

## Kernel Networks

⇒ Advantages:- ① even if reduced dim. coords  
are complex func<sup>cs</sup> of our original  
coords, they are still explicit func's.

$$Y = f(X)$$

$\downarrow$   
 $J$

→ ' $f$ ' is invertible, so decomposition can be computed

→ ' $f$ ' is differentiable, in principle,  
so - no. of times depending on  
our choice of NNs.

(in principle) with some caveats

- (2) We introduced Kernel methods to treat  
highly nonlinear & complex data manifolds.  
If we introduce a kernel, we implicitly  
decide to treat the neighbourhood of our data  
in a given manner (determined by the  
kernel),

For ex, if we say the kernel is Gaussian  
in distance

$$K(x^i, x^j) = \exp\left(\frac{-\|x^i - x^j\|^2}{2\sigma^2}\right)$$

→ only varitional parameters we have is  $\sigma$

→ so, we have few parameters to play with & in kernel methods we sort of assume we know the structure/relationship b/w  $x^i$  &  $x^j$  over data / distances.

⇒ In NNs, there are many variational parameters.

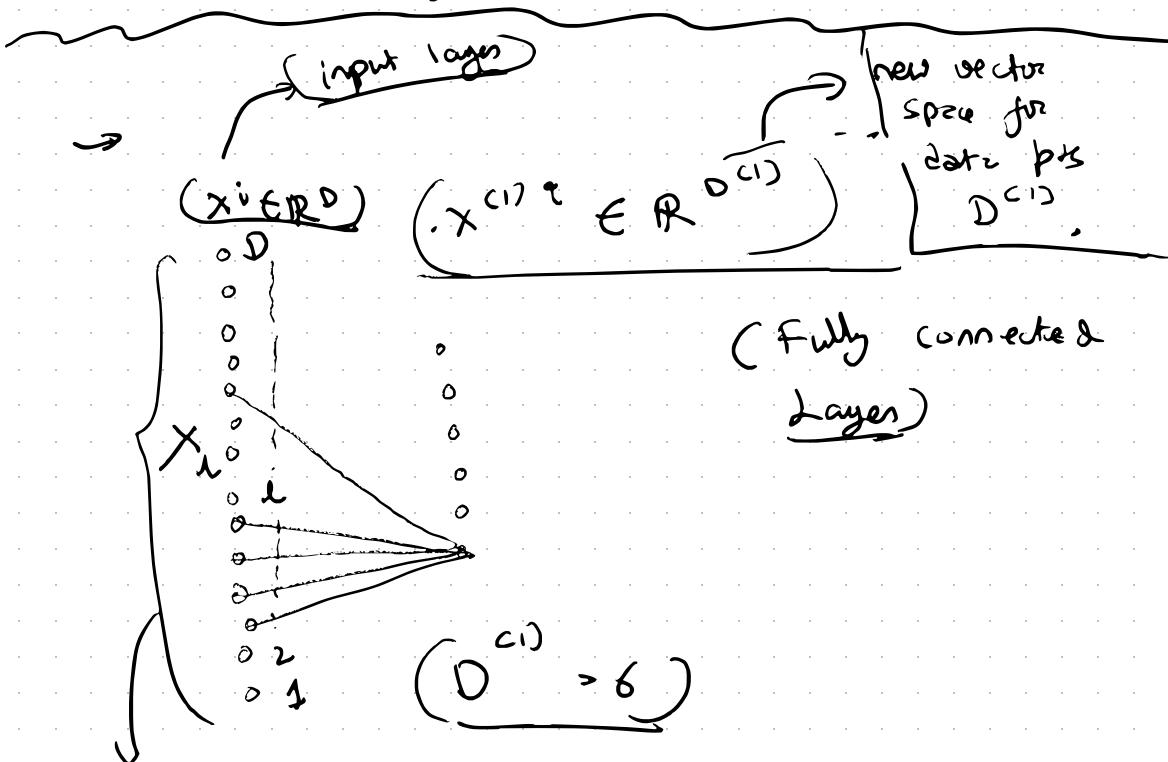
$$Y = f(X | \pi)$$

→  $Y$  depends explicitly on a very large # of params ( $\pi$ ). Ordinary NNs are tuned to reproduce ground truth predictions/represent our data manifold.

In many applications it is really hard to choose / find a suitable kernel. This is where NNs are powerful.

## → Disadvantages of NNs:-

- ① really complicated process + training a NN is like alchemy + is poorly understood.



→ No. of input nodes in our network = D

→ The input node will take as input  $(x_d)$  i.e., the  $l^{th}$  feature of our data point.

$\Rightarrow$  We are mapping our data from

$x^{(0)}_i \in \mathbb{R}^{D^{(0)}}$  in input layer to

$x^{(1)}_i \in \mathbb{R}^{D^{(1)}}$  which is the next layer  
of our autoencoder

$D^{(1)}$  is the new vector  
space of our data

$\rightarrow x^{(0)} \rightarrow$  input data

$\rightarrow D^{(0)} = D \rightarrow$  initial dimensional space  
of our data.

$\Rightarrow$  So this transformation from  $x^{(0)} \rightarrow x^{(1)}$   
is in general performed with an explicit  
func<sup>n</sup> & there are many possibilities.

$\Rightarrow$  architecture in NN : func<sup>n</sup> def<sup>2</sup>  
that maps data across the layers

- The example given so far has
- Fully Connected layers architecture

→ In a "Linear" Fully Connected layers :-

- (1) we do something identical to PCA

$$(1) \quad x^{(1)} = W \cdot x^{(0)}$$

(matrix)

↓

$(= A \text{ in PCA})$

→ same operation as PCA  
 transformation done independently for  
 each data point (i)

→  $x^{(1)}$  has dim =  $D^{(1)}$

→  $x^{(0)}$  " " =  $D^{(0)}$

→ ∴ ( $W$ ) will be a rectangular matrix  
 with dimension =  $\underbrace{(D^{(1)} \times D^{(0)})}$

→ If our NN is a Linear func<sup>2</sup>  
is a fully connected layer :-

$D^{(1)} = 5$  (as in PCA)  
then the functional form of this single  
layer NN will be identical to PCA.

→ if we then optimize the NN to reproduce  
the distances we EXACTLY get PCA.

② if we have a fully connected nonlinear  
layer :-

→ we apply "pointwise nonlinearity"

→ i.e., to each element of  
 $x^{(1)i}$  we apply independently  
a nonlinear function of our  
choice,

mathematically

$$\underline{x_n^{(1)}}$$

$$= g \left( \sum_{m=1}^{D^{(0)}} w_{lm} x_m^{(0)} \right)$$

$$\left\{ \begin{array}{l} \text{component } (l) \\ x^{(1)} \end{array} \right\}$$

nonlinear  
func<sup>=</sup>

representing

$w x^{(0)}$  in terms of  
components.

$$(w x^{(0)}) = \sum_{m=1}^{D^{(0)}} w_{lm} x_m^{(0)}$$

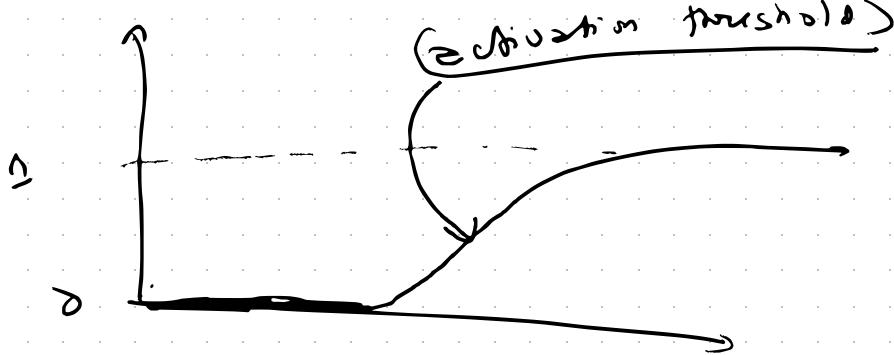
This nonlinear func<sup>=</sup> ( $g$ )

is applied to each component  
of the generated vector,

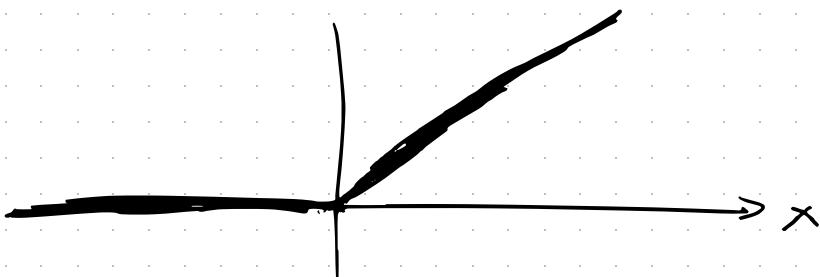
In principle, ( $g$ ) can be diff. for  
each  $(l)$ , but in practice, ( $g$ ) is same for  
each layer.

→ Functions for us of nonlinear func<sup>e</sup>(g)

①  $g(x) = \tanh(x)$   
(switches from 0 → 1)



② RELU (Rectified Linear Unit)



$$\text{RELU}(x) = \begin{cases} x, & x > 0 \\ 0, & \text{else} \end{cases}$$

→ In both cases we introduce an entire variational parameter, called "shift"

→ So we typically optimize a function of the form:

$$g(\lambda - \lambda_e)$$



"shift" parameter

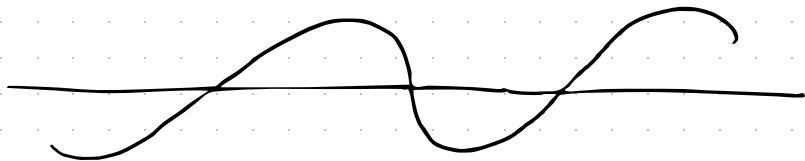


shift  $\lambda_e$   $x^{(0)}$  by a quantity

which we try to optimize

NOTE  $\lambda_e \rightarrow$  so can vary from component to component.

$$\textcircled{3} \quad g(x) = \cos(x)$$



\textcircled{4} ... there can be many different forms of  $g(x) \rightarrow$  any nonlinear form is mostly used if it suits our data.

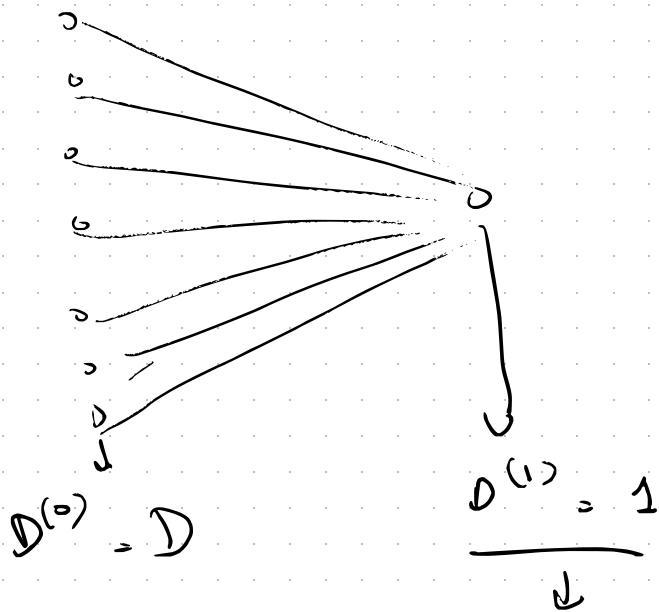
$\Rightarrow$  Why does RELU/tanh get used so much?

\textcircled{5} ; NNs try to mimic how neurons in the brain work, neurons fire only if total activity is above = certain threshold (activation threshold) & RELU/tanh mimic this property very well.

$\hookrightarrow$  "Weights"

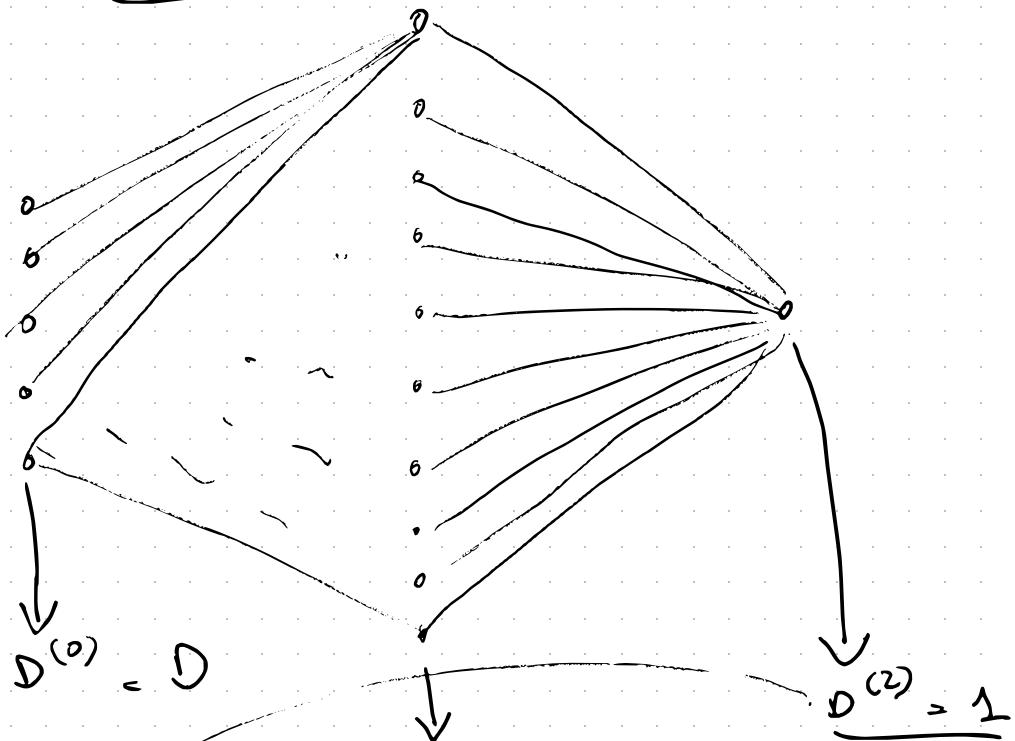
→ applying these transformations successively  
in multiple layers made NNs very popular

⇒ Single "Perception" model from  
the '20s



We want a single  
real value result.

→ "Perception" model in practice:-



Intermediate layers

$$D^{(1)} = \infty \quad (\infty \text{ # of nodes})$$

We want an output that is only a single real number

"We can represent any nonlinear func<sup>2</sup> of the input with RELU/tanh etc. activation func<sup>3</sup>"

(Theorem" relevant to Perceptions)

$$\rightarrow Y = \underbrace{f}_{\downarrow}(x)$$

complicated func<sup>2</sup> of the input ( $x$ )

IFF, it has a single value,  
the perceptron architecture allows  
us to represent this func<sup>2</sup> with  
a given activation func<sup>2</sup> of our choice.

Case 1  $\rightarrow$  even if this representation  
is known & it exists (as the theorem  
states) algorithmically it is very difficult  
almost impossible to find it.



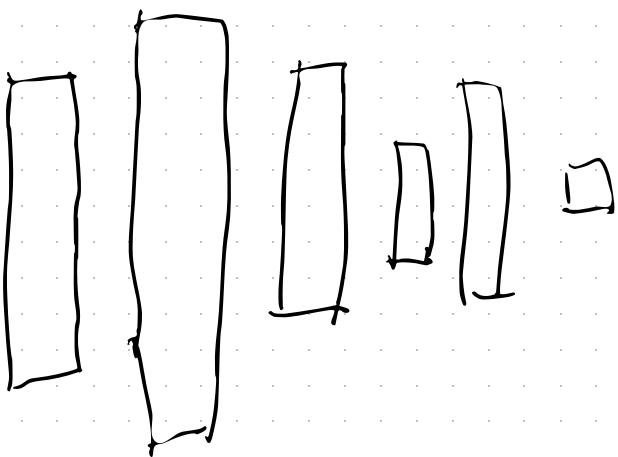
③ Breakthrough in MLPs :-

Multi-layer Perceptrons

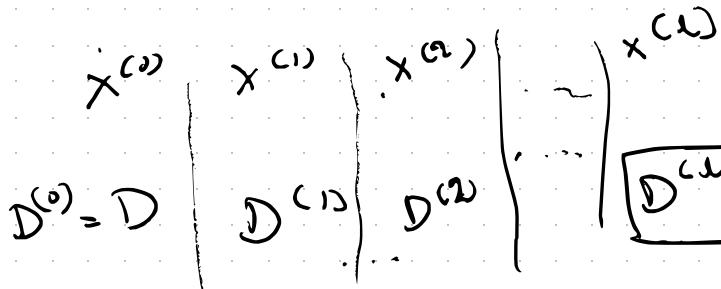
Perceptrons)

( $\sim 2010s$  or so)

Clinton  
de Cun et al.)



They realized that concatenating many layers of perceptrons finding one soft becomes easy.



$$D^{(L)} = 1$$

simple case, where output is a single no.  $\in \mathbb{R}$

$\rightarrow L = n_o$ . If layers in our deep NN.

$\rightarrow$  each representation is a func<sup>?</sup> of preceding representation.

$$\Rightarrow \text{e.g. } X^{(L+1)} = f_L(X^{(L)} | W^{(L)})$$

Set of weights / Parameters specific to a layer.

→ In principle,  $f(\cdot)$  can be diff. for each layer, & usually is

→ But for simplicity, let's look at the case where  $f(\cdot)$  is same for all layers.

→ The nature of this approach is that it implements a Markovian Process.

$x^{(L+1)}$   
depends ONLY on  $x^{(L)}$

→ There can be skipped links / even bidirectional links but the reason it is useful to implement our NN as a Markov Process is that there is a simple way to optimize it

$\begin{matrix} l \rightarrow l+2 \\ l+2 \rightarrow l+4 \text{ etc.} \end{matrix}$

↳ Backpropagation

In USL,  $D^{(L)} \rightarrow D$  is taken, no dim reduction is performed in their output

## Training a NN :-

→ For the case of binary classification —

(i)  $X^i$

(images of cat/dog)

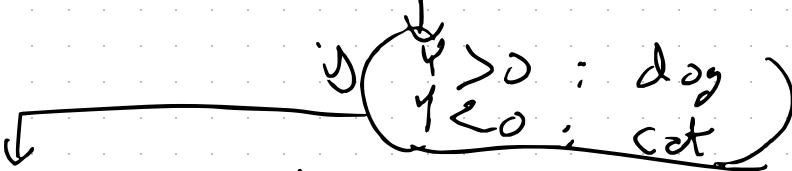
↓      input data

$$Y = f(X | W)$$

parameters

$\rightarrow (x)$

$f(\cdot)$ , given an image will return  
= red no.,  $Y$  which is the "dogishness"  
of an image.



We want our NN to be using only differentiable func —.

→ For each image :- we have a ground truth label → binary label →  $b^i = 1$ , if  $i = \text{dog}$ ,  
 $= -1$  if  $i = \text{cat}$

→ What we want is to find:

$$\tanh(y^i) = \tanh(f(x^i | w)) \\ \approx b^i$$

by  
training  
the NN

We find  $w$ , s.t. this is  
true for all  $(i)$

$\tanh()$  is chosen as  $f(\cdot)$

⇒ How do we find these desired  
parameters?

→ By minimizing p.r.t. the parameters,  
a "loss func" → which is a func<sup>2</sup> of  
our parameters &  
data.

$$L(\mathbf{w} | \mathbf{x}) = \frac{1}{N} \sum_{i=1}^N (b^i - \tanh(f(x^i | \mathbf{w})))^2$$

minimize this  
loss func<sup>2</sup>  
w.r.t.  $\mathbf{w}$

L2 - loss func<sup>2</sup>

(loss func<sup>0</sup>)

measures  
# errors  
we make.  
if all predictions  
are correct for binary classification,  $\lambda = 0$ ,  
if  $\lambda = 0.01$ , 1% errors are made.

other possible loss  
func<sup>2</sup> can be the  
Cross-entropy loss func<sup>-1</sup>.

$y = f(x^i | \mathbf{w})$   $\Rightarrow$  chain of func<sup>2</sup> ; we  
 $= x^{(i)}$  describe the NN as a  
Markov Process

$$\begin{aligned} \Rightarrow \therefore y &= x^{(i)} = f(x^i | \mathbf{w}) = f(x^{(i-1)} | \mathbf{w}^{(i-1)}) \\ &= f(g(x^{(i-2)} | \mathbf{w}^{(i-2)})) | \mathbf{w}^{(i-1)} \\ &= \dots \quad (\text{so on till we get the input}) \end{aligned}$$

∴ Basically  $\hat{x}^L$  is written as a  
 $\text{func}^2 f = \text{func}^2 f = \text{func}^2$

⇒ So basically, we take a derivative of  
the loss  $\text{func}^2$  w.r.t. the parameters.

→ If the NN has a Markovian str., this  
is easy & such a NN is called

"feed-forward NN"



$\frac{\partial L}{\partial w} \equiv$  compute the derivatives  
explicitly.

→ train the NN by (plain  
gradient descent.)

Start :- ①  $w \sim$  random gaussian no. <sup>(initial guess)</sup>

Then :- ②  $w = w - \epsilon \frac{\partial L}{\partial w}$

$$\rightarrow W = V - \epsilon \frac{\partial L}{\partial V}$$

plain gradient descent

$\epsilon \rightarrow$  Learning Rate

$\Rightarrow$  This process is called "Backpropagation"

diff. manners of  
computing derivative  
of  $e$  func<sup>2</sup> of  $e$   
func<sup>2</sup> of  $e$  func<sup>2</sup> of  $e$   
by chain rule

→ "Deep" NN → it's called deep because many layers are concatenated one after another



① Consequence of this :- In a typical SL/ML setting, we try to learn a model that depends on data ( $x$ ) & set of parameters ( $w$ )

$$\therefore Y = f(x|w) \quad \rightarrow \text{typical setting of ML training protocol}$$

→  $N$  data points on which we train our model.

Say,  $N = 10,000$

→ Another thing that defines our model is # of free params. If we do lin. reg., (# of params = dim. of data)

∴ → "m" ≠ # of parameters

→ " $\eta$ " → dimension of data

① What if for  $N = 10^4$ , we want

$$M = 10^6?$$

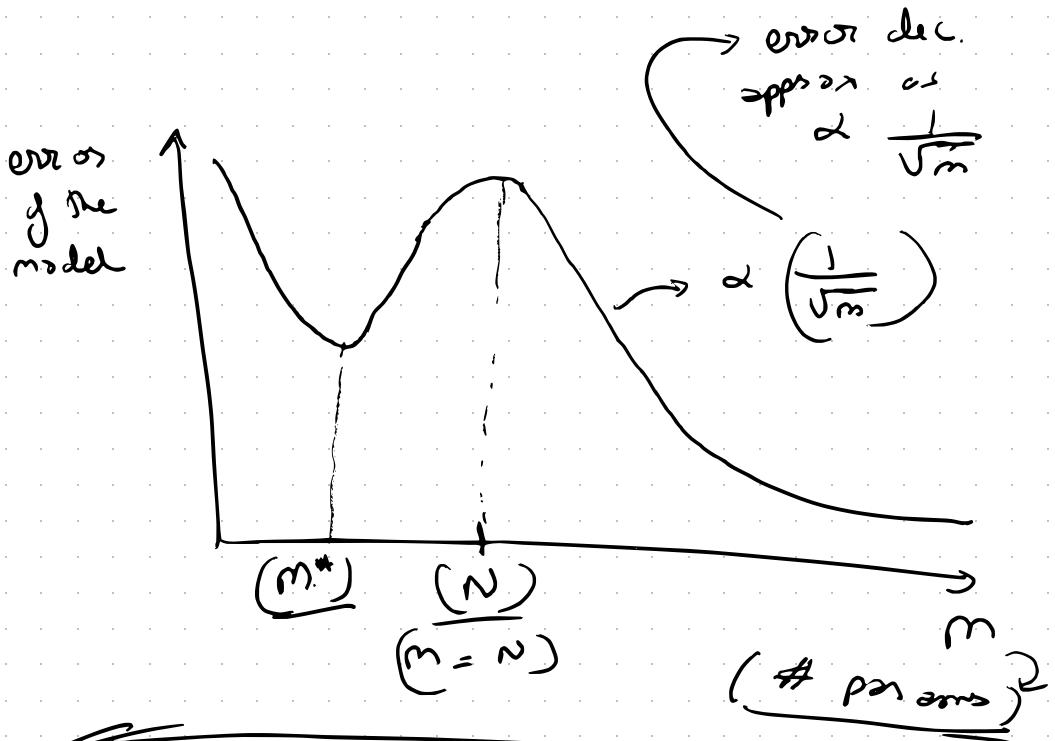
→ normally, we would expect

( $m$ ) to be 1 order of magnitude  
less than ( $n$ ) → data

→ instead in  $NN_s$  we go at odds  
with this paradigm & choose  $M = 10^6$

for  $N = 10^4$ ,

→ i.e., we choose an insanely large  
number of parameters to describe our  
data, in fact (# params) much  
greater than (# data pts).



### Double - Descart curve

- few params → high errors
- increase ↑ → errors reduces
- even more increase → overfit the data
  - ↓
  - prediction in ← reproduce training data
  - garbage
  - error increases.

(i)  $m^*$  = best no. of parameters in classical statistical learning.

- Beyond  $(m^*)$ , error  $\uparrow$  due to overfitting.
- error peaks as  $(m \sim N)$   
 $\# \text{params} \sim \# \text{data}$
- then error miraculously  $\downarrow$ .

(double descent phenomenon)

observe first in NNs, but also  
in other models

## (i) Law's interpretation of Double Descent

- classical statistical regime
- $N$  data,  $m$  params
- if training instead of being deterministic is stochastic

↳ stochasticity can be trivial. For eg  
loss landscape has many local minima  
we initialize params in a statistically  
independent manner & then do gradient descent

→ we do optimization since to obtain  
the model :-

$$f(x|w)$$

→ if loss landscape is complex, we can initialize 1000 times & so instead of having 1 model, we have 1000 models.  
Finally, our best guess for ( $\hat{y}$ ) is chosen as a statistical avg. of the output of 1000 models.

$$\therefore \hat{y} = \underset{m \in \text{models}}{\sum} f(x|w^m)$$

So our optimization procedure having a degree of stochasticity improves the quality of our prediction.

$\rightarrow \because m$  is initialized randomly

$$y_i = \sum_{\text{models}} f(x_i | w^m)$$

has  $(1000 \times m)$  # of persons.

$\rightarrow$  This can also be looked at as a single model with  $1000 \times$  more persons.



Simple case if we can improve reliability of our prediction by inc. the

~~# of params.~~

$\rightarrow$  error dec  $\propto \frac{1}{\sqrt{m}}$  ) it

implies that a deep NN is implicitly

defining a large # of statistically indep.

models.  $\therefore$  A deep NN can be considered to be averaging over statistically indep. representation of our data.

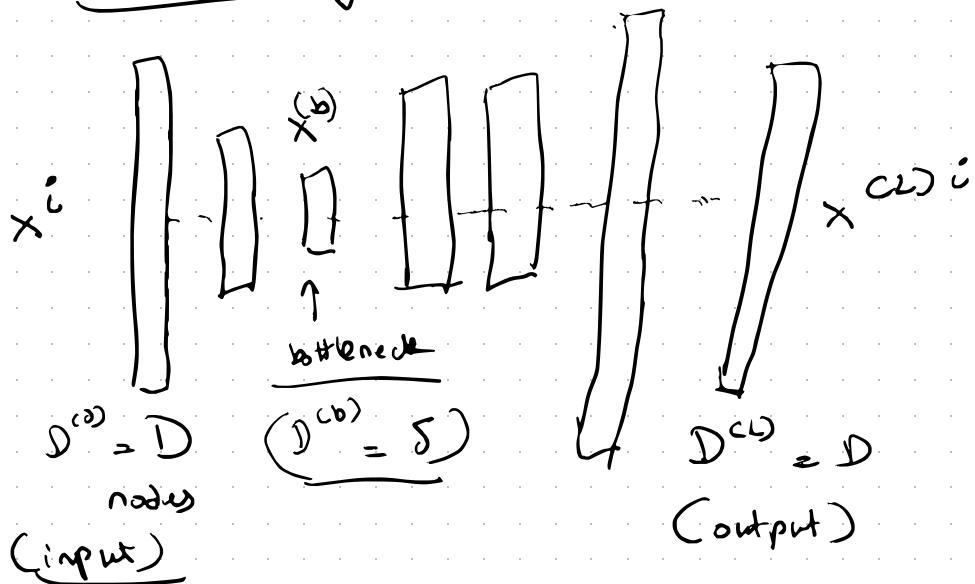


## Auto encoders (AE)

- all real world auto encoders work in the over-parameterized regime.
- having an insanely large # params allows us to find a sol' :: having so many params allows us many diff. statistically indep. methods of fitting our data.
- That's why NNs & AEs work well in the over-parameterized regime.

"Distillation" → train in over-parameterized regime then kill unused params.

# ① Structure ↴ AE ↵



- input layer takes as input our original data
- output layer also ALWAYS has (D nodes)
- $L \geq \# \text{ of layers}$
- $x \neq x^{(L)}$  belong to same space,
- in this  $\#$  we have a chain of others layers.

→ no reason for layers to be symmetric.

→ In AEs, there's a layer called a bottleneck, which is the smallest intermediate layer with the smallest dimension (least # of nodes).

→ Dimension of bottleneck = 5 (like in PCA)

$$\rightarrow \boxed{Y = X^{(b)}}$$

low dim. representation of our data that we try to learn  $\equiv Y = \underline{X^{(b)}}$

→ Basically, our AE is typically a FF NN (feed forward NN) that performs a Markov Process on our data, tries to reduce the dim. of our data representation upto a minimum value chosen by the user.

→ Train this FFNN / AE :-

① define a loss func :-

$$\begin{aligned} L(\omega) &= \sum_i \left( \frac{x^{(l)}_i}{\text{↓}} - x^i \right)^2 \\ &= \sum_i \left( \underbrace{\int (x^i(\omega) - x^i)}_{\text{↓}} \right)^2 \end{aligned}$$

Concatenation of several  
complex funcs.

↓  
Output is an explicit func<sup>2</sup> of the  
input, that is produced by concatenating  
many layers.

→ each layer's transformation = some

①  $\int x^{(l-1)i} \underbrace{w^{(l-1)}}_{\text{Layer specific pars.}}$

∴  $L$  depends on all the params  $(\omega)$

→ finally  $\mathcal{L}(w)$  depends on ALL these persons.

→ implication of loss func<sup>2</sup>: We want the A, E, to learn the identity of our original data points.

$$\rightarrow x^{(w)i} = \underbrace{\gamma(x^i | w)}_{\downarrow}$$

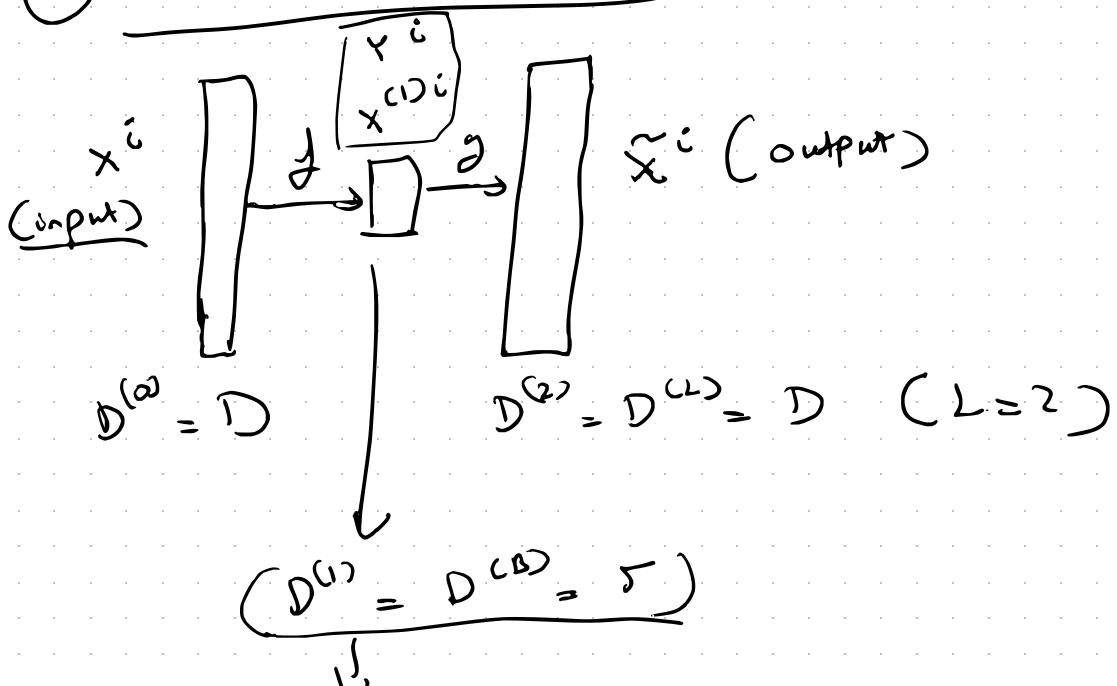
$x^{(w)}$  is a func<sup>2</sup> of original data, given the persons ( $w$ )

Goal of the A, E<sub>o</sub>: Find  $w$

that minimizes  
 $\mathcal{L}(w)$

□ Simple eg. of A.E. :-

○ Linear Autoencoder :-



→ only 1 hidden layer

→ 2 func  $\rightarrow f \leftarrow g$

both are linear func  $\Downarrow$

$$\rightarrow \gamma^i \text{ (bottleneck representation)} = x^{(1)i}$$

↓

$$\tilde{\gamma}_w^i = w_l \cdot x^i ; l=1, \dots, \delta$$

( $\ell^{\text{th}}$  component of vector  $\gamma^i$ )

Then we can variationally define

$$\underline{w_l}$$

↓

Vector of dim = D

⇒ Determine free ( $\delta$ ) vectors to minimize loss func<sup>2</sup>,

→ from our representation of data in the bottleneck we have to reconstruct original data. This is linear projection → first project on  $W$ , & then we reconstruct original  $x$ . We have to represent  $x$  on the basis of  $W$ .

$$\therefore \hat{x}^i = \sum_{l=1}^5 Y_l^i w_l$$

$$= \sum_{l=1}^5 (w_l \cdot x^i) w_l$$

$$[\because Y_l^i = w_l \cdot x^i]$$

→ Rewriting the generalized loss func<sup>2</sup> for this specific case :-

$$L = \sum_i (x^i - \hat{x}^i)^2$$

$$= \sum_i \|x^i\|^2 + \sum_i (\|\hat{x}^i\|)^2 - 2 \sum_i x^i \cdot \hat{x}^i$$

→ Let's choose ' $w$ ' that preserves the

"norm" of the data, i.e., we choose

$$w, \text{ s.t. } \|x^i\|^2 = \|\hat{x}^i\|^2$$

same con<sup>2</sup> as ← { norm of  $x^i$  = norm of  $\hat{x}^i$   
 for deriving MDS  
 can be done with  
 Lagrange multipliers }

→ To minimize  $\lambda$ , we need to

maximize the scalar part.  $\sum_i x^i \tilde{x}^i$

$\therefore \|\tilde{x}^i\| = \|x^i\|$  is already fixed

+ constant.

→

→ Let's rewrite the scalar part:

$$\sum_i x^i \tilde{x}^i$$

$$= \sum_i x^i \cdot \sum_{l=1}^d (w_l \cdot x^i) w_l$$

$$= \sum_{l=1}^d \sum_i (w_l \cdot x^i)^2 \quad (\underline{x^i \rightarrow \text{scalar}})$$

(rewrite in components)

$$= \sum_l \sum_{m m'} w_{lm} x_m^i w_{l m'} x_{m'}^i$$

↓  
summing over  $i$ , we get  
the covariance matrix (ref: PCA)

$$= \sum_m \sum_{m'} C_{mm'} W_{lm} W_{lm'}$$

$$\therefore L = \text{Tr} (W C W^T)$$

$$\begin{aligned} & \xrightarrow{\text{trace}} \sum_l \\ & \xrightarrow{\text{multiply by } l} \sum_m \\ & \xrightarrow{\sim \sim W^T} \sum_{m'} \end{aligned}$$

this is basically  
the same as loss  
func<sup>2</sup> in PCA



the only difference is the normalization cond<sup>2</sup>  
 which in PCA is implemented by adding the  
Lagrange Multipliers.

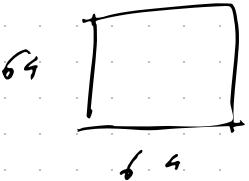
⑤ TL; DR :- A linear autoencoder in which  
 the dimension of the bottleneck is 5 &  
 the encoder & decoder weights (fwd &  
 bwd) are the same, is equivalent to PCA.

- $\rightarrow \therefore$  PCA can be seen as = very special  
 f simple case of an A. E.  
 $\rightarrow$  In PCA we project our data on a  
low dim. Space, & we assume this is done  
 without significant information loss, & so  
 we assume we can reverse our transformation  
 & go back to the original data.
- $\rightarrow$  if we are NOT able to project back,  
 then PCA is NOT opt.
- $\rightarrow$  if rank of C is  $\leq n$ , then  
 $\delta = \text{rank}(CC^T)$ ,  $\lambda = \sum_{i=1}^n C_i^T x_i$   
 $A \tilde{x}_i$  are  
 collinear &  
 when subtracting  
 terms  $\lambda x_i^m$ , we  
 get 0)  
 here we ignore  
 the non zero terms  
 in final represent.  
 of loss func<sup>2</sup> for  
 convenience.

$\rightarrow \text{if } \text{rank}(C) > \sum m - d \neq 1$

& rows' heights loss,  $\therefore$  some info is lost.

$\rightarrow$  this shows why in PCA we try to maximize the trace of the covariance matrix, which is the same as minimizing the loss func<sup>2</sup>.

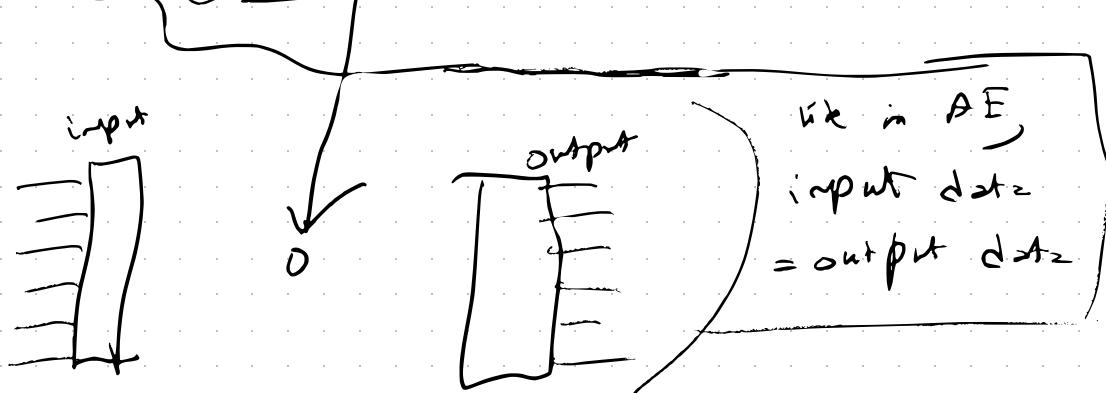
- paradoxical manner of defining an A.E.
- dataset of images →  $\frac{64 \times 64 \text{ pixels}}{\text{black \& white}}$
- 
- $\downarrow$
- $16 \text{ levels of gray}$
- $(2^4)$
- pics of 100 diff animals doing diff. stuff. We D.

→ architecture :-

(1) input → pixel of image

(2) output → " " "

(3) bottleneck → single node



→ Is this possible?

- ① It seems No ∵ we think (so far) that dim. of bottleneck should be  $\approx 1D$  of dataset.
- ② However, this is not only possible it's also quite easy.

↓  
↓  
(HDD?)

- we made some assumptions to simplify some things.
- It levels of gray  $\rightarrow 2^4$  possible values of each pixel
- $64 \times 64$  pixels  $\rightarrow (2^6 \times 2^6)$   
∴ any possible such image with this level of accuracy + this gray level can be uniquely encoded by assigning  $(2^4, 2^6, 2^6 = 2^{16})$  bits.

→ So, if this is represented by a single double precision no, (contains  $2^{16}$  bits)

the info we store in each bit of the double precision # allows us to reconstruct the image.

→ However this doesn't mean we have learnt any info of our dataset, rather we can "compress" the image, which is a different task.

→ i.e. To obtain an A → E. that can learn meaningful features we need some stochasticity in the training process.

Why? In the architecture we defined now, an image of donkey & cat can be mapped quite close to each other in ( $Y$ ) space.

- ;" To train on A.E. on just information compression, we violate the rule of USL that data which is semantically close should be close in reduced dim. representation too,
- ;, Dim value is totally useless ;  
we understand nothing

Now to fix this problem?

→ Instead of asking the A.E. to reproduce exact input image, we ask the A.E. to reproduce an ensemble of images that are very similar to input image with some tiny added noise.

e.g., if input image = Cat,  
output = ensemble of cats with  
slightly longer/shorter whiskers etc.

→ General idea :- train A.E. to

reproduce not exact input imagine,  
but something which is almost the same  
even if we perturb it.

→ So, our previous idea of using a single  
 bottleneck node can't work any more.

"Denoising Auto encoders"

$$L = \left\langle \sum_i \left( f(x^i + \underbrace{\text{noise}}_{\downarrow} | w) - x^i \right)^2 \right\rangle$$

(can be GAN etc.)

→ empirical avg. over all realizations  
of the noise.

$\Rightarrow$  This forces the bottleneck size to be AT LEAST  $> ID$  of the dataset ; here we force the A,E. to learn the neighbourhood of each pixel,  $\therefore$  similar images can be categorized effectively by the A,E. now

### Problems of Denoising A,E. :

- choosing appropriate noise is hard.
- choosing indep. SDN for each pixel is quite arbitrary.
- $\because$  we would prefer noise that is semantically connected, & the image after adding noise is transformed s. that it makes sense to us physically, & is not just something arbitrary.

## Variational AutoEncoders (VAE) :-

→ helps solve this problem with "denoising"

A.E. <sup>'</sup>

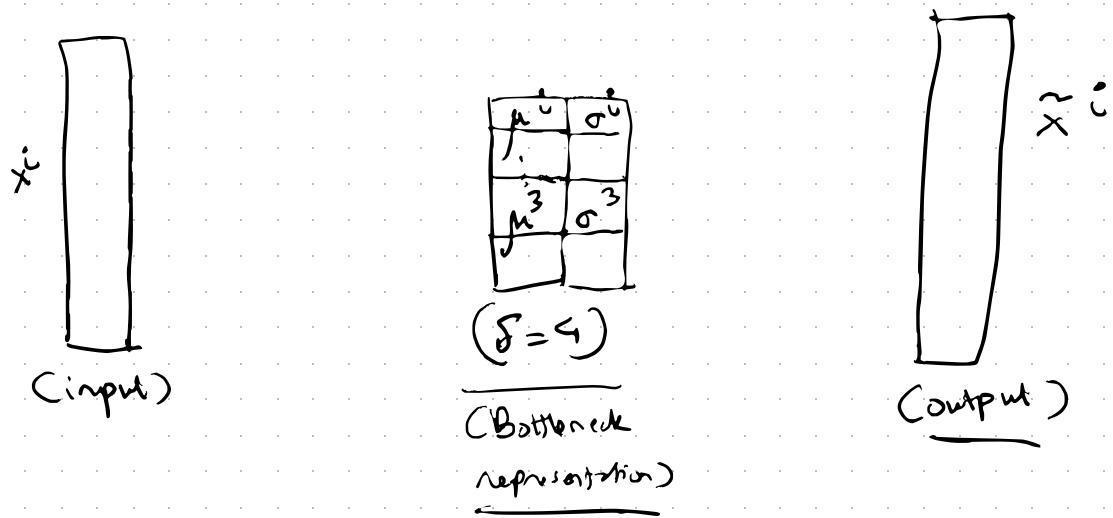
→ Key idea :- Adding the noise in the bottleneck.

→ in denoising A.E., we add ( $\eta$ ) to  $x^i$  i.e. we add noise to the input, then we have something arbitrary.

→ We assume for a meaningful A-E, the intermediate representation in the bottleneck will have some semantic meaning. So, images of cats should be similar somehow in S-space. So we add noise in the bottleneck step. After this it reproduces something close to original image we can say that the VAE learnt the relevant semantic features.

→ Now to do this in practice ?

## ① Architecture of VAE :-



→ assume: what we learn is not a single representation of our input  $x^{(b)i} = y_i$  (bottleneck representation) but rather it is a probability density in  $\{y^i\}$  space,

→ Say,  $\delta = 1$  ( $D^{(B)} = 1$ )

→ We learn  $\hat{y}^i = \mu^i + \sigma^i$  ( $i = 1, \dots, 5$ )  
→  $\mu^i + \sigma^i$  plays role of  $\langle y^i \rangle + \text{std dev. of } y^i$  ( $\delta$ )

$\rightarrow \therefore$  We can sort of see component ( $i$ )  
 $y^i \rightarrow y_u^i$ , as harvested from  
 $\Rightarrow$  normal distribution of average ( $\mu_u^i$ )  
+ variance ( $\sigma_u^i$ ).

$$\therefore \boxed{y_u^i \sim N(\mu_u^i, \sigma_u^i)}$$

so, instead of learning directly the representation of  $y^i$ , we learn the parameters of  $\Rightarrow$  Gaussian distrib<sup>2</sup> from which we can then harvest ( $\mu^i$ ) and generate a representative representation.

$\Rightarrow$  Let's try to implement this architecture in practice. For simplicity, we take  $\sigma = 1$  & so we have 2 nodes in the bottleneck layer (1 for  $\mu$ , 1 for  $\sigma$ ).

Bottom node

Encoder

layer

Decoder

$x_i^c$

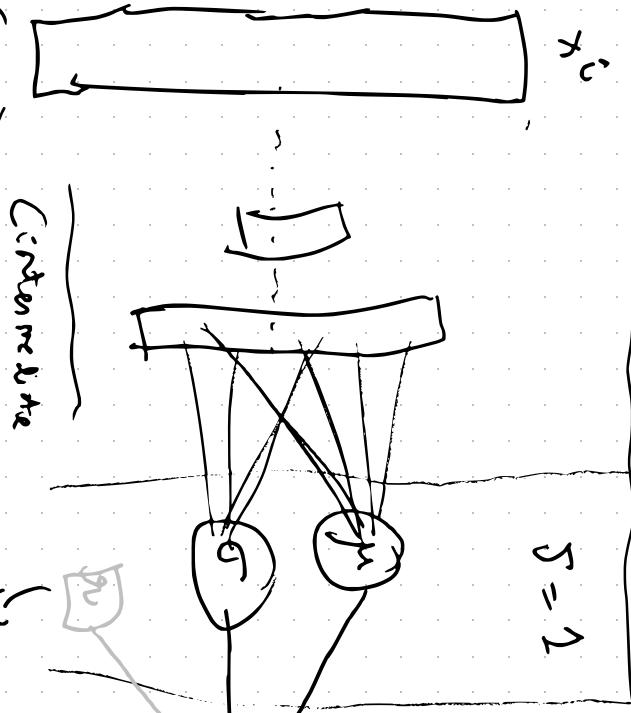
$J = 1$

$x_i^c$

Comput

single  
random  
var.

$$\mu \sim N(0, 1)$$



Input

Intermediate  
Layers

drain independently over  
several cycles. Can be diff. for each input.

$\rightarrow (x)$  is harvested from a Gaussian  
 / Normal distrib<sup>2</sup> with mean =  $\mu$ ,  $\text{Var} = 1$   
 s. that we transform  $y^i$  to a Normal  
 Distrib<sup>2</sup> with mean =  $\mu^i$  &  $\text{Var} = \sigma^2$ .

$$\Rightarrow L(\omega) = \sum_{i=1}^n \left( f(x_i, r | \omega) - x_i^i \right)^2$$

Added in the bottleneck layer

So, minimizing this loss func<sup>2</sup>, we learn  
 not only a 1000 dim. representation for each  
 data pt. (or image), we also learn how  
 variable the representation is

$\rightarrow$  In denoising A.E., noise is chosen by  
 us. So, for large ( $\lambda$ ), we get crap  
 results; we blur all the image.

→ In VAE, the level of noise is learnt by the network itself, ∵ the network must produce output (after decoding) at least similar to input.

→ We must add some constraints on the optimizer, ∵ if  $\sigma^i = 0$  there's no noise term, & we have same problem as before.

→ How to fix this?

① We add to loss a prior on the values of  $\mu + \sigma$ .

expected prob. distrib<sup>2</sup> we'd like to see on the data.

→ We have many different choices.

$\rightarrow$  Simplest possible choice of Prior :-

$$\begin{aligned}\mu^i &= 0 \quad \forall i \\ \sigma^i &= 1 \quad \forall i\end{aligned}$$

arbitrary  
choice for all  
prior

$\hookrightarrow$  DE can use additional info.  
to drive our choice of prior.

$\rightarrow$  Key criteris is to ensure that  $\sigma^i \neq 0$   
regardless of other choice of prior.

$\rightarrow$  Another choice of prior :-

$$\rightarrow \mu^i \rightarrow \text{uniform } \forall i$$

$$\rightarrow \sigma^i \rightarrow \text{forced to be } \neq 0 \quad \forall i$$

∴ loss func<sup>=</sup> is re-written as :-

$$\mathcal{L}(\omega) + \frac{1}{2} \sum_{i=1}^{N_l} \sum_{l=1}^{\infty} \left( (\mu_i^i)^2 + (\sigma_i^i)^2 - \log \sigma_i^i \right)$$

↓  
(as before)

↑

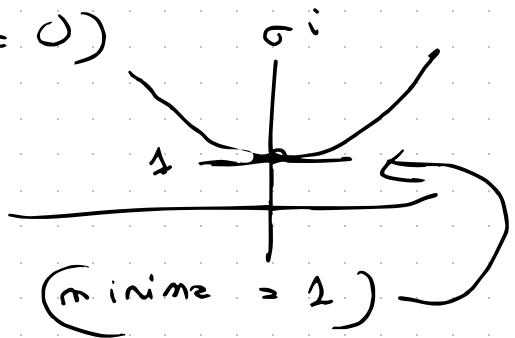
favourite  $\mu_i^i$  that are close to 0  
the func<sup>=</sup> has  $\rightarrow$  minimum  
for  $\sigma_i^i = 1$

$\Rightarrow$  KL divergence :- (look it up)

$$KL(N(\mu, \sigma) || N(0, 1))$$



(minim<sup>=</sup> = 0)



(minim<sup>=</sup> = 1)

→ This breakthrough of VAEs made them a competitive technique for dim. reduc<sup>?</sup>.

### Q) Take - Home Message :-

- (1) Relation b/w = AE & PCA
- (2) manner of introducing noise is a very relevant principle

### D) VAE v/s SNE :-

- (i)  $y^i$ 's are explicit func<sup>=</sup> of  $x^i$  in VAE
- (ii) VAE's dim reduc<sup>?</sup> can be inverted.

So, VAE is better for interpretability than SNE & other methods (except for PCA)  
-  $(y^i)$  is explicit func<sup>=</sup> of  $(x^i)$