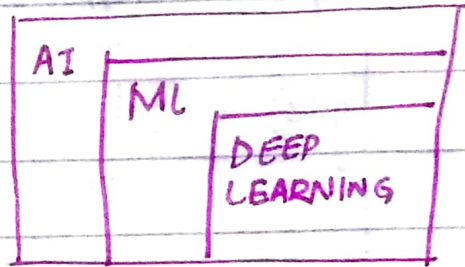# MIT : Introduction to deep learning 6.S191

→ Tensorflow Implementation



AI : Any technique that enables computers to mimic human behavior.
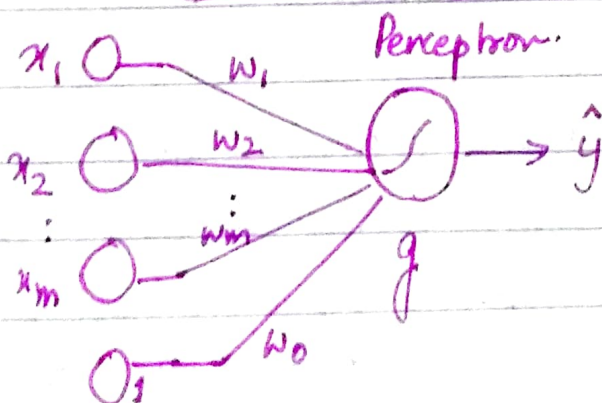
ML: Ability to learn without being explicitly programmed.

Deep learning (DL) : Extract raw patterns directly from data.

### Core Idea of deep learning : Can we learn underlying patterns directly from data ?

Why Now? → Data available (collected 4 stored)

→ Availability of computing power

→ Improved software tools ( Tensorflow)

## Perceptron : The structural building block of deep learning.
(Forward propagation)



$$\hat{y} = g \left( w_0 + \sum_{i=1}^{m} x_i w_i \right)$$

Non linear activation function

linear combi of weights 4 inputs

'g': Non linear activation functions

Examples: Sigmoid, tanh and relu.

$$g(z) = \frac{1}{1+e^{-z}}$$

SIGMOID
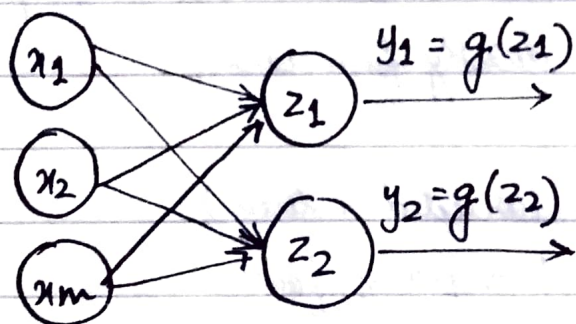
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

TANH

$$g(z) = max(0, z)$$

RELU.

: purpose is to introduce non-linearity into the network.

★: **Multi Output Perception** (2 outputs)
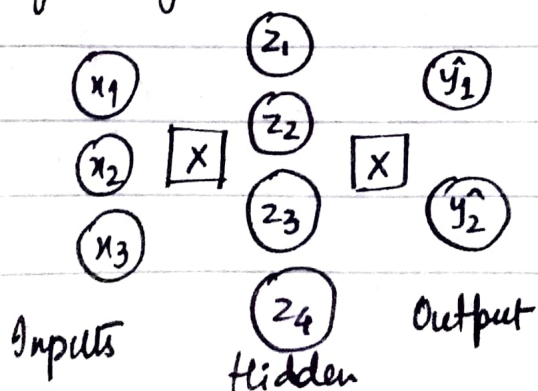


$y_1 = g(z_1)$

$y_2 = g(z_2)$

where

$$z_j = w_{0j} + \sum_{i=1}^{m} x_i w_{ij}.$$

*Dense layers*: all inputs densely connected to all outputs.

★ **Single Layer Neural Network** (1 hidden layer)



Inputs

Hidden

Output

```
import tensorflow as tf
model = tf.keras.Sequential([
    tf.keras.layers.Dense(n),
    tf.keras.layers.Dense(2)
])
```

# Deep Neural Network



$(z_1)$ $(z_{K,1})$
$(x_1)$ $(x_2)$ $(x_3)$ $x_m$
$(z_1)$ $(z_2)$ $(z_3)$ $(z_4)$
$X \ldots X$ $X \ldots X$
$(\hat{y}_1)$ $(\hat{y}_2)$
$(z_{K,n_{1k}})$

model...

```
[
tf.keras.layers.Dense(n_1),
"       "      "   .Dense(n_2),
⋮
tf.  "       "   .Dense(2).
]
```

## ✭ Applying Neural Networks

eg: Will I pass this class.

$x_1$ : # of lectures attended

$x_2$ : # of hours spent in final project.

Training: using the objective function / loss function as

$$\mathcal{L}\left(\underbrace{f(x(i); W)}_{\text{predicted}}, \underbrace{y(i)}_{\text{Actual}}\right)$$

for all instances in data :, total loss : average loss across instances

$$J(w) = \frac{1}{n}\sum_{i=1}^{n}\mathcal{L}\left(f(x^{(i)};W), y^{(i)}\right)$$

* <u>Training Objective</u> : find the best set of weights that minimises total loss.
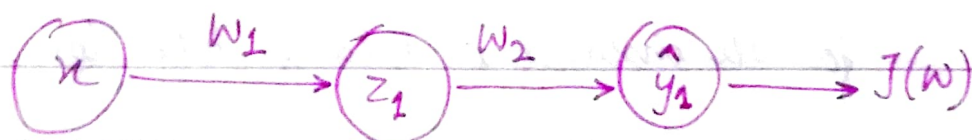
* <u>Gradient Descent</u> :

    Algorithm :

      1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
      2. Loop until convergence
      3.     Compute gradient, $\dfrac{\partial J(w)}{\partial \mathbf{3}(w)}$
      4.     Update weights (in the direction of downward slope and update weights)

$$W \leftarrow W - \textcircled{$\eta$}\dfrac{\partial J(w)}{\partial \mathbf{3}W}$$

      5.   Return Weights

                  <u>eta : learning rate</u>

* <u>Computing Gradients : Using Backpropagation</u>

Computing gradient : How much a differential change in $w$, r.t $W_2$        $W_2$ affects the loss.

$$\frac{\partial J(w)}{\partial w_2} = \frac{\partial \hat{y}_1}{\partial w_2} \times \frac{\partial J(w)}{\partial \hat{y}_1}$$

Similarly, computing gradient w.r.t $w_1$.

$$\frac{\partial J(w)}{\partial w_1} = \frac{\partial J(w)}{\partial \hat{y}_1} * \frac{\partial \hat{y}_1}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

* <u>Setting the learning rate</u> :

- decides how fast the algorithm converges
- too low value : can get stuck in local minima
- too large value : can overshoot the minima & diverge.
- setting $\boxed{eta}$ value
    → fixed : trial and error
    → adaptive - lr changes with landscape.
        • SGD - vanilla GD        • Adagrad
        • Adam                    ◦ RMS Prop.
        • Adadelta

<u>Stochastic Gradient Descent</u>
→ Gradient descent is computationally expensive. Usually a random point i picked & gradient is computed for that iter.

another random point is selected in the next iteration and the gradient is computed again. But it is very noisy. Taking a middle ground - Batch computation.

$$\frac{\partial J(w)}{\partial w} = \frac{1}{\boxed{B}} \sum_{k=1}^{B} \frac{\partial J_k(w)}{\partial w}$$
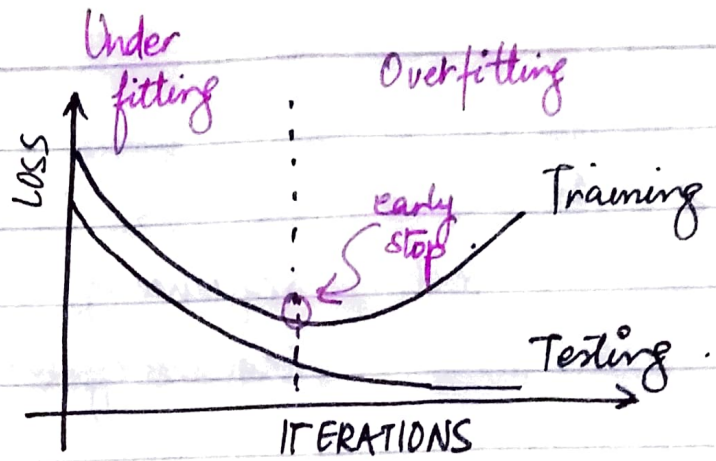
$\downarrow$
size of the batch

→ Batch computations helps by
 - giving ability to set larger learning rates
 - smoother convergence.
 - parallelize computation.

⭐ Overfitting :
Contain overfitting using regularization.

1. Dropout : Randomly setting some activations to '0'.

tf. keras. layers. Dropout (p = 0.5)
 - drops 50% of activations in layer (sampled)
 - does not memorize data.

## 2. Early Stopping :

Under fitting

Overfitting

LOSS

Training

early stop

Testing

ITERATIONS

# LECTURE: 2 : DEEP SEQUENCE MODELLING

1: Audio

2: Text :  Sequence of Characters

sequence of words.

**Example**: Sequence modelling and predicting the next word given a set of words.

**Problem 1:** Sentences can be of variable lengths.

**Idea 1:** Used fixed window ( last two words -> predict next).

**Problem 2:** Contextual info from prev. sentence or longer sentences would not be captured

**Idea 2:** Use bag of words representation (complex set of vocabulary).

**Problem 3:** Counts don't preserve order:

The food was good, not bad @ all
The food was bad, not good @ all

Same sentences (length 4 words) - diff contextual meaning.

**Idea 3:** Use really long fixed window

**Problem 4:** No parameter sharing. A phrase that starts at the start of sentence could not be distinguised or would

be considered a different meaning if that phrax starts coming at the end of sentence when parameters applied to a feed forward NN.

eg: this morning took the cat.  (1)

. . . . . . this morineng  (2)

(this morning : means the same in two sentences but could be mis represented as something else by feed forward neural network)
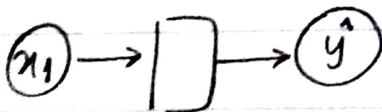
## Sequence Modelling Design Criteria
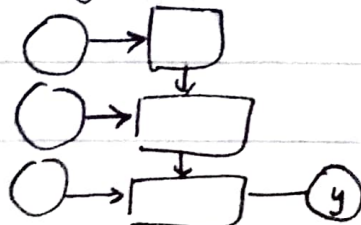
Develop sequence models that

1. can handle variable length sequences.
2. can Track long term dependencies
3. Maintain information about the order.
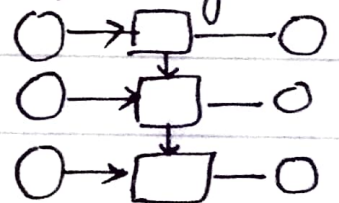4. Share parameter across the sequence.

## RNN - Recurrent Neural Network.

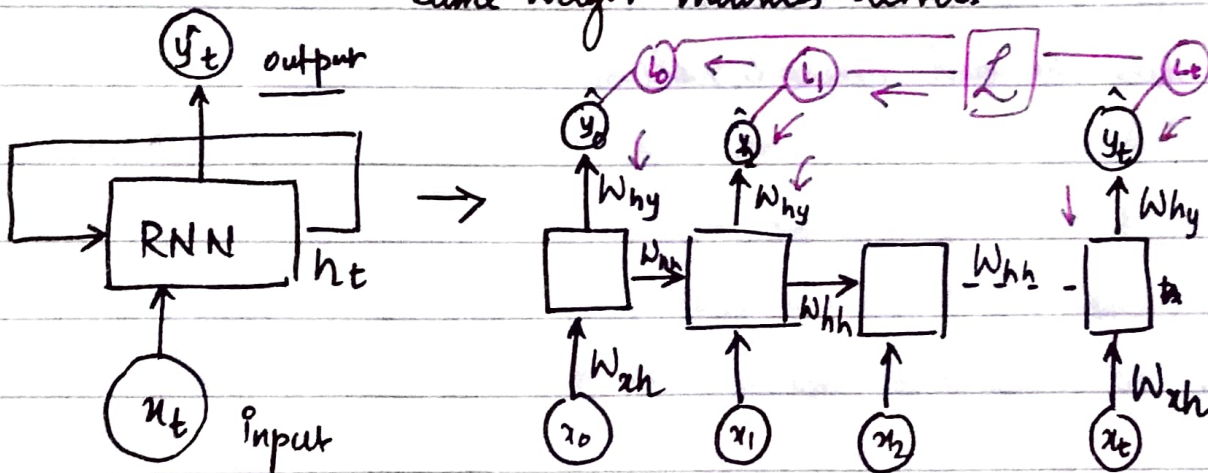| Vanilla NN | Many to One (Sentiment) | Many to Many (Music) |
|---|---|---|

Define state of the cell (Recurrent cell).

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

↓        ↓             → inputs @ $t'$

State at time $t'$    state @ prev time step

$f_W$ - parameterised by $W$

     - Same weight matrices across



Total loss $= L_0 + L_1 + \dots L_t.$     → forward pass

                            ← Backward pass.

Given    Input vector → $x_t$

update hidden state → $h_t = \tanh\left(W_{hh}^T h_{t-1} + W_{xh}^T x_t\right)$

output vector → $\hat{y}_t = W_{hy}^T h_t$

# Back Propagation Through Time (BPTT)

→ Gradient computation ②
  - each time step
  - across time steps.
→ Repeated multiplication of gradients 4 cov with weight matrix and multiple gradient computations.
→ Problems while computing Gradient
  1) Exploding Gradients - many gradients become more than 1,
     Solved by Gradient clipping, or sca i.e scaling back gradients to smaller values.
  2) Vanishing Gradients - gradient values become exceedingly smaller. 3 approaches to solve this problem.
     → Activation function
     → intelligent weight initialization
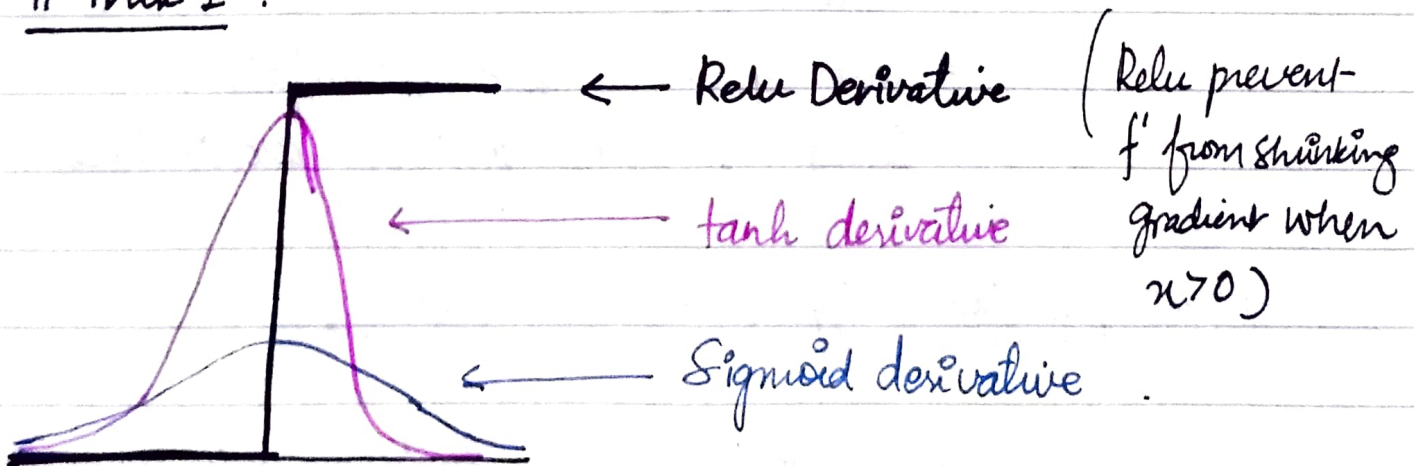     → Designing robust network architecture.

→ <u>Why is Vanishing Gradient a Problem?</u>

  → Vanishing gradients cause problems with long-term dependencies.

→ Computing gradients requires multiplying repeatedly with the weight matrices.

→ When we multiply many small numbers together
 - leads to gradient values becoming very small.
 - the error is not propagated back to the network
 - leads to capture more short-term signals and thus bias towards recent values seen by the network.
 - this is not always a problem, but could be a problem where reference or content from earlier in the sentence is required to make a prediction.

## Solving vanishing gradients problem

# Trick 1 :



← Relu Derivative   (Relu prevent-
                     f' from shrinking
← tanh derivative    gradient when
                     $x > 0$)
← Sigmoid derivative .
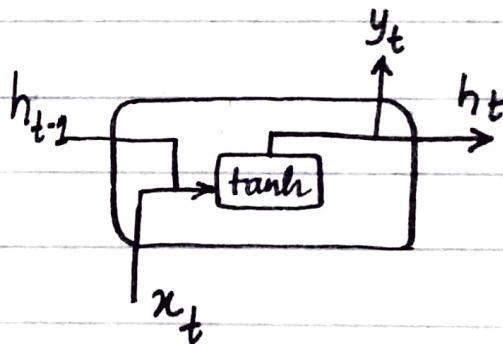
Trick #2 :    1) Initialize weights to identity matrix.
              2) Initialize biases to zero.

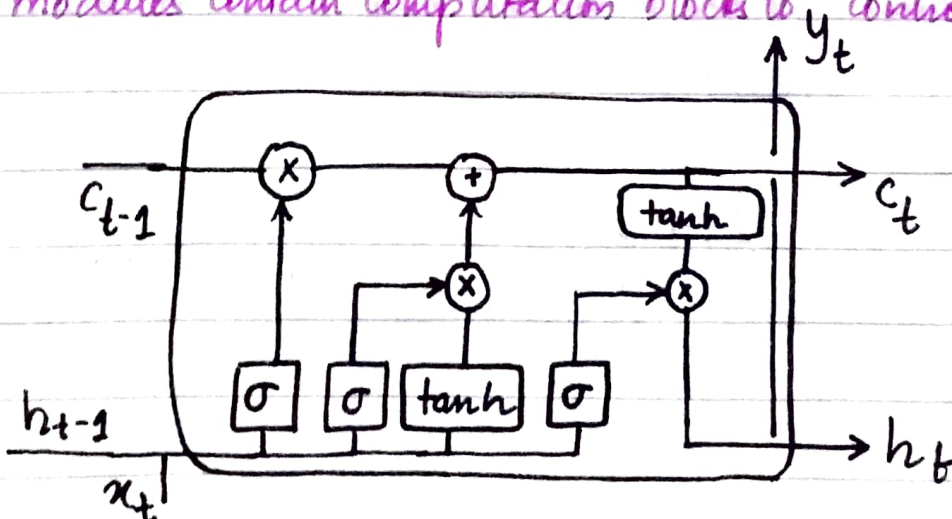Trick #3 :    Gated Cells : More Robust

→ using a more complex NN with gates to control
  which information is passed through.
  eg: LSTM, GRU

\* **long Short Term Memory (LSTM) Networks**

Standard RNN Cell



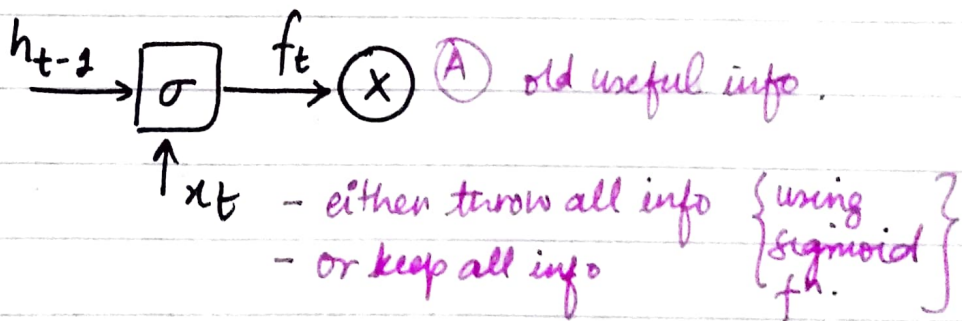LSTM : modules contain computation blocks to control flow of info.

→ Information is added or removed through gates.

# How do LSTMs Work?

1) Forget  — irrelevant information of previous state
2) Store  — relevant information by performing computation.
3) Update  — collectively update cell state values using 1)+2)
4) Output  — control what to send in the next time step.

## 1) Forgetting:

function of prev. state & current input

$$h_{t-1} \rightarrow \boxed{\sigma} \xrightarrow{f_t} \otimes \quad \text{(A)} \quad \text{old useful info.}$$

$\uparrow x_t$

— either throw all info   { using sigmoid fn.
— or keep all info

## 2) ~~B~~ Store: Decide what part of new information is relevant.

$i_t$

$$h_{t-1} \rightarrow \boxed{\sigma}$$
$$x_t \rightarrow \boxed{tanh} \rightarrow \otimes \rightarrow \oplus \quad \text{(B)} \quad \text{new useful info.}$$

## 3) Update:
using old & new info.

$$c_{t-1} \rightarrow \otimes \cdots \rightarrow \oplus \rightarrow c_t$$

**4) Output:** decide what to send out.

Key thing about LSTMs is that they use gating mechanisms to create and update internal state 'c' which leads to <u>unitterupted gradient flow</u>

## LSTM Key Concepts

→ Maintain seperate cell state than what is outputted.
→ Uses gates to control flow of info
    → forget, store, update & output.
→ Backpropagation through time for uninterrupted gradient flow.

## RNN Applications

→ Music generation
→ Sentiment clasification of tweets / text.
→ Machine translation - (one language to another).
    - Backbone of google translate RNN.
    - Attention Mechanisms to make avail
    all memory steps in all time steps.
      states

→ Self driving cars.
→ Environmental modelling - predicting wind flow
    directions.
→

# LECTURE 3 : CONVOLUTIONAL NEURAL NETWORKS

Computer Vision applications transformed by deep learning.
- → facial recognition and detection.
  - emotions, features detection.
- → Medicine, biology and healthcare
- → Self driving cars.

**★ How does computers see ?**
- Images as two-dim. numbers. - Grayscale
- RGB - 3 X 2Dim Matrices - stacked

## High level features

→ Manual feature extraction

| using domain knowledge | → | define features | → | detect features to classify |

- this is difficult in practice, since to classify an image, the pipeline should be agnostic of variations in features. These variations could be deformity, brightness, viewpoint, illumination, background & in-class variations.

→ **learning feature representation**

- learn <u>hierarchy of features</u> directly from the data
  to make a representation of what makes up the
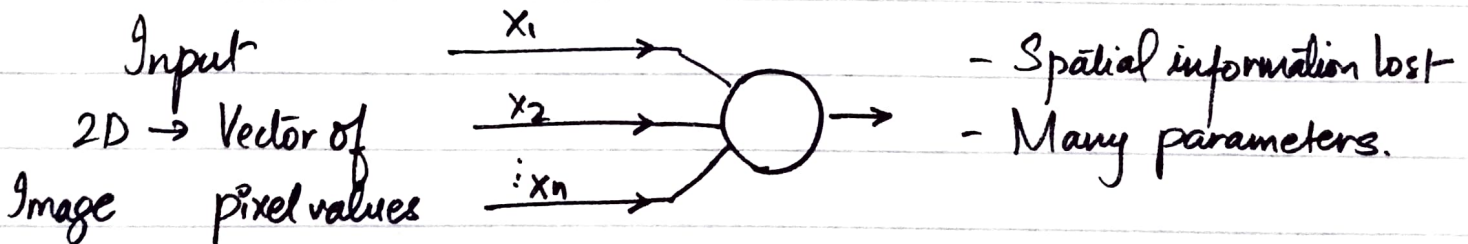  final class label.
- the hierarchy of features could include
    - low level features - darkspots, edges.
    - mid level       "    - eyes, ears, nose
    - high level      "    - facial structure.

⭐ <u>How to learn features?</u>

Input
2D → Vector of
Image    pixel values

$x_1$
$x_2$
$\vdots x_n$

- Spatial information lost
- Many parameters.

→ <u>Preserving spatial structure using 'convolution'</u>

Instead of passing entire 2D image as array of numbers,
pass patches of input and connect to neurons in hidden layer.
using sliding window.

→ learning feature representation

- learn hierarchy of features directly from the data to make a representation of what makes up the final class label.
- the hierarchy of features could include
  - low level features - darkspots, edges.
  - mid level " - eyes, ears, nose
  - high level " - facial structure.

* How to learn features?

Input
2D → Vector of
Image    pixel values

$x_1$
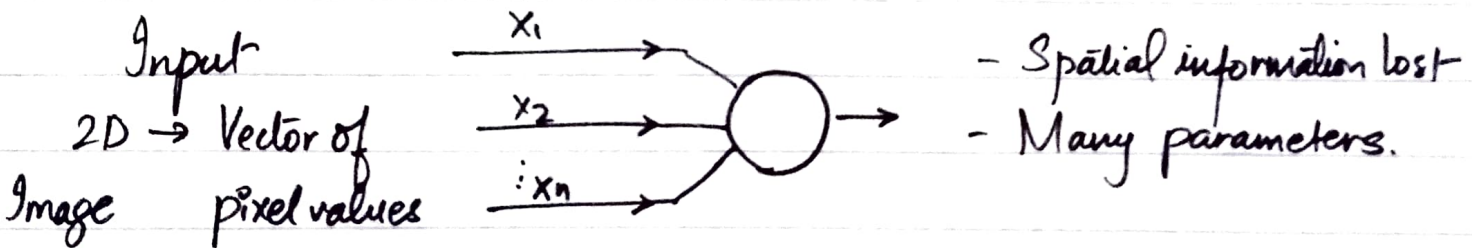$x_2$
$: x_n$

- Spatial information lost
- Many parameters.

→ Preserving spatial structure using 'convolution'

Instead of passing entire 2D image as array of numbers, pass patches of input and connect to neurons in hidden layer, using sliding window.
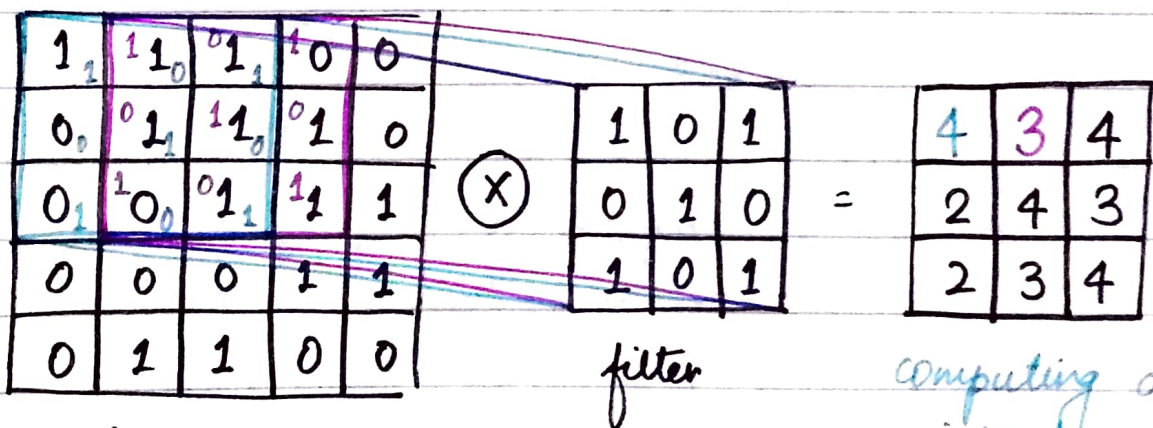
→ Weighting the patch to detect particular features.

## Feature Extraction by applying filters.

→ apply - a set of weights to extract local features.
→ use multiple filters to extract different features.
→ spatially share parameters of each feature.

## The Convolution Operation

Computing a convolution of $5\times5$ image using $3\times3$ filter.

| | | | | |
|---|---|---|---|---|
| $1_1$ | $1_1$ | $1_1$ | 0 | 0 |
| $0_0$ | $1_1$ | $1_1$ | $0_1$ | 0 |
| $0_1$ | $0_0$ | $1_1$ | 1 | 1 |
| 0 | 0 | 0 | 1 | 1 |
| 0 | 1 | 1 | 0 | 0 |

Image

$\otimes$

| | | |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

filter

$=$

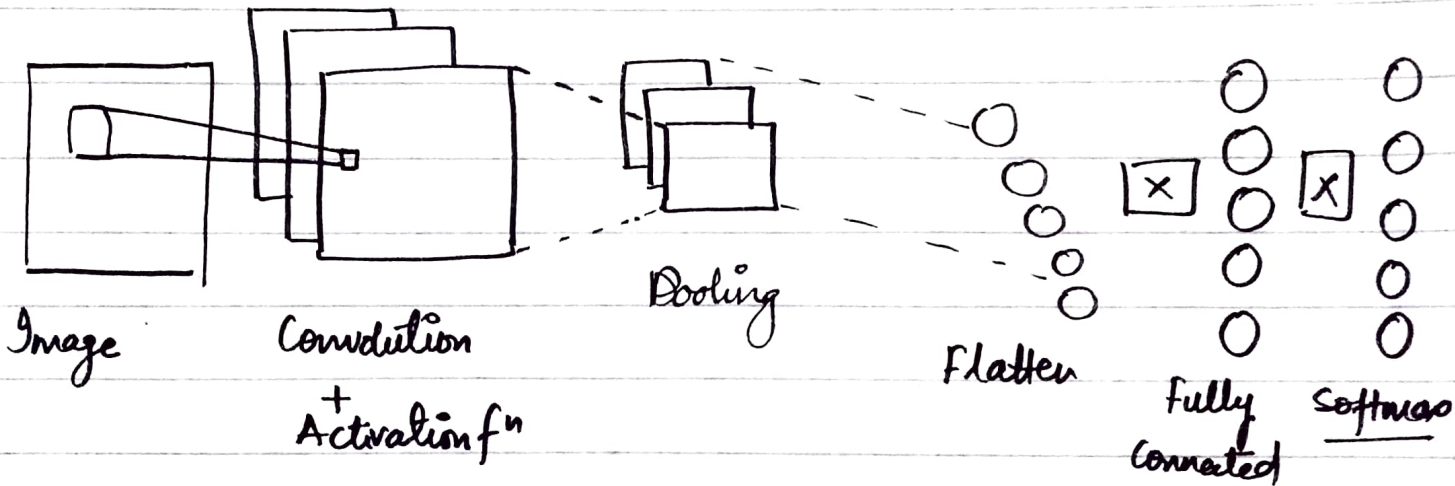| | | |
|---|---|---|
| 4 | 3 | 4 |
| 2 | 4 | 3 |
| 2 | 3 | 4 |

computing a weighted sum by element-wise operation.

The diagonal gives max. output at the diagonals. This means that this filter gives maximum overlap on this image along this diagonal.

→ changing the weights in the filter can give you a different output feature map.

\* <u>Convolutional Neural Networks</u>

CNN for image classification.



Image    Convolution    Pooling    Flatten    Fully    Softmax
         +                                    Connected
         Activation fⁿ

3 broad steps to classification.

→ <u>Convolution</u> - applying filters to extract features
    - using multiple filters to extract a hierarchy of features

→ <u>Non Linearity</u> - apply non-linearity

→ <u>Pooling</u> - Downsampling the learned features
    - reduces the dimensionality
    - makes pipeline spatially invariant.

1) **Convolution :** - Take inputs from a patch  
                - compute weighted sum    } For a neuron  
                - apply bias.             in hidden layer.

→ Apply a window of weights.  
→ compute linear combinations.  
→ Activating with non-linear function.

layer dimensions - $h \times w \times d$ → depth (# of filters)

            Spatial dim of image

stride - filter step size  
Receptive field - locations in input image that a node in path connected to.

2) **Non linearity** - ReLU      $g(z) = \max(0, z)$.  
     - pixel-by-pixel operation that replaces all negative values by zero.  
     - applied after every convolution operation.

3) **Pooling:**

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

⟹

| 6 | 8 |
|---|---|
| 3 | 4 |

maxpool = 2×2 filters.

The image classification pipeline consists of two parts
= <u>Feature learning</u> + classification task.

|

can be plugged into tasks like
1) <u>Image detection</u>
2) <u>Semantic segmentation</u> - identifying class of pixel
   (whether pixel belongs to foreground, background,
   object or different objects). -
3) <u>End-to-end robotic control</u> - self driving cars.

# LECTURE 4: DEEP GENERATIVE MODELS

<u>Generative modelling</u> : unsupervised learning

<u>Goal</u> : Take as input training samples from some distribution and learn a model that represents that distribution.

learn Model $P_{model}(x)$ similar to $P_{data}(x)$

1. Density Estimation - learning prob dist
2. Sample Generation - generate new samples of data from learned training data.

Applications to
  → generate representative datasets for debiasing models.
  → outlier detection - generative models for outlier detection. train on outlier generated to improve model performance on rare circumstances.

<u>Two classes of Modes</u>: Latent Variable Models
  → GAN - General Adverserial Networks.
  → Autoencoders and Variational Autoencoders (VAEs)
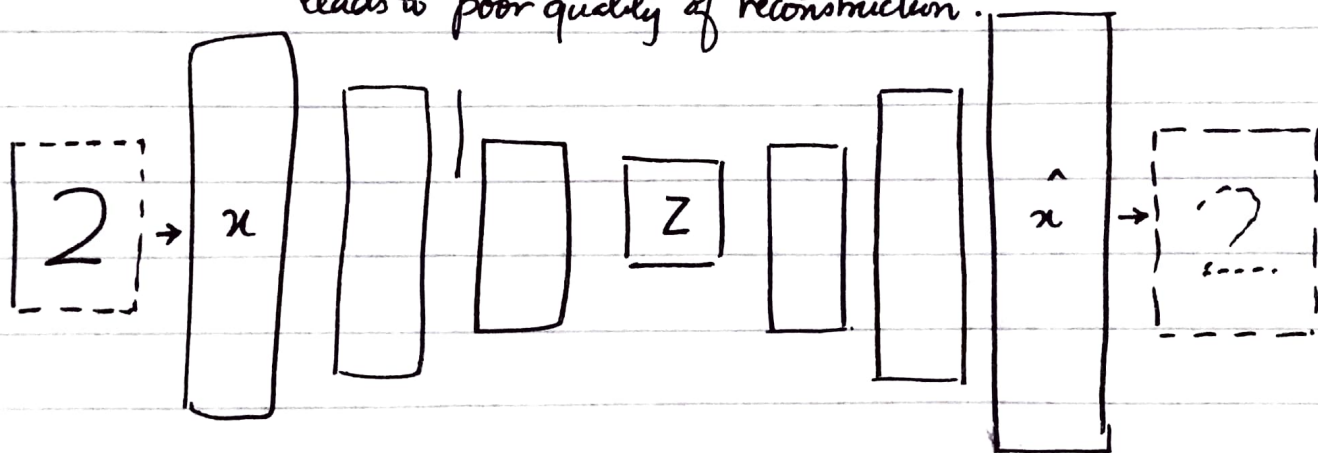
## Latent Variable vs Observed Variables

↓

variables that are ~~not~~ true explanatory factors
that are /~~can be only~~ observed in the data · only inferred
cannot be  through a model from other obs · variables

Goal of Generative modelling - is to learn these
latent variables from only observed data.

---

**\* AUTOENCODERS** - Simple generative model.

→ autoencoding is a form of compression · lower dimensions
leads to poor quality of reconstruction.



Input

ENCODER

Encodes information in Input
into lowdimensional output '$z$'.

DECODER.

Train CNN/Fully connected layer
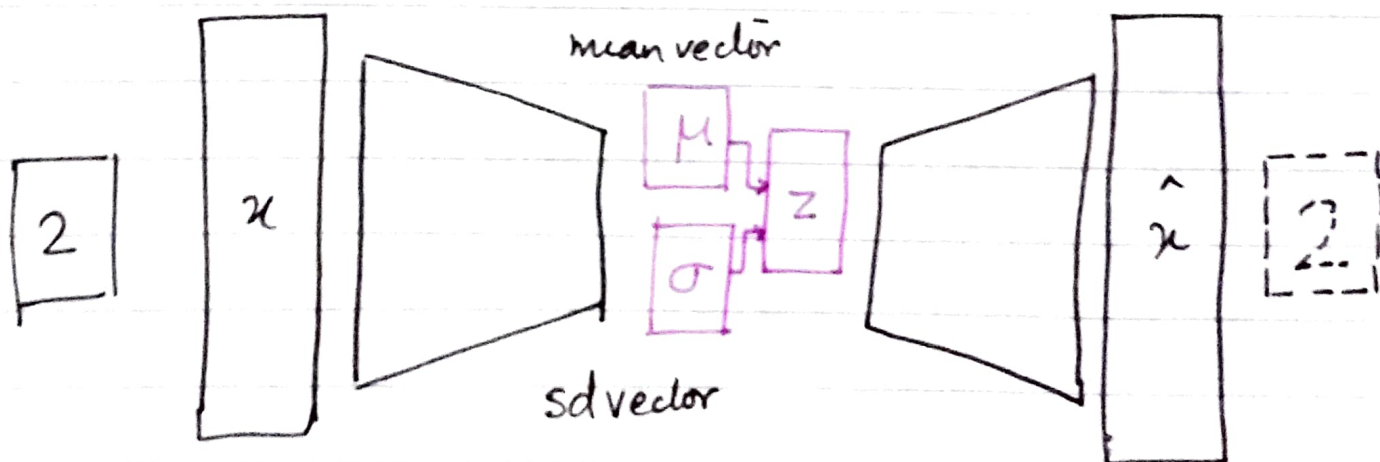to predict-/reconstruct training
data.

Reconstruction loss

$$\mathcal{L}(x - \hat{x}) = \|x - \hat{x}\|^2 \qquad (MSE).$$

In Summary, Autoencoders

→ use <u>bottleneck hidden layer</u> to compress input and learn latent representation.

→ Reconstruction loss forces latent representations to encode max. possible info about the data.

## Variational Autoencoders (VAEs).

- replaced bottleneck layer (deterministic) → stochastic sampling.
- for generating similar data.
- Instead of learning latent variable 'z', the model learns mean and std. dev for each variable, that parameterises a prob. distribution for each of those latent variables.



mean vector

sd vector

**Encoder**: computes prob. dist ~~of latent~~ $\phi$ of the latent space $z$ given an input $x$.

$$- P_\phi(z|x)$$

**Decoder**: Computing new prob dist $q_\theta$ of $x$ given latent variables $z$.

$$- q_\theta(x|z)$$

$$\mathcal{L}(\phi, \theta, x) = (\text{reconstruction loss}) + (\text{regularization term}).$$

$\downarrow$ Same as before $\quad\quad\quad\downarrow$

$$||x - \hat{x}||^2 \quad\quad D\left(P_\phi(z|x) \,||\, p(z)\right)$$

$\downarrow$ Inferred latent dist

- Some initial hypothesis of what $z$ looks like
- to help network not overfit.

<u>Priors on the latent distribution</u>

$$D\left(P_\phi(z|x) \,||\, p(z)\right)$$

Inferred latent distribution $\longleftarrow \quad\quad \longrightarrow$ Fixed prior on latent distribution

Common choice of prior — Gaussian distribution.
$$\mu = 0; \sigma = 1.$$

$=$ KL Divergence when gaussian. (with $\mu = 0, \sigma = 1$)

$$-\frac{1}{2} \sum_{j=0}^{k-1} \left( \sigma_j + \mu_j^2 - 1 - \log \sigma_j \right)$$

[6] Basically learning latent variables with a known distribution?.

{Machine learning mastery}

- Kullback-Leibler divergence - a measure of how one probability distribution is different from a second one.
- also called relative entropy
- can be used for feature selection (information gain).
- or cross entropy used as loss function.
- P and Q, score for divergence $\neq$ Q and P score.

  • KL $(P \| Q)$

  ↑
  'divergence' of P from Q

$$\text{sum}\left( p[i] * \log_2 \left( \frac{p_i}{q_i} \right) \right) \text{ for all values in } p.$$

available in scipy library. use rel_entr() instead of kl_div()

# Jenson-Shannon Divergence

→ JS divergence is symmetrical i.e
$$JS(P||Q) == JS(Q||P)$$

→ $$JS(P||Q) = \frac{1}{2} KL(P||M) + \frac{1}{2} KL(Q||M)$$

$$M = \frac{1}{2}(P+Q)$$

→ More useful measure as it provides more smooth and normalized version of KL.

→ Score = 0 (Identical distributions)
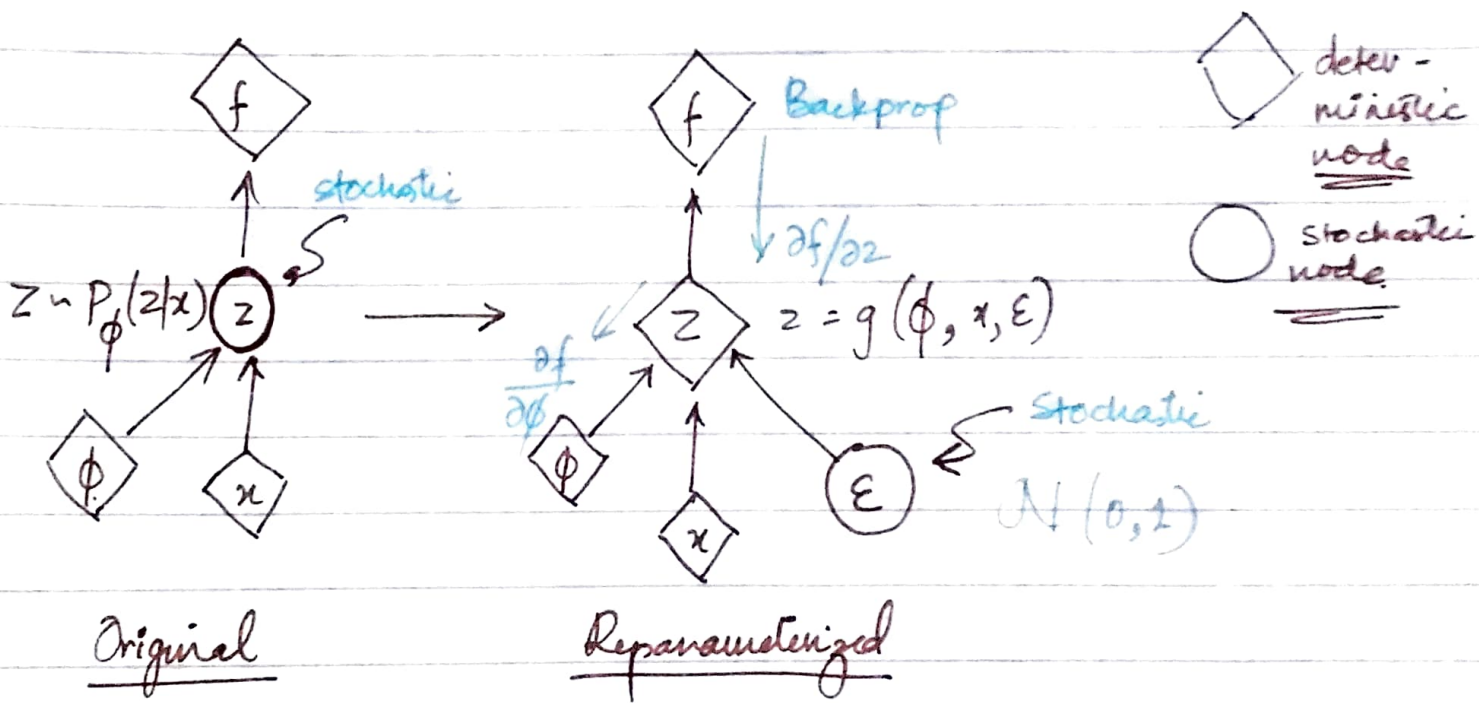= 1 (Maximally different)
when using log base 2.

# Training VAEs using Backpropagation

- $z \leftarrow$ is a result of stochastic sampling.

- **problem** : cannot compute gradient over non-deterministic nodes.

- **Solution** : use stochasticity as an input to a determinant node and compute gradient.

**\*** <u>Reparameterize the sampling layer</u>

$$\boxed{z = \mu + \sigma \odot \varepsilon} \quad \text{instead of} \quad \boxed{z \sim \mathcal{N}(\mu, \sigma^2)}$$

$$\varepsilon \sim \mathcal{N}(0, 1)$$



Original

Reparameterized

**\*** <u>What do these latent variables mean?</u>

→ Keeping all latent variables same, sample one variable from the prior distribution and tune the variables to find out what they have learnt

→     tuned variable ⟶ Decoder ⟶ output

      (face)                     (orientation of faces)

                                eg <u>latent? – orientation</u>

→ latent variables should be independent / uncorrelated.
 called disentanglement
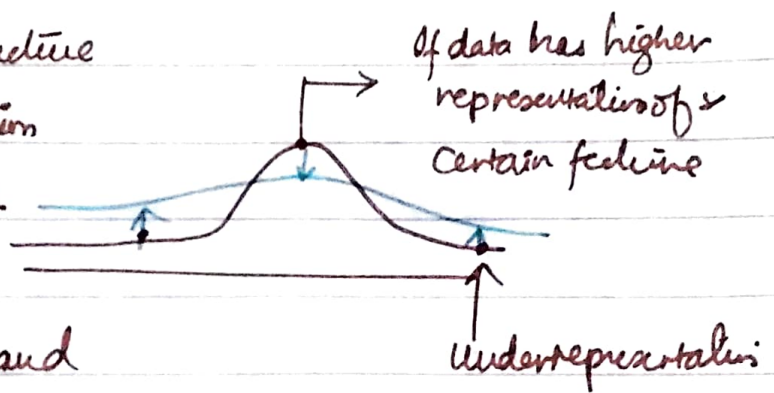
\* <u>Mitigating Bias through learned latent structure</u>

Step 1: Learn latent structure
Step 2: Estimate distribution
Step 3: Adaptively Resample
Step 4: ~~to ge~~ Generate
 diverse set of data and
 fit a Fully connected layer
 to the data.

Of data has higher
representation of x
certain feature

underrepresentation

\* <u>General Adversarial Network</u> (GAN)

<u>Idea</u>: Generate new instance similar to the training data.
 → Not learn the underlying dist 4 latent var. distribution.
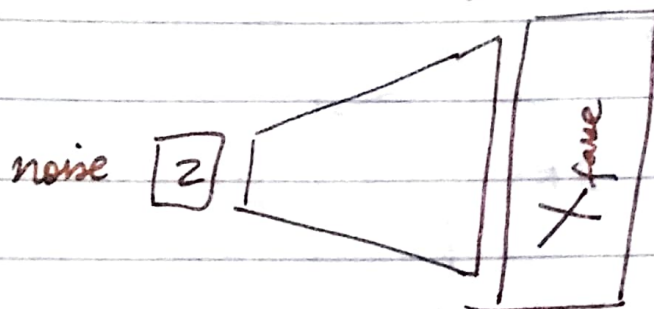
<u>Problem</u>: Cannot directly sample from complex dist

<u>Solution</u>: Generator network starts from random noise.
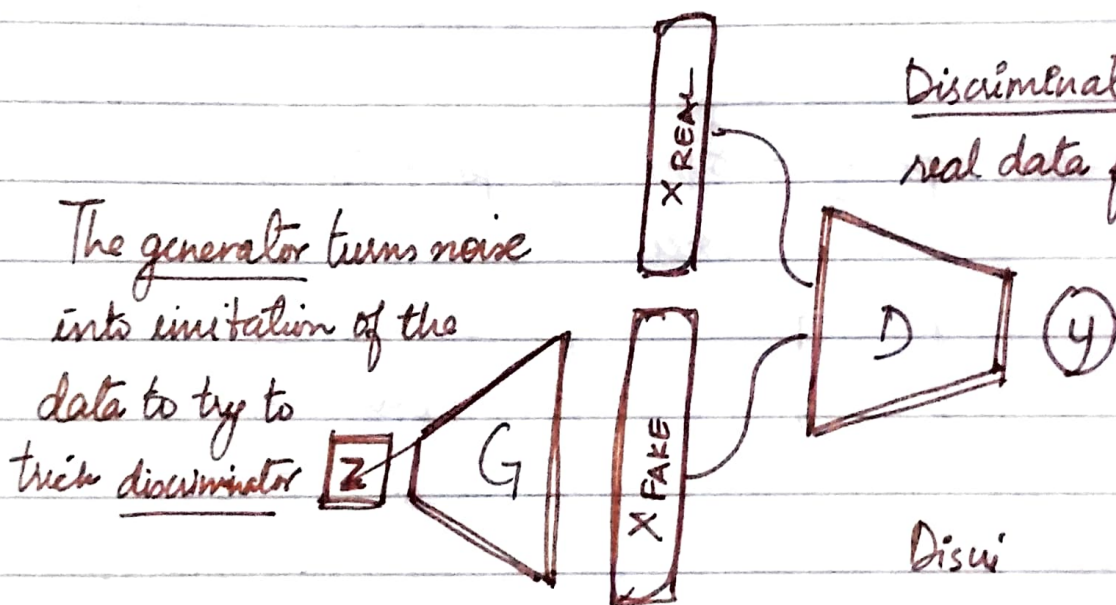 then " trained going from noise to learning a

transformation to the training distribution.

Goal:

noise $\boxed{Z}$ $X_{fake}$

To generate a data that is very similar to the training data.

The generator turns noise into imitation of the data to try to trick discriminator $\boxed{Z}$ G
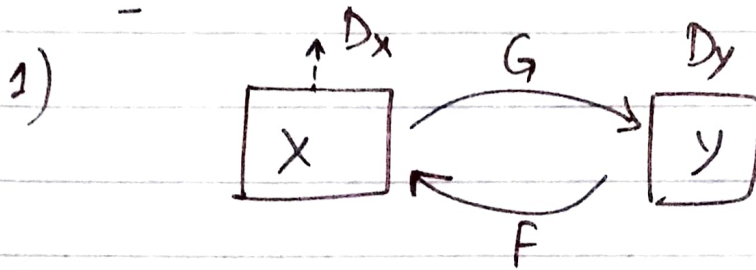
$X_{REAL}$

$X_{FAKE}$

D $\textcircled{y}$

Discrimenator tries to identify real data from fakes created by generator

Discri

After training — Trained generator network to sample and create new data never seen before.

Applications → 1. Progressively growing GANs.
- low resolution starting
- incrementally add resolution during training

## 2. CycloGAN : Style Transfer

1)



learns from domain X and Y. — Two generators and discriminators Working.

2) Transfer Speech — Synthesise speech in someoneelse voice



Audio waveform

Spectogram (A) ⇄ Spectogram (B)

cycleGAN