random  $t$  from 0 to 1000 and that's how we'll add random noise. I mean max noise, 1000 means minimum, min

Noise  $\beta$  can also be denoted as  $(\beta)_t$ , more commonly.

" Each time you create a mini-batch to pass into the model, you randomly pick an image, you randomly pick amount of noise or timestep ( $t$ ) and then you use that amount of noise to each one.

Then you fix that mini-batch to model to train it. That trains weight in model so it can learn to predict results.

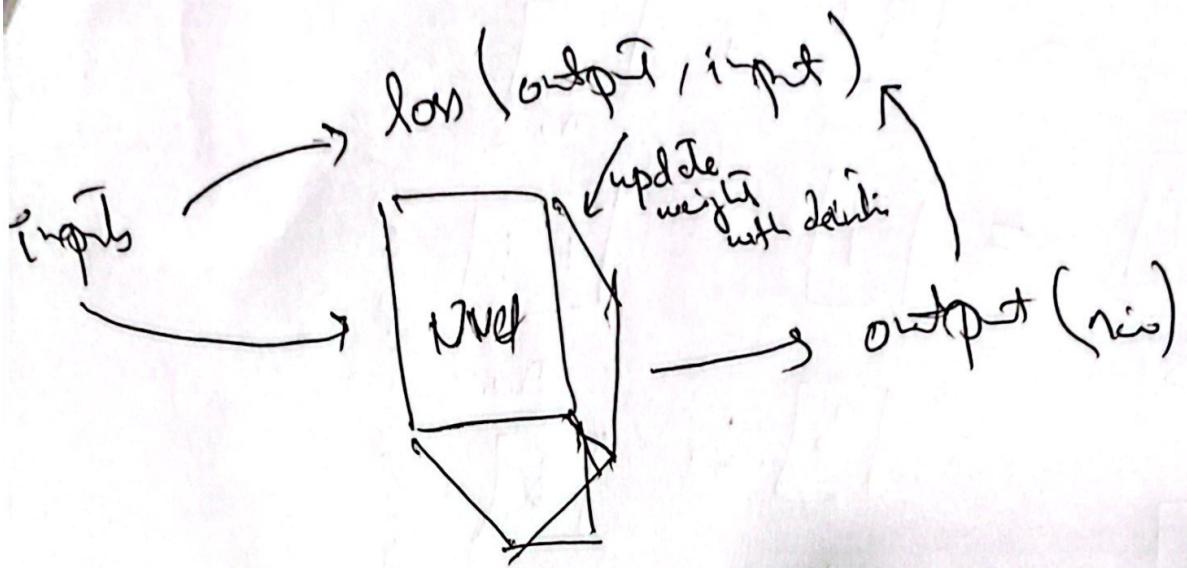
3. Let predict how much noise

$$\begin{array}{rcl} \boxed{9} & = & \boxed{9} + \boxed{\text{noise}} \\ \boxed{11} & = & \boxed{11} + \boxed{\text{noise}} \\ \boxed{12} & = & \boxed{12} + \boxed{\text{noise}} \\ \boxed{15} & = & \boxed{15} + \boxed{\text{noise}} \\ \boxed{16} & = & \boxed{16} + \boxed{\text{noise}} \end{array}$$

So the amount of noise is telling how much of a digit it is.

Something with less no noise is a very much like a digit and something with more noise is less like digit.

So amount of noise tells how much of like a digit it is.



$$x_3 = x_3 - \nabla_{x_3} p(x_3)$$

noise / 0.0001

$$\boxed{9} = \boxed{8} + \boxed{\square} \rightarrow N(0, 0)$$

$$\boxed{17} = \boxed{17} + \boxed{0} \rightarrow N(0, 0.1)$$

$$\boxed{100} = \boxed{3} + \boxed{97}$$

$\underbrace{\hspace{2cm}}$        $\underbrace{\hspace{2cm}}$

$\downarrow$                    $\downarrow$

Inputs to model                   $\rightarrow N(0, 1)$

So, why don't we use the variance to predict amount of noise or maybe even actual noise.

6)

So when you come to iterative.  
the Generative image from pure noise.

Mean now model is starting at  
 $\beta=2$  (max noise) and you want it  
to learn to remove noise.

Cursors ~~noise~~  $\rightarrow$

We can not do it in single step  
because, we don't have noisy images  
in dataset.

$$\text{noise} = \text{noise} - c \nabla p(\text{noise})$$

~~p(noise)~~  
noise  
 $c$  is hyperparameter

This do look like "Optimizers".

$\Rightarrow$  remember, if you change the same parameter  
by similar amount by multiple times in multiple  
steps, maybe you should increase the  
amount by which you change them.

→ momentum

→ adam

↳ what if  
variable  
changes.



Optimizer

→ Yes we can use these tricks.

→ If all come from differential Calc.

which is more about taking these little steps and trying to figure out how to take bigger steps.

→ Differential Eq solvers tend to take 't' as input. So ~~but not~~ Stable Diff models take  $(x_i, d, t)$

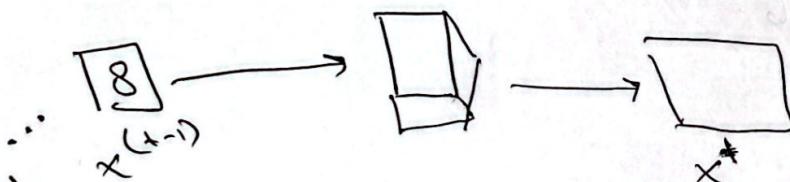
[  
-  $x_i$  - pixels  
- caption( $d$ )  
-  $t$  : timestep.

't' because model would runne more better if you know when it is.

(12)

Q11State Diff (Metho):
 $q_{x^{(t)}}(x^{(t)})$  - data distribution
↳  $x$  input variable. $x^{(t)}$  - instance or sequences of  $x$ .↳  $q$  - prob density function: outputs probability.
 $q(x^{(t)} | x^{(t-1)})$  - conditional probability

more formal would be,

 $q_{x^{(t)} | x^{(t-1)}}(x^{(t)} | x^{(t-1)})$ 

 $q_{x^{(t)} | x^{(t-1)}}(x^{(t)})$ 

or

 $q_{x^{(t-1)} | x^{(t)}}(x^{(t)})$ 
 $q_{x^t | x^{(t-1)}}(x^t | x^{(t-1)})$

So, what q

$$N(x^{(t)}; \underbrace{x^{(t-1)} \sqrt{1-\beta_t}, I^{\beta_t}}_{\text{Parameter of probability dist}}, \mu, \Sigma)$$

↓  
prob dist

Normal dist / Gaussian

↳ thin tails  
↳ mean you only need to describe their behavior in small region of space.

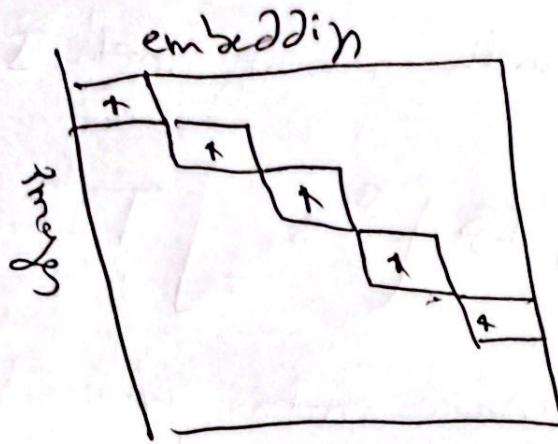
6 - for scale variable

$\Sigma$  - for covariance (multiple variables)

$\sigma^2$  - variance ;  $G$  - std deviate

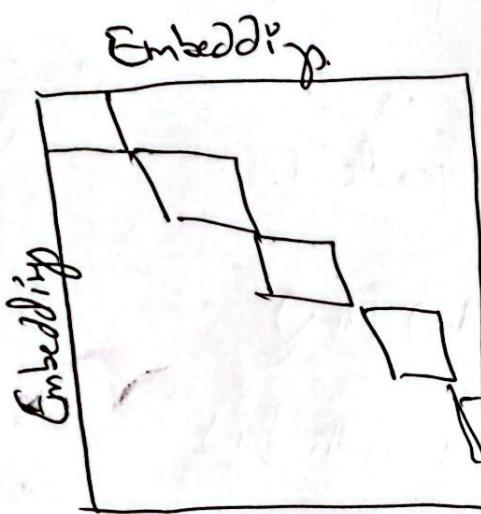
In CLIP

(14) (15) (13)



Above for embedding a both side

if we start  
near from  
diagonal entries

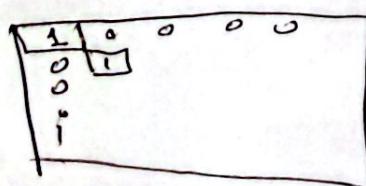


then diagonal  
are  
 $(X - E(x))^2$

↓ that is  
variance  
for each  
vector.

There dist  $q_{(r)}()$  is of  
each pixel diff.

so we are assuming  $c_i r_i = 1$ , so all other place  
 ~~$c_i r_i$  is the~~ across that pixels  
i.e.  $c_i r_i$  and  $c_i r_i$  are 0.



Var is spread from  
 $(x - E(x))^2$   
For all 1s at diagonal only  
var=0

so that mean all 0 result pixels are independent.

That's the assumption we are making,  
& if we start with identity matrix I.

$$I = [1 \ 0 \ 0], I\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

↑  
So, covariance b/w  
pixels 0  
b/w value &  
pixels (itself) is  
1.

so  $I\beta$  is a covariance matrix where  
for each individual pixel  $\beta_i$  is a  
covariance of that pixel and covariance  
(relation b/w pixels) is 0 & they're  
expected to be independent.

$$\beta_1 = 1 \text{ for (identity). (img)}$$

$\beta_{ij}$  is variance of pixel 1 with pixel 2

$$\left( \frac{x - E(x)}{\text{avg}} \right)^2 \text{ would be Variance.}$$

↳ for each pixel there would be different  
distribution

74

$$\text{mean} : \quad X^{(t-1)} \sqrt{1-\beta_+}$$

$$\text{covar} = I\beta_+$$

lets see how it behaves at edges (extreme cases).

$$\rightarrow \beta_+ = 0:$$

$$\begin{aligned} & \text{So, for } t=1 \\ \text{mean } (1) &= X^{(t-1)} \sqrt{1-\beta_+} \\ \text{mean } (1) &= X^{(0)} \sqrt{1-\beta_+} \\ \text{mean } (1) &= X^{(0)} \sqrt{1-\beta_+} \end{aligned}$$

$$\begin{aligned} & \text{mean } X^{(t-1)} \sqrt{1-\beta_+} \\ & = X^{(0)} \sqrt{1-\beta_+} \\ & = X^{(0)} \end{aligned}$$

| X  
example 2  
ke pixel

$$\text{mean}(2) = X^{(t-1)} \sqrt{1-\beta_+}$$

$$\text{mean}(1) = X^{(t-1)} \sqrt{1-\beta_+}$$

$$= X^{(0)} \sqrt{1-\beta_+}$$

$$\text{for } \beta_1=0 \Leftrightarrow \beta_+=0$$

:  $X^{(t-1)}$  is previous.

there will  
be other  
items for  
 $X^{(t)}$

$$\boxed{\text{mean } (t=1) = X^{(0)}}$$

- so we have  
mean of our previous  
img.

$$\text{var} = I\beta_+$$

$$\text{var}(1) = I(0)$$

$$\boxed{\text{var} = 0}$$

So, we have a normal dist  
with mean + previous img  
and variance  $\beta_+$  which mean we have  
noise.

Same image & var=0 mean no  
noise.  
at mean it gives image

when  $\beta_t = 1$

$$= X^{(t)} \sqrt{I-1}$$

$\text{mean} f_{t+1} = 0$

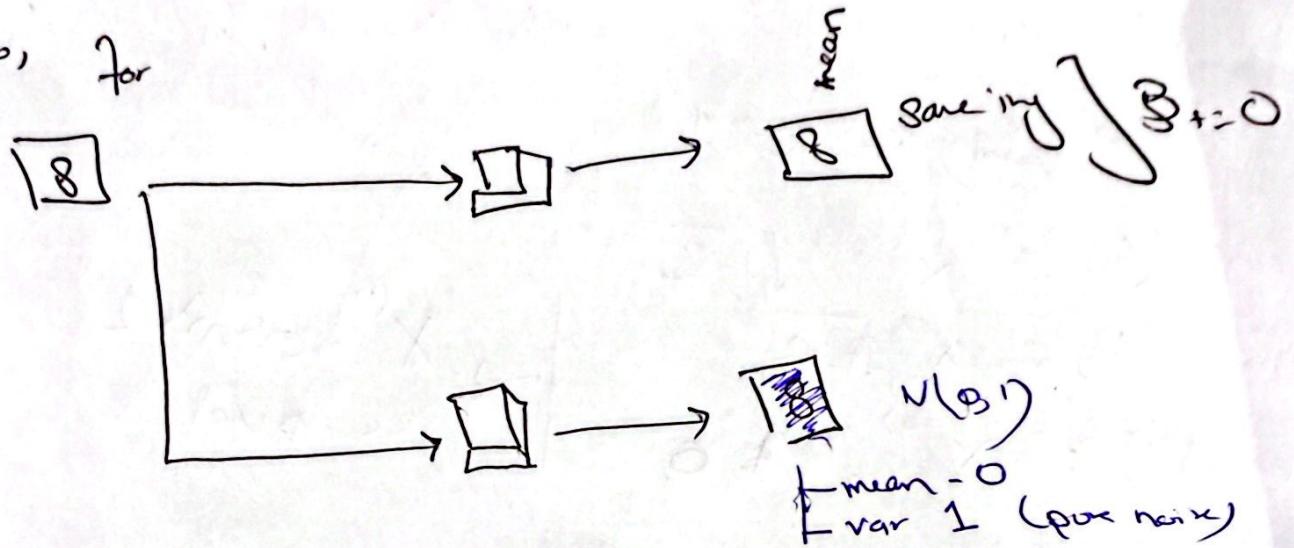
$$\text{var} = \sum \beta_t$$

$$\boxed{\text{var} = I}$$

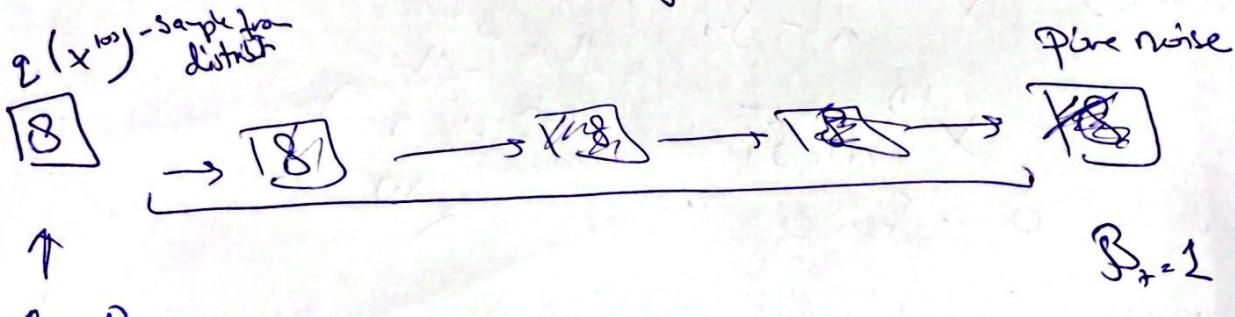
$\hookrightarrow I$  - identity matrix

mean variance of 1.

So, for



So, b/w these two cases there would be some mixtures (plus noise going)



$\rightarrow$  this is called Forward

$\rightarrow$  diffusion process,

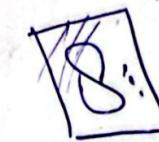
$\tilde{w}_t$

## Forward Diffusion Process



$$q(x^{(0)})$$

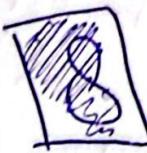
↳ sample from  
dist



$$\frac{q(x^{(1)} | x^{(0)})}{\text{conditional prob}}$$

.....

Rule  $\#$



This is called

Markov Chain Process

↳ transition depends on initial state

↳ Markov Process with Gaussian transition  
↳ noise

$$q(x^t | x^{t-1})$$

↳ transition is the  $f_{\theta}$ .

→ Sampling from distribution

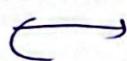
↳ trying to find some random data samples  
that maximizes the likelihood (or has high likelihood)

$$\prod P(x_i | \theta)$$

To modify normal distribution to any value

To get any particular noise/distribution =  $\mu + \sigma N(0, 1)$   
← constant distribution

$$N(\mu, \sigma^2) = \mu + \sigma N(0, 1)$$



Process of going backwards (same form but different parameters)

$$P\left(\overset{\text{previous}}{x^{(t-1)}} \mid \overset{\text{current}}{x^{(t)}}\right) = N\left(x^{(t-1)}; \text{mean?}, \text{variance?}\right)$$



q-Forward process	P-backward/reverse process.
-------------------	-----------------------------

P is what we're figuring out,  
~~q is input~~

To figure out  $\mu \& \sigma \neq P(x)$ ,  
we'll use TALB

(6)

Find  $\mu \{ \sigma \}$  of  $P(\cdot)$  by MLE.

Likelihood is not diff

But we can do thi exactly b/c we'll need to do some integrals that'll be dimensional (continuity) So, you can't do thi.

There are thousand of steps & there will be many possible values, in the reverse p/p process, so it will be tough to evaluate it for over all these steps for all possible values.

So, log likelihood:

$\log L$  is for scaling.

monotonic fn.

$\rightarrow$  product  $\rightarrow$  sum

$\rightarrow$  Normal dist have exp function that disappear with log.

we'll still can not optimize log likelihood

directly because it is not differentiable w.r.t. of parameters of models, so we can not use gradient descent

$\Rightarrow$  Reason for ELBO is used in stable diffusion; (Evidence lower bound) with log likelihood

↳ Two main aspects

- ↳ Regularization
- ↳ Tractability.

1- ELBO in regularization: without ELBO, maximizing the log likelihood alone can lead to overfitting.

So, ELBO introduces a regularization term that penalizes complex models, encouraging them to learn simpler and more generalized representations.

This is achieved by adding the KL div b/w the prior and posterior dist to the objective function.

1500 AED

## 2- Tractability.

↳ Directly optimizing the log likelihood can be computationally expensive, especially for large models and complex details.

ELBO provides a tractable way to approximate log likelihood by introducing a latent variable (the noise variable) in stable diffusion and performing variational inference. This allows the model to learn a lower bound on the true log likelihood.

7)

Imagine you have a huge pile of dirty dishes to wash. Washing each dish individually would take a long time & would be inefficient. Instead, you might use a dish washer to clean them all at once. This is similar how ELBO works in state diffusion.

Directly optimizing the log likelihood, like washing each dish by hand, is accurate but slow and computationally expensive.

ELBO is like using a dish washer. It introduces a shortcut by using a latent variable (think of it as a special container for the dirty dishes).

This container holds all the info about the ~~ing~~ in a compressed and manageable ~~form~~ way.

Then, ELBO uses a technique called variational inference to estimate the log likelihood

based on this container. This is much faster and easier than calculating log likelihood for each individual ~~ing~~. The log likelihood is the ideal result, I want to achieve. ELBO might not give you the exact result but it provides a 'lower bound' which means it guarantees that the final output will be at least as good as ~~as~~ ~~and in most cases lower bound is good~~

So, instead of directly calculating the log likelihood for every possible  $iy$ , ELBO introduces a noise variable that represents the image in a simplified way. This allows the model to learn about the rules of good image, without having to generate & analyze every single possible  $iy$ .

Think of it learning how to draw a portrait. Instead of trying to draw every single detail of a face at once, you might start by sketching a basic outline and then add in more details.

ELBO is like that first outline, allowing the model to learn basic structure of  $iy$  before focusing on finer details.

18)

## QB (contine)

we do Jodh leam  
B2

for backward process, we are using same  
but parameters are different.

$$= N(x^{(+)}; \text{mean?}, \text{var?})$$

but we can not use to compute by  
likelihood directly to find these parameters, we  
need to use ELBO (Evidence lower bound).

ELBO can be calculated easily so, we  
can update our function to predict  
mean and variance of the reverse process.

so important to clarify: You can recover whole  
distribution but you can not a single step,  
convert it to sparse noise and recover it.  
So, it's on distribution level.

So, approach ELBO is like:

As we have forward process (inferred to dist.)  
going all the way to pure noise. with

ELBO we are trying to match the our  
distribution that we are trying to optimize to  
that distribution that we saw in forward  
process.

KLDiv is to compare two distrib

So, we are trying to minimize the difference  
b/w forward process dist and inference process  
dist.

---

DDPM (2020), (muy like Jerry)

'Denoising Diffusion Probabilistic Model'

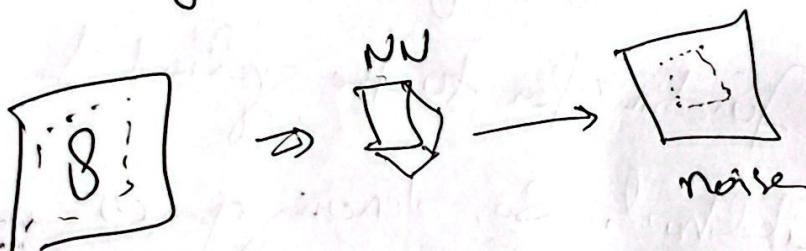
↳ Let's assume that  $\text{var}(\text{in noise process})$  is  
just a constant & we do not learn it.

And we do assume that step size from earlier  
( $\beta_+$ ) - the variance of the noise we added at each step  
is also constant. So, we are just predicting the  
mean - there are set constant values.

18) then the loss turns out to be, that you predict the noise.

loss  $\Rightarrow$  Noise

So, you need to train NN that takes an image & it tells you what I think my noise.

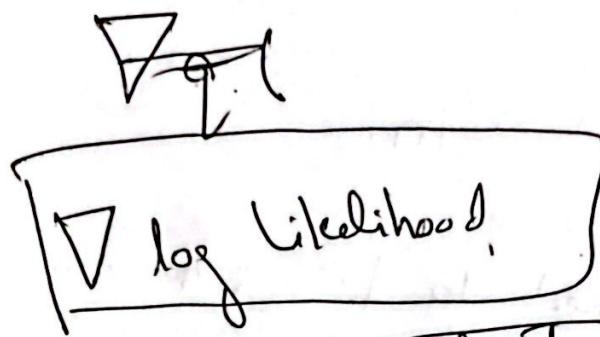


$\Rightarrow$  The idea of adding & removing noise from from distribution is that adding noise to the images we have, that takes us away from  $\hookrightarrow$  (max likelihood).

our real images to noisy images (<sup>decreases</sup> low Likelihood.)  
So, we want to learn how to get back to high Likelihood images.

The Denoising process is equivalent to learning score function. Score function is gradient of log of the likelihood.

Score function of log of likelihood.



~~Score Function particularly refers to  $\nabla \log$  of gradient of log of likelihood~~

Score Function refer to the gradient of  $\log$  Likelihood. So, denoising process allows us to learn the score function.

So, when we are doing the noise prediction, we know we have a probabilistic framework w.r.t. the likelihood and it comes all down to predict the noise.

Now what DDPM proposed is its equivalent to calculate the score fn

With DDPM - you just predict noise

Noise

20)

before #10

(Distillation)

To reducing steps from  $60 \rightarrow 4$ .

Distillation:

↳ You take a teacher network that knows how to do something but is slow and big and then teacher is used by a student network which learns to do same thing but faster or with less memory. Here we want our to be faster  
60 steps:

16

12

18

24

30

36

42

48

54

60

It is taking about 18 steps to go from [42] to [60]. It seems like something that can be done in single step.

So, think of a model that takes input  $\text{img at } \boxed{t=2}$  and outputs  $\text{img at } \boxed{t=60}$  excluding all intermediate steps



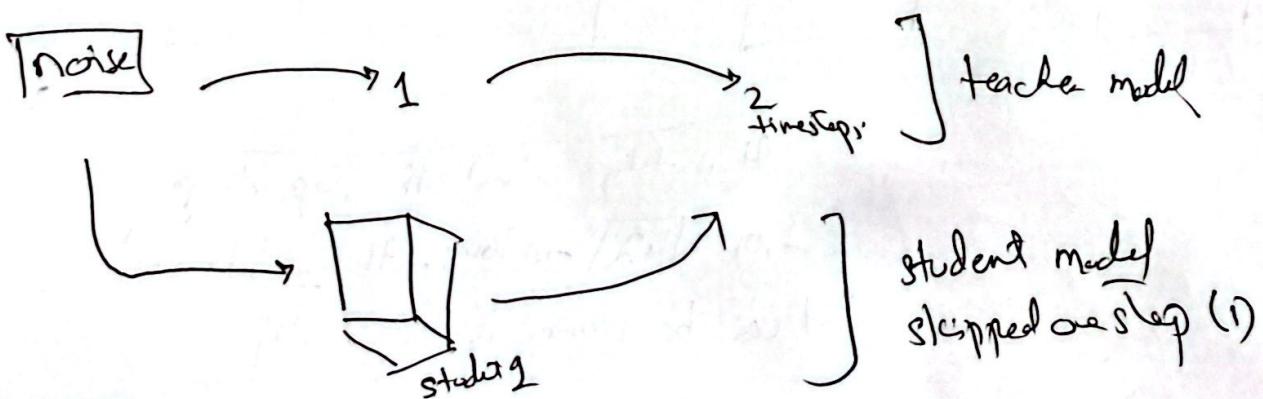
compare  $\hat{\boxed{t=60}}$  predicted output with original output from teacher model  $\boxed{t=60}$ .

$$\text{loss} = \text{MSE}(\boxed{t=60}, \hat{\boxed{t=60}})$$

So, this UNET will learn how to do it correctly.

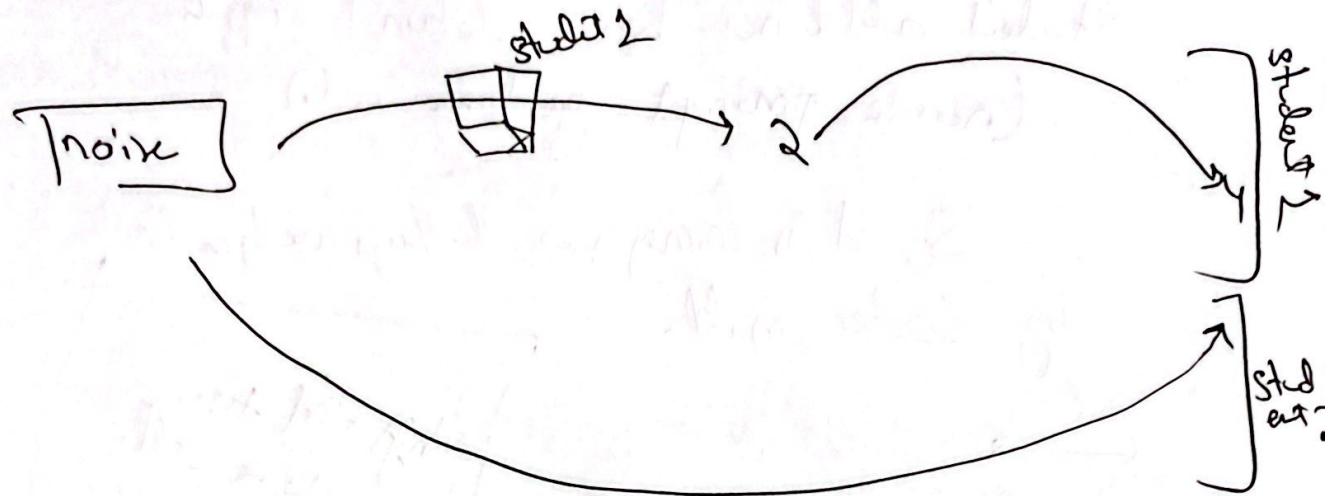
teacher model - already heavily trained model

Timedep is adding noise to style img (here we add noise into style img in different steps)



TIMESTEPS ARE JUST GIVING THE MODEL SAME INPUT AGAIN.

2) Now treat student<sup>1</sup> as a teacher.

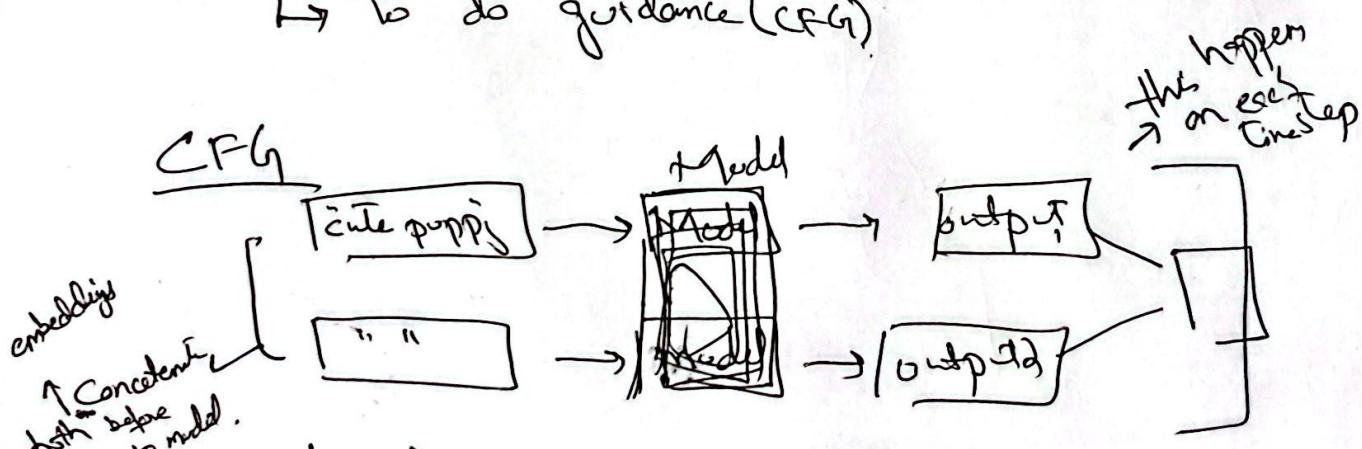


we treated student<sup>1</sup> as a teacher.  
and student<sup>2</sup> as a student & go on.



## ON DISTILLATION OR GUIDED DIFFUSION MODEL.

↳ To do guidance (CFG)



In CFG we take weighted avg of  
output<sub>1</sub> and output<sub>2</sub>.

[ Just like previous paper but add embeddings to.  
here ]

S<sub>3</sub>  
Student model now has additional inputs  
(noise, prompt, guidance scale).

S, it is learning how all things are handled  
by teacher model.



JMDGIC

Shift + M to  
combine cells

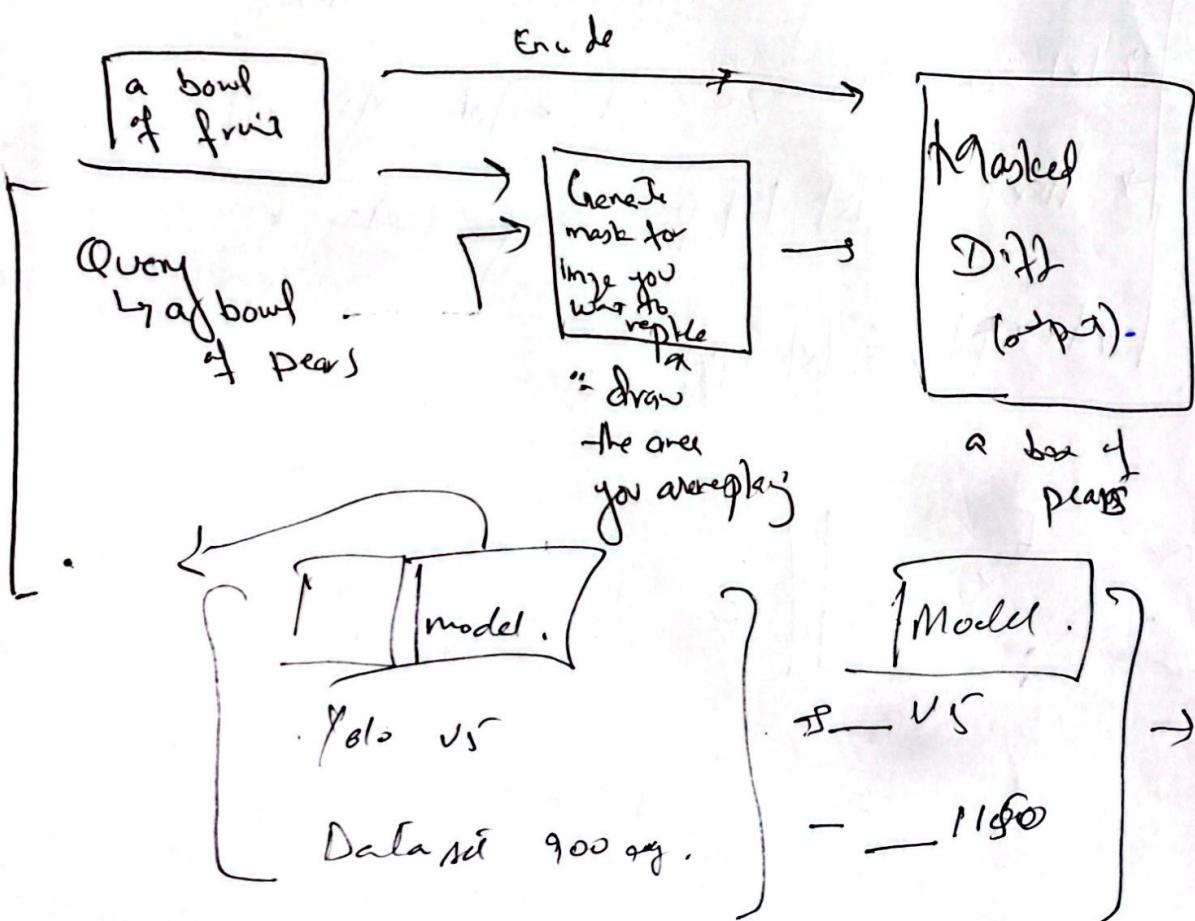
22)

## Build from foundations

### DIFF EDIT

#### Semantic Image Editing

Extension of my Generator. We'll edit image based on text query. It is Semantic Image Editing with additional constraint that generated image should be as similar as possible to given ~~opt~~ input.



Norm - to get  
non-ve values

| - vector norm  
 $\| \cdot \|$  - matrix norm

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{1/p}$$

$$e \sim N(0, I)$$

I. identity  
matrix  
mag 1

So far

$$\|x\|_2 = \|x\| = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

$\|x\|$  is a non-negative norm defined such that

$\rightarrow \|x\| > 0$  when  $x \neq 0$  and  $\cancel{\rightarrow} \|x\| = 0$  iff  $x = 0$

$\rightarrow \|kx\| = \|k\| \|x\|$  for any scalar  $k$ .

$\rightarrow \|x+y\| \leq \|x\| + \|y\|$   $\because x$  or  $y$  can be  
-ve so  $\leq$  sign

$$h = \underbrace{\sum_{x_0, t, \varepsilon} \left\| \varepsilon - \varepsilon_0(x_0, t) \right\|_2^2}_{g + \text{just mean sum of squares}}$$

$\therefore g + \text{just mean sum of squares}$

of  $\varepsilon - \varepsilon_0(x_0, t)$

Expected value  $E_p(x)$   $\downarrow$   
 $\therefore$  (value)(prob of that value)

23)

S,

$$L = \underset{x_0, t, \epsilon}{\mathbb{E}} \left\| \epsilon - E_0(x_t, t) \right\|_2^2$$

noise  
noise estimator (Prediction of noise)

$E_0$  is the noise estimator which aims to find the noise  $\epsilon \sim N(0, 1)$  that is mixed with an input  $x_0$  to yield  $x_t = \sqrt{\alpha_t} * x_0 + \sqrt{1-\alpha_t} \epsilon$ .

$$\alpha_t \in [0, 1] \quad \begin{matrix} \text{(works} \\ \text{as a} \\ \text{mask)} \end{matrix}$$

coeff  $\alpha_t$  defines the level of noise and is a decreasing fn of Time step  $t$ , with  $\alpha_0 = 1$  (no noise) and  $\alpha_{t \rightarrow \infty} \approx 0$  (almost pure noise).

So, we got noise  $\epsilon$ , predict  $\{x_t\}$  the noise  $E_0(x_t, t)$ , we subtract one from other and square it to get the expected value. In other words its MSE.

$$\text{here } x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \epsilon$$

↳ reduce value of each pixel and noise to it.

The effect of mask is that original pixels in black areas won't be changed i.e. ignored. So, the scenery would be same.

$$= \underbrace{M_t}_{\text{mask}} + \underbrace{(1 - M_t)}_{\text{Expr.}} \underbrace{x_t}_{\text{mask}(x_t)}$$

\*  $x_t$  = image  
 Inversion  
 $1 - \text{mask} = \text{horse}$

$\downarrow$

new obj (i.e. zebra) zebra world be discard

$\rightarrow$  mask binary  
 $M \in [0, 1]$

only horse/zebra      only scene with object contact

mask is binarized & it is actually a bit bit bigger than the object i.e. it loosely captures object. That is fine.

24)

Automatically Generate a Mask].



Matrix Multiplication  $\rightarrow x_0, i \rightarrow$

$$h = E \underset{x_0, t, \epsilon}{\|} E - E_0(x_+, +) \|_2^2$$

$$E \underset{x_0, t, \epsilon}{(E - E_0(x_+, +))^2}$$

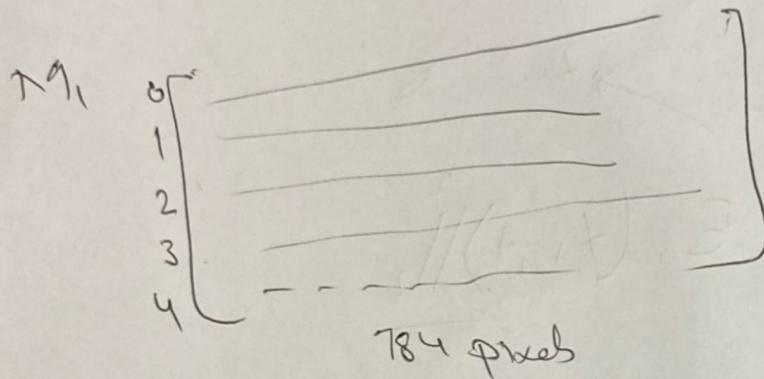
$$= \frac{\{(E - E_0(x_+, +))^2\}}{n}$$

$$\epsilon \sim N(0, I)$$

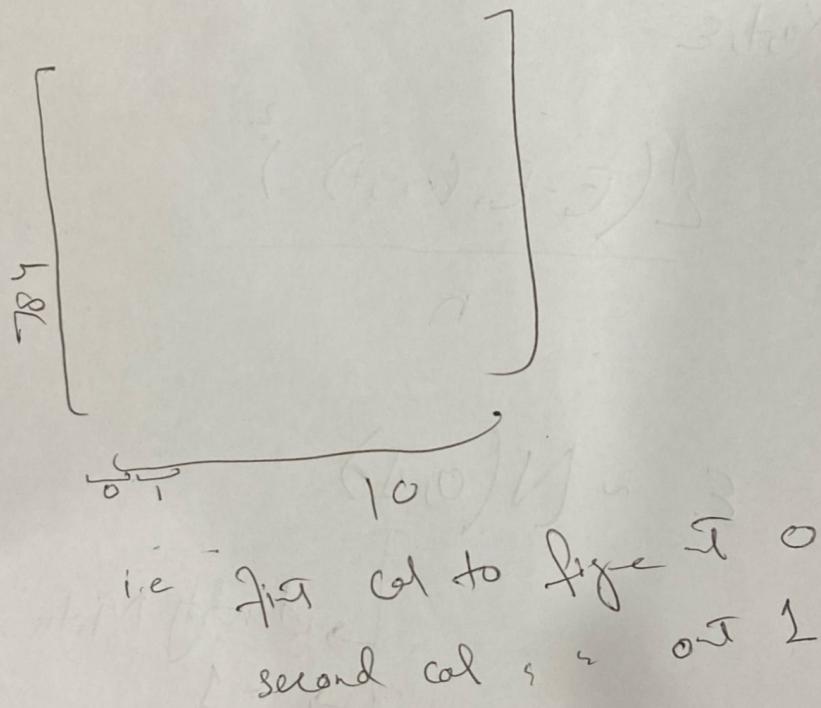
$\left[ \begin{array}{l} \text{Identity Matrix} \\ = 1 \end{array} \right]$

# Matrix Multiplication

$M_1$  : matrix (first five MNIST digits) =  $5 \times 784$



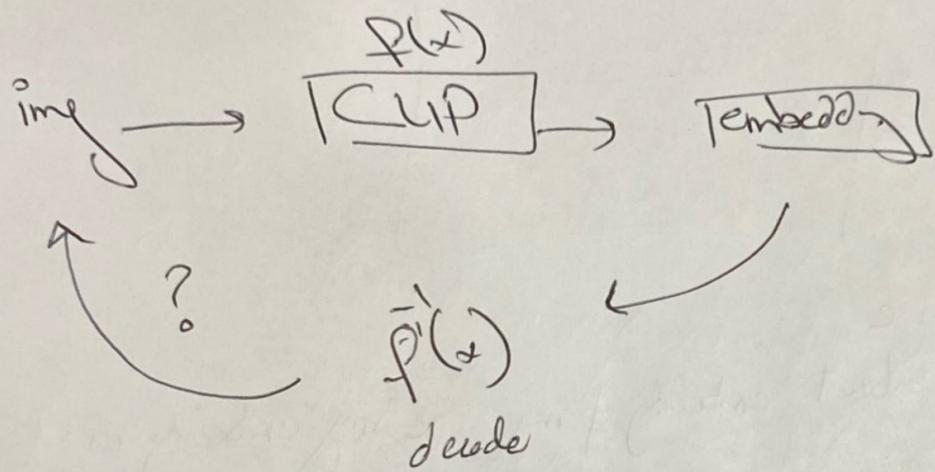
$M_2$  : matrix (weights) =  $784 \times 10$



$$\Rightarrow M_1 \times M_2$$

Decoding CLIP | Can we get text prompt <sup>from</sup> generated img? No

→ You got CLIP embeddy of an image  
can you decode it back to img.



Think of def  $f(x)$ :  
return 0

could you get  $x$  back? No

but what if we could put it through diffusion process?

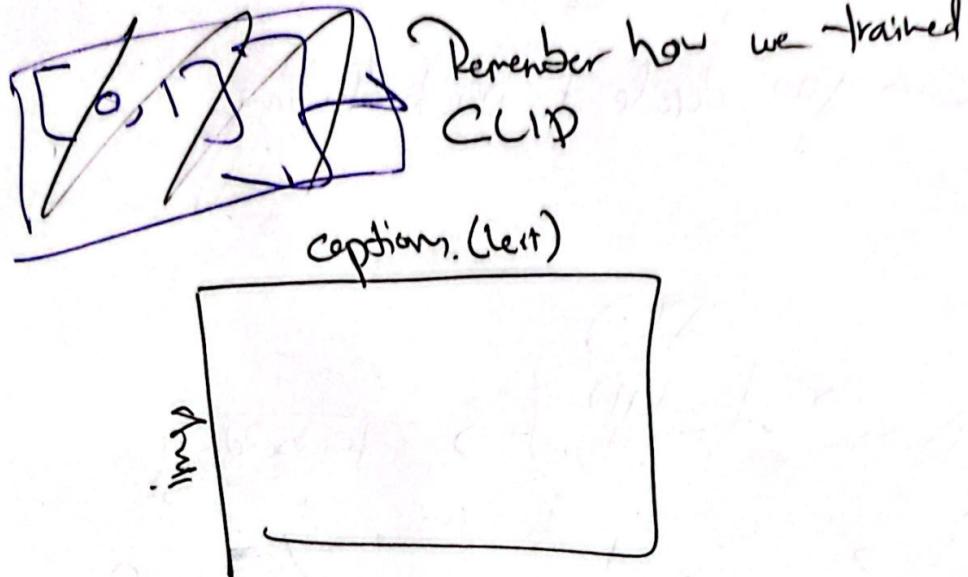
so replace

$f^{-1}(x)$  with  $\boxed{\text{noise}} + \boxed{\text{embedding}}$

and we can pass it n times.

It might be - bit similar but not exact

But in prompts we don't get "image embeddings" instead we get "text embeddings", but that does not matter.



i.e

text embedding of my  $\approx$  image embedding

CLIP embedding = Text image embedding of this picture  
should be similar to actual embedding of my.

So  $f^{-1}(x)$  is basically inverse embedding fn.

These kind of problems are generally referred to as  
inverse problems.

So, Stable diffusion attempts to approximate  
the solution to an inverse problem.

27) So, CLIP injector is not actually 'inverted'  
the picture to give us back the  
"INVERSE PROBLEM"

If we got an image embedding, try to  
undo that to get back to picture & try to  
undo that to get back to a suitable prompt  
prompt is equally infeasible. Both of  
them require inverting an encoder that  
doesn't exist.

All we can do is approximate  
that big difficult process.

[CLIP is not inverse of CLIP  
encode  
↳ captioning]

## Einstein Summation

Here we are using indices not only as subscript but also superscript.

$a_i, a^i \dots i = \text{index}$

$$\left[ a^i \times a^i = (a^i)^2 \right] \text{ "2 power."}$$

$a^i$  is a superscript  $i$

In tensor calculus, a lot of expressions include summation over particular indices.

$$\sum_{i=1}^3 a_i x_i = a_1 x_1 + a_2 x_2 + a_3 x_3$$

In einstein summing we don't write summation ( $\Sigma$ )

$$\sum_{i=1}^3 a_i x_i \xrightarrow[\text{summation}]{\text{in einstein}} a_i x_i$$

Index refers to indices of a tensor.

Einstein Notation Rule:

Rule 1: If you have twice-repeated index in a single term is summed over.

Index: 1, 2, 3, ..., n

typically n=3 :: because of three dimensions

$$a_{ij} b_j = a_{i1} b_1 + a_{i2} b_2 + a_{i3} b_3$$

$\downarrow$   
 $j=1, 2, 3$

Where j is a dummy index because it is repeated twice in a term.

i = free index (can take any value that j takes on,

e.g. i=1, 2, 3 but only one value. If is not summed over).

So, I can be 1, 2 or 3 but it can only be one of these values.

j occurs only once in expression but it can not be replaced by another free index. Dummy index can be replaced.

$$a_{ij} b_j + a_{kj} b_j$$

Dummy index is summed over

→ appears twice

→ can be replaced by another dummy index.

Rule 2

Free index is not summed over

→ occurs once

→ can not be replaced by another free index

Rule 3

→ No index can occur 3 or more times in given index term

$a_{ij}b_{ij}$  ✓       $a_{ii}b_{ij}$  X

$a_{ij}b_{jj}$  X

we can have only three indices that occur once or dummy indices that occur twice in a single term. So, no index occurs more than twice.

## 28) "Practical Tensor Manipulation"

### MEAN SHIFT

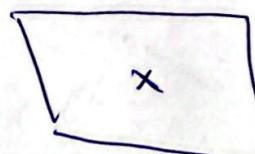
→ CLUSTERING:

↳ alternative to k-means clustering theorem

For each datapoint, take weighted avg from all other data points in sample X.

Remember avg of all points would

be centre of

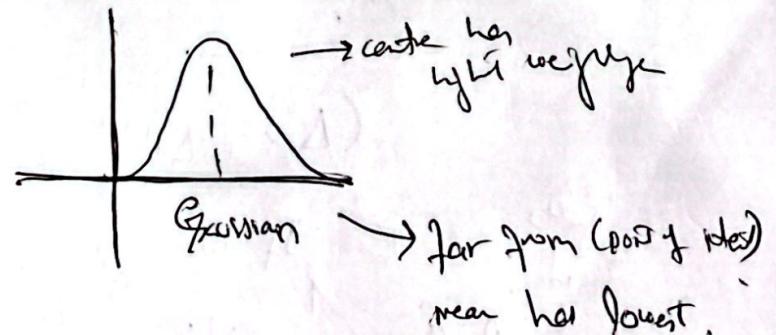


So, we

are taking weighted avg.

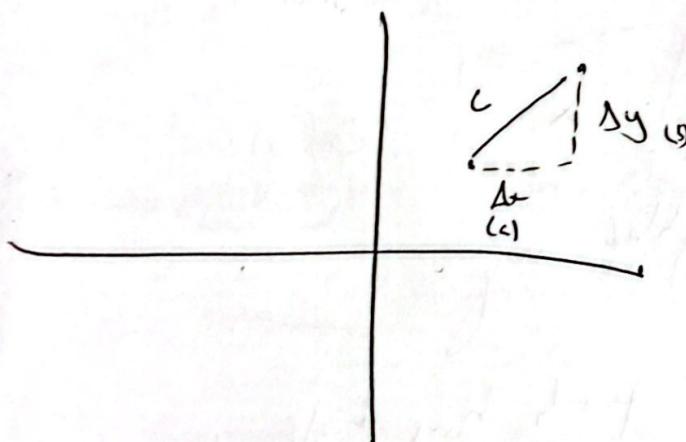
→ closer data points have higher weight & further away has lowest.

→ So, we create weights for every point compared to one we're interested in using Gaussian Kernel.



# If you have shorter two shapes of diff lengths, we can use shorter length if it will add unit axis to the front to make it long as necessary.

## Normals :



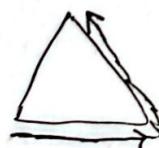
using Pythagoras

$$a^2 + b^2 = c^2$$

$$c^2 = (\Delta x)^2 + (\Delta y)^2$$

$$c = \sqrt{\Delta x^2 + \Delta y^2}$$

but if without using hypotenuse we first cover Δx and then Δy

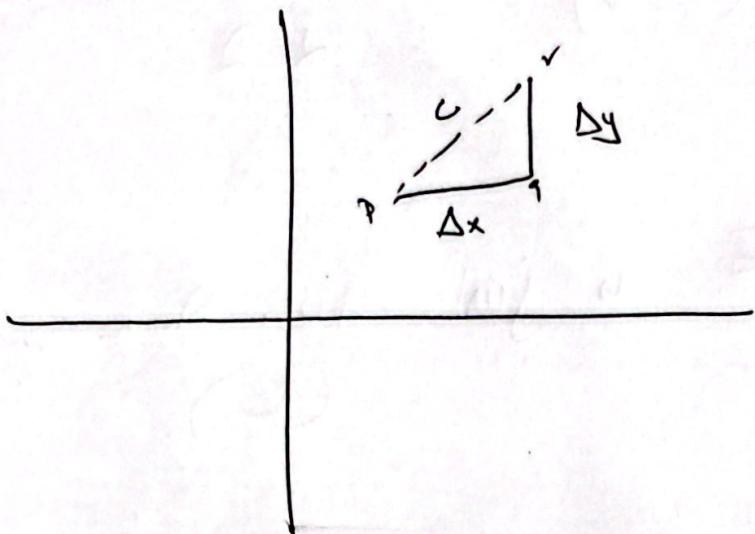


$$(\Delta x' + \Delta y')^n \approx \Delta x + \Delta y$$

So, if we get list of numbers we can add them up for each power of  $\Delta$  & take that sum & then Norm.

$$\left( \sum v^\alpha \right)^{1/2}$$

Norms



to find  $c$ , using pythagoras,

$$a^2 + b^2 = c^2$$

$$c = \sqrt{\Delta x^2 + \Delta y^2}$$

but what if we just cover  
 $\Delta x$  first then  $\Delta y$ , then its sum  
would be -

$$c = (\Delta x + \Delta y)^{1/2} \propto \Delta x + \Delta y$$

So general of norm is

$$\left( \sum |V|^{\alpha} \right)^{1/\alpha} \quad \text{where } V \text{ is vector / list of numbers.}$$

↳ because dist can't be -ve

$$\|v\| = \sqrt{\sum v_i^2} \quad \therefore \alpha=2$$

$$\|v\|_2 = \left( \sum v_i^2 \right)^{1/2} \quad \therefore \alpha=2$$

$\Rightarrow$  RMS is just two-Norm ( $L_2$ )

Root Mean Square Err.  $\rightarrow$  Euclidean Distance

whereas

MAE is just 1-norm ( $L_1$ )

Mean  
Abs  
Error -

So,  $\|v\|_2^2 = D_x^2 + D_y^2 ]$  or Seen  
in prev  
paper

Generally weighted avg is:

$$\frac{\sum x_i}{n}$$

if weighted avg is

$$\frac{\sum w_i x_i}{\sum w_i} \quad \text{"divide by sum of weights"}$$

$X = [1500, 2]$

weight = 1500

In Matrix Multiplication, filters are compatible if

- One of them is 1
- they are equal

$$\left. \begin{array}{l} [1500, 1] \\ [1500, 2] \end{array} \right\} \rightarrow \begin{array}{l} \text{One } 1500 \text{ will match} \\ \text{Two } 1500 \text{ are equal} \end{array}$$

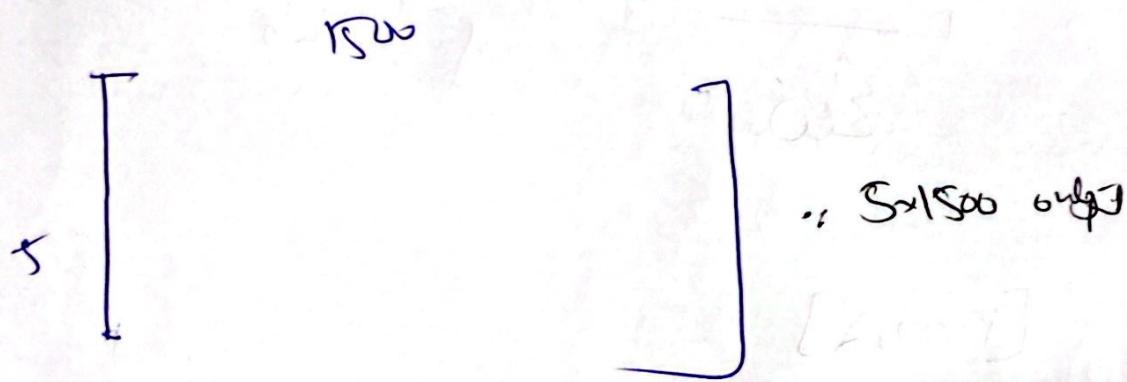
In case of min batch

$$\text{batch size} = [S \times 2]$$

$$i.e. X = [S \times 2]$$

the distance will be b/w every  $\in [bs]$  with  
every instance in  $X$  and vice versa so

result matrix distance metric will be



$$A[\text{None}] = [1, 1500, 2]$$

↳ one demand  $\Rightarrow [1500, 2]$

# Calculus

$f(t) = \text{Distance}$

$f(t)$ :

$$\frac{f(t_2) - f(t_1)}{t_2 - t_1} = \frac{\Delta d}{\Delta t}$$

but what if slope is concave.



here comes limit.

$$\frac{f(t_1 + \Delta) - f(t_1)}{f(t_1 + \Delta) - f(t_1)}, \Delta \text{ is a small value}$$

$$\Rightarrow \frac{f(t_1 + \Delta) - f(t_1)}{\Delta}$$

if  $\Delta \rightarrow 0$

$\frac{dy}{dx} \rightarrow y \text{ is result of the small change}$

$\left. \begin{array}{l} \downarrow \\ \text{is just a very small number} \end{array} \right. \begin{array}{l} \downarrow \\ x \text{ is very small number} \end{array}$

LIMITS & DERIVATIVE  
CALCULUS

## FORWARD & BACKWARD PASS:

lets consider a linear model of

Pick a single pixel from MNIST digits and that would be our  $x$

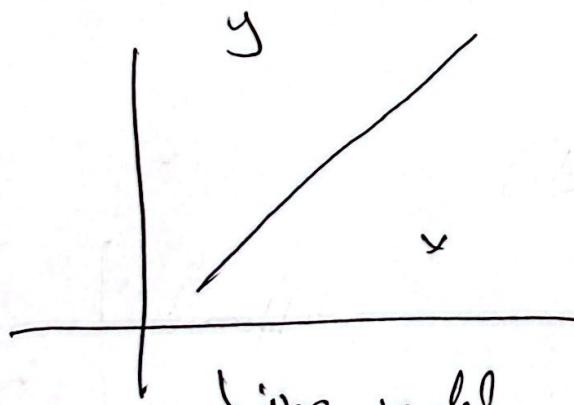
$y$ -value is how likely it is #3 based on value of this one pixel.

— pixel value is  $x$

— prob of it being 3 is  $y = P(x)$

For linear model it says, the brighter the pixel is, the more likely it is that its number 3 we are trying to see. i.e. the

but go get that from linear model

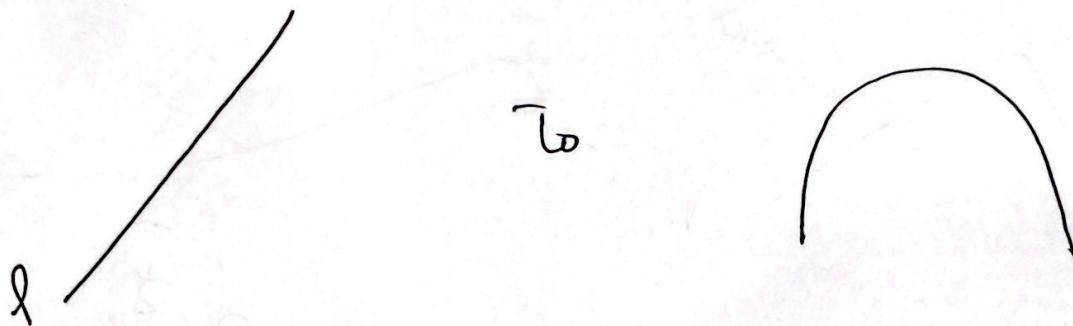


convert straight line to curve.

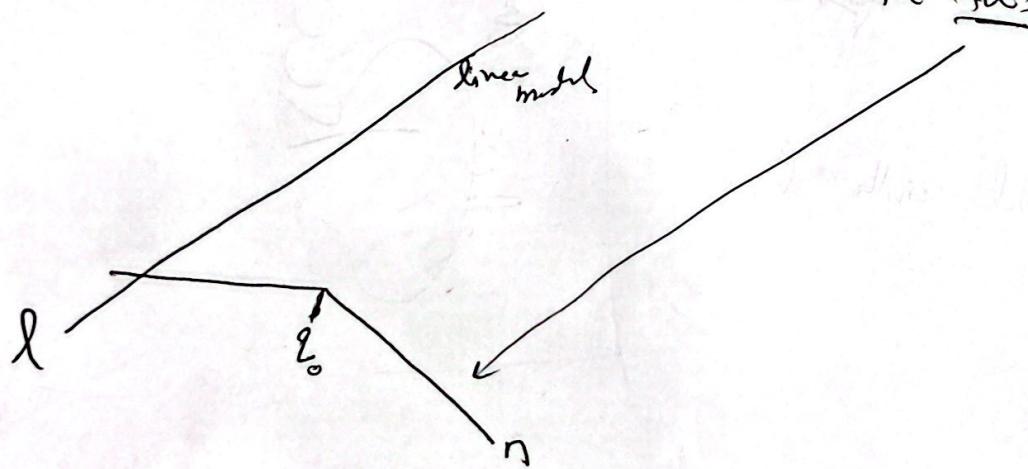
33)

we can use activation functions for that!

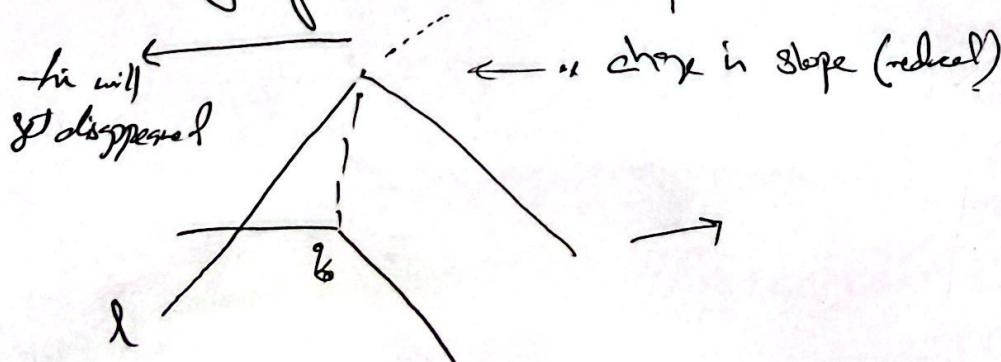
convert



Instead we can create a line like this,

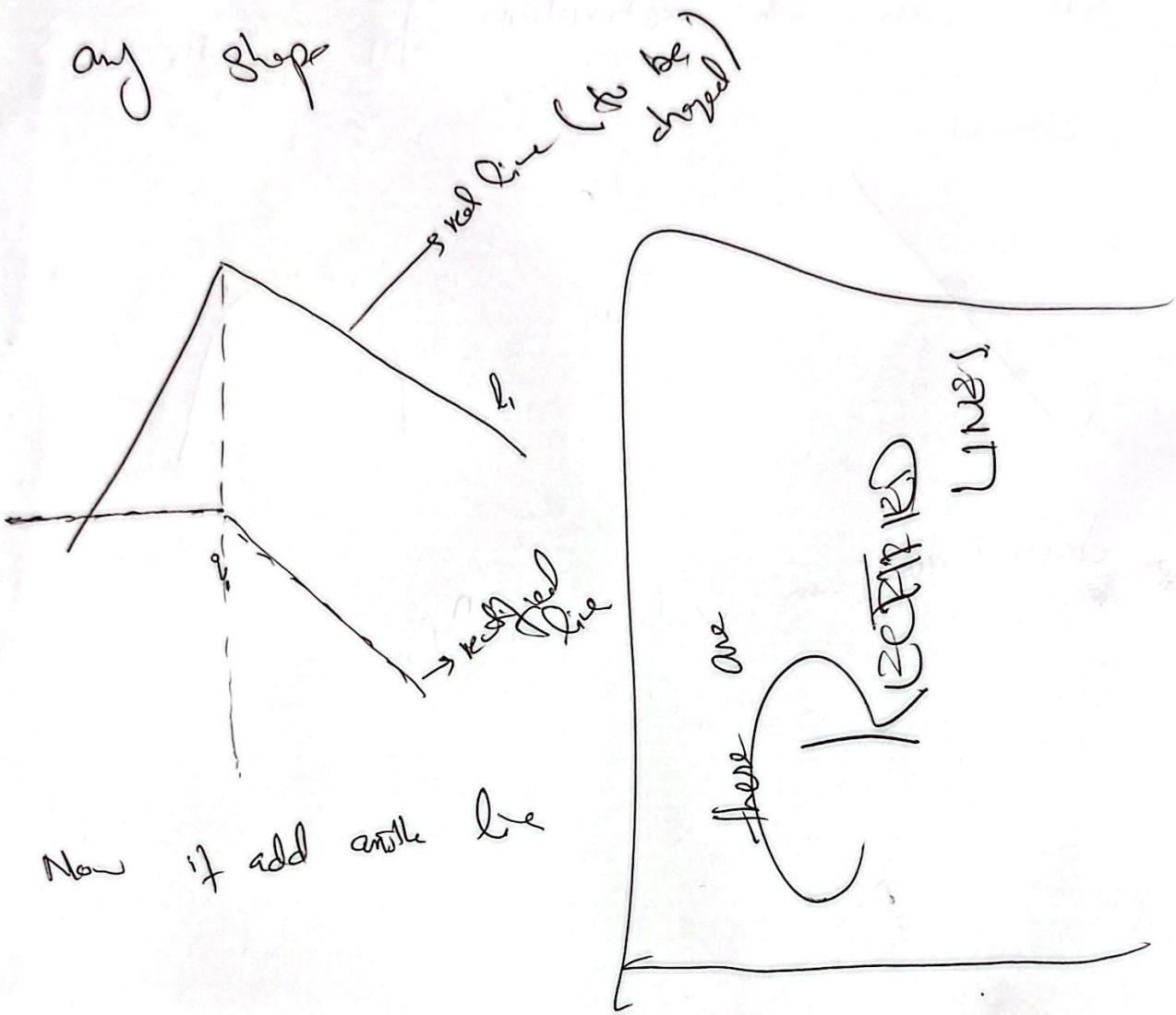


Energy to the left of point  $g_0$ , is not going to change if you add these two lines together because this is zero - but enough on right of it is going to be reduced.

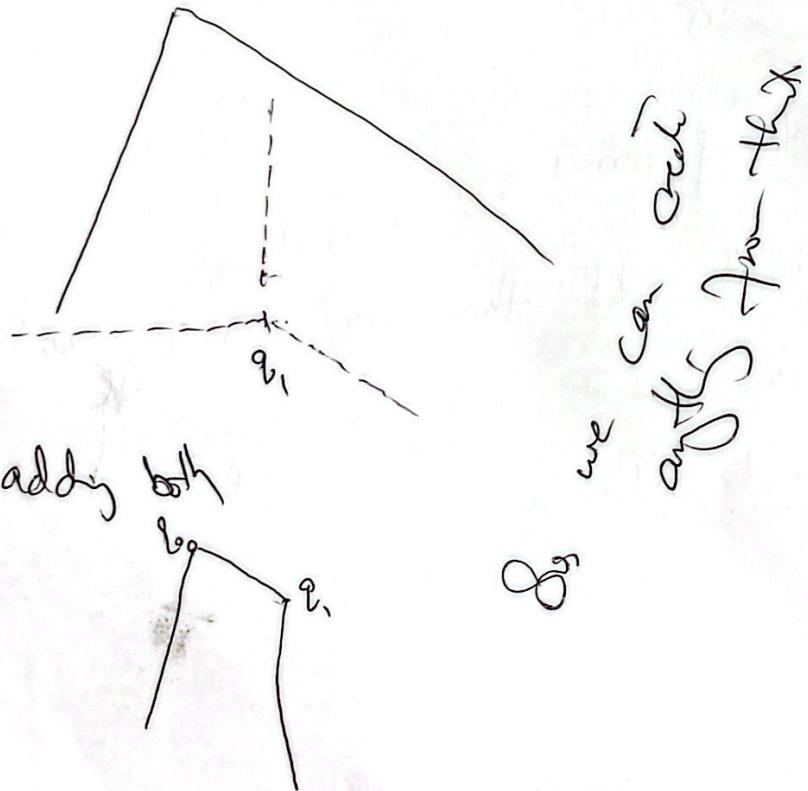


we can keep doing this to get

any shape



Now if add another line



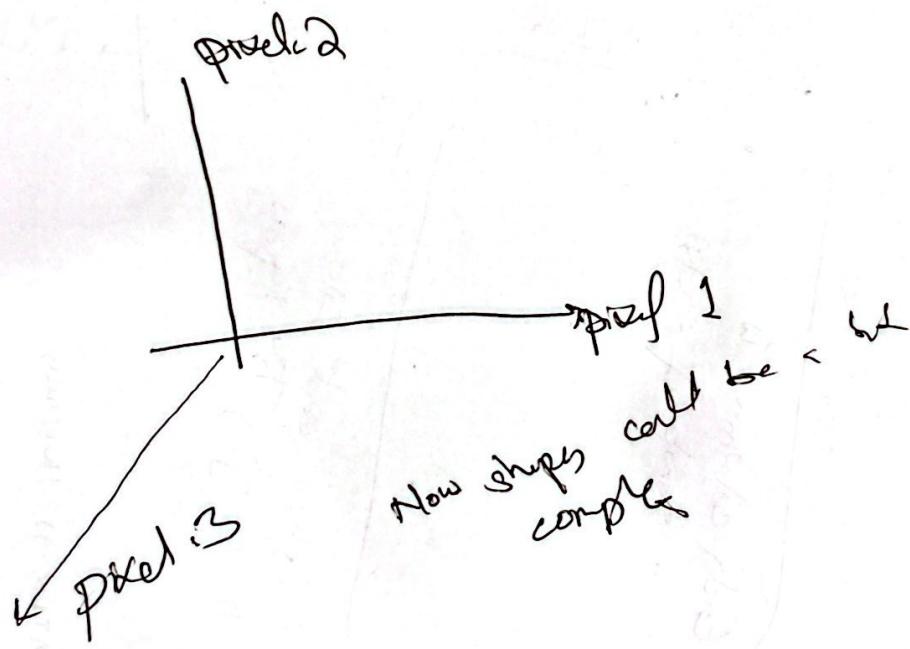
34)  
a)

we are going to have lots of pixels,

Adding n rectified lines would make  
it rectified plane.

↳ If its smaller/less than zero, truncate  
it to zero.

We are not going to have just one pixel  
but many lots of it



hidden layer = # of Rectified lines to add up  
activation

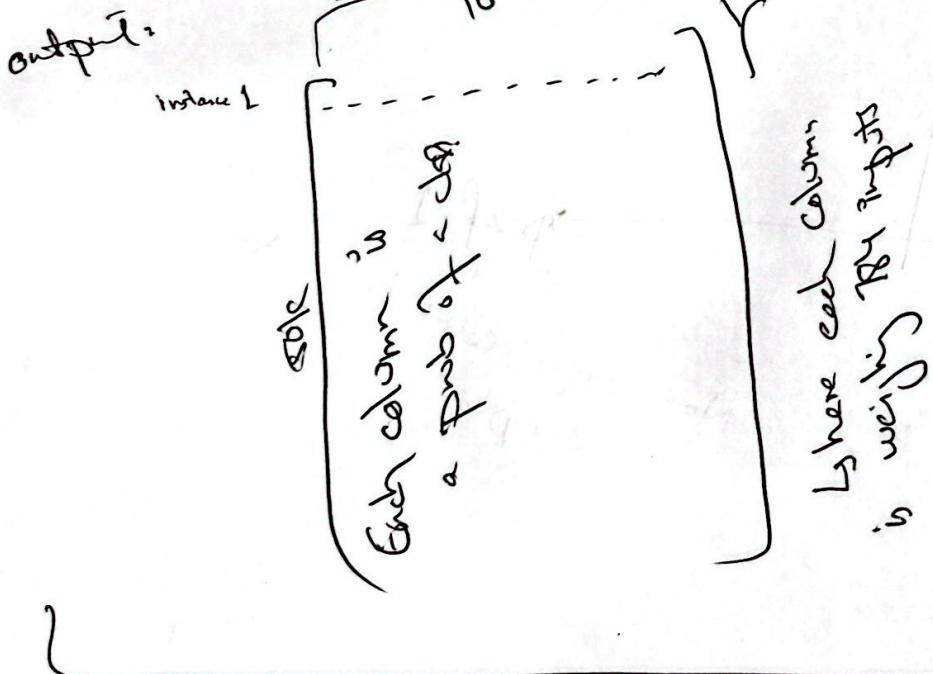
without hidden layer

$$\text{size } \begin{bmatrix} X \end{bmatrix} @ \begin{bmatrix} W \end{bmatrix}^{(784, 10)} = (505784) \quad \text{bias}$$

first column  
by putting 0

↳ where  
each col  
will be  
part of  
by a class

$$\begin{aligned} &\rightarrow \text{i.e.} \\ &\rightarrow 1^{\text{st}} \text{ col} \\ &\approx 40 \\ &\rightarrow 2^{\text{nd}} \text{ col: } 5 \end{aligned}$$



This was linear model

Weight matrix is initialized  
filled with random values.

34)

b) linear:

$$\text{res} = x @ w + b$$

$$(\text{res} - y) =$$

Gradient (slope)

$$f(x, w) \quad - \text{model}$$

$$\text{loss}_j = L(f(x, w), y_j) = \text{loss function}$$

If we could get derivative of loss wrt one particular weight

$$\frac{\partial \text{loss}}{\partial w_0}$$

"If I increase weight a little bit what happens to loss"

- If it increases loss, I want to decrease it
- If it decreases, I want to increase it

SGD

$$\text{SGD} = \theta_j - \eta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

for

multiple inputs  $x_{i,j}$  250 outputs

so  
derivative  $\left[ \frac{\partial z}{\partial y}, \frac{\partial z}{\partial x} \right]$

↳ for 784 inputs, we'll have  
784 of these

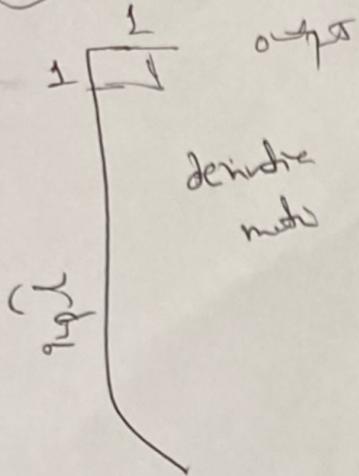
shaded with  $\rightarrow$

$$\frac{\partial L}{\partial w} \rightarrow \text{for all weights}$$

inputs       $\frac{\partial p_i}{\partial w_j}$   
784 pixels      weight

So, it will be huge matrix of derivatives,  
that says for every input that you change, by  
a little bit, how much does it change every  
output.

35)



& (1) is this matrix will  
will have changing output  
input I will smooth down  
output. b1

The well at end need a scalar. (Max)

2. in organized way

$$\nabla f(x,y) = \left[ \frac{\partial f(x,y)}{\partial x}, \frac{\partial f(x,y)}{\partial y} \right] : f(x,y) = 3x^2y$$

$$= [6xy, 3x^2]$$

Jacobian is Matrix of Derivative

Chain Rule; In chain rule, we multiply derivatives

- 1)  $l_1(x, w_1)$  " input layer
- 2)  $r(l_1)$  " relu at layer 1
- 3)  $l_2(r, w_2)$  " layer 2
- 4)  $\text{loss}(l_2, y)$  " loss computation

So to compute derivative of 'loss( $l_2, y$ )' but

(4) is  $\sim f_n$  so, then we'd need to multiply  
its derivative with derivative of (3) but  
that's also  $\sim f_n$ , so we'd need to multiply  
(3)'s derivative with (2) and so on.

This is back propagation

$$\frac{\partial \text{loss}^{\text{out}}}{\partial l_2} \times \frac{\partial l_2}{\partial l_1} \times \frac{\partial l_1}{\partial w_1}$$

$l_2$  [lin-grad ( $l_2, \text{out}, w_2, b_2$ ) # out. g are gradients  
of loss]

$$\frac{\sum(x^*)^2}{n}$$

$$\boxed{2 \times \frac{d}{dx} x}$$

helps forward pass values  
 helps gradient  
 backpropagation  
 backward pass.

$\Rightarrow$  lin-grad ('inp', 'out', 'wb'),

'inp' g: gradient of the input w.r.t the loss.

This line calculates the gradient of loss w.r.t  
input 'inp'. Here, the gradient is computed  
 by taking the dot product of the gradient  
 of the output ('out', 'g') and transpose of  
 weight matrix. This process/operation  
 essentially backpropagates the error from the  
 output to input through the weight  
 matrix.

$$\frac{\partial L}{\partial \vec{v}} = \frac{\partial L}{\partial \text{out}} \cdot W^T$$

→ 'w.g' is  $\text{lin\_grad}(\text{inp}, \text{out}, w, b)$

↳ Gradient of the weight matrix w.r.t loss  
It calculates how the loss changes w.r.t 'w'.

w.r.t weight

$$\frac{\partial M}{\partial w_2} = w_2 \frac{\partial (w_2 x_2)}{\partial w_2}$$
$$= (x_2)$$

↓<sub>o</sub>

Multiply inputs, weights & bias  
with output gradient

P - output      | w - where am I      | r - return  
c - continue.      | n - next line

[ Derivative of sum of fn  $\propto$  sum of derivatives of fn. ]

F: O

37)

backward -  $\text{q\_forward}(\text{inp}, \text{tag})$

→ # forward pass

$L = \lim ($

forward-pass =  $\text{call}()$

↳ to process forward path

↳ backward()

endum

$i^o, j^o, k^o \rightarrow i^k$

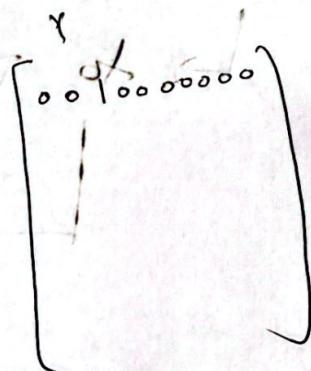
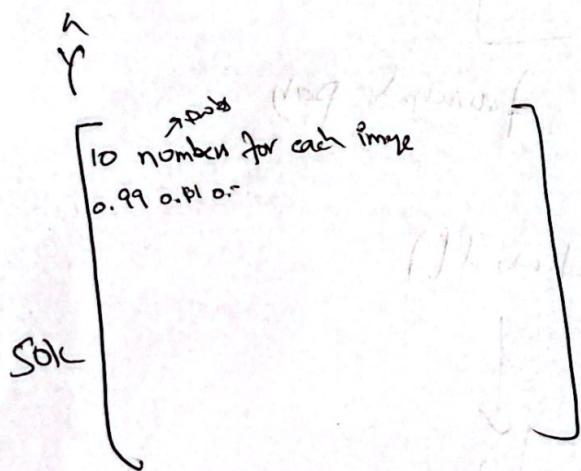
$i^o$   
inp  
set  
→ predicted  
error

predicted  
error

Cross entropy loss

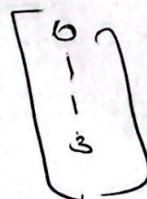
$$\text{softmax} \left[ \begin{bmatrix} x \end{bmatrix} @ \begin{bmatrix} w \end{bmatrix} + b \right] = \hat{y} = \begin{pmatrix} 0.99 & 0.01 & 0 \end{pmatrix}^T$$

here we are



Targets would be one-hot encoded

or we can just store actual target variable



38)

	Model's output	$\exp$	Softmax	Actual	
Cat	-4.89	0.01	0.00	0	6
Dog	2.60	13.43	0.97	1	-0.65891
Plane	0.59	1.81	0.12	0	0
Fish	-2.07	0.13	0.01	0	0
Building	-4.57	0.01	0.00	0	6
		<u>15.38</u>	<u>1</u>		

Softmax output  
will be compared

$$-\sum_{j=1}^M y_i \log \phi(\hat{y}_i))$$

we can just look at actual label and only look at that value i.e value of 1 by dog -0.65 matters only, else is discarded

So, For CE loss;

so,  $\hat{p}$  (after softmax word)

$$[50,000, 10]$$

so for each of 50000 examples  
we have an actual value (label), so  
we only need to consider that  
index value from <sup>each</sup> prediction row.

Softmax [actual, index]