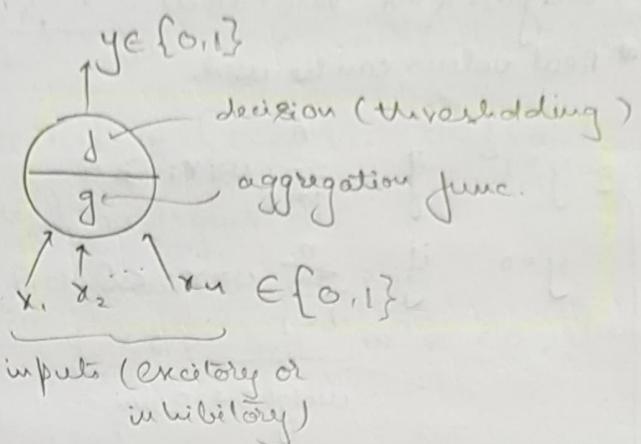
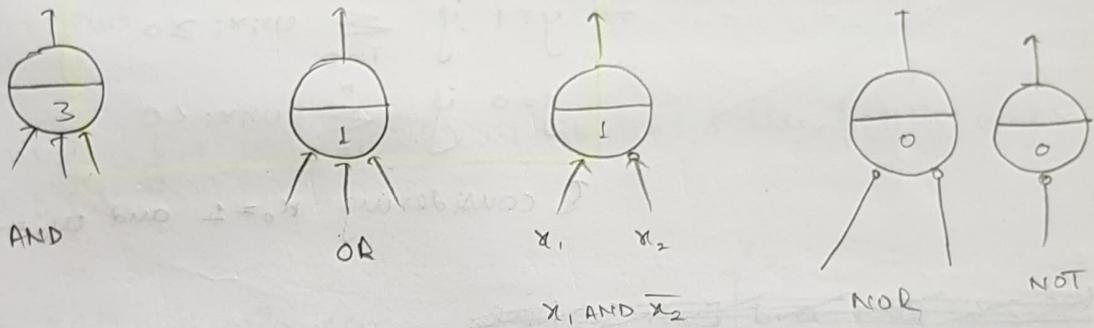


Lec 2.2: McCulloch Pitts Neuron (Thresholding logic)

$$g(x_1, x_2, x_3, \dots, x_n) = \sum_{i=1}^n w_i$$

$$y = f(g(x)) = 1 \text{ if } g(x) \geq 0 \\ = 0 \text{ for } g(x) < 0$$



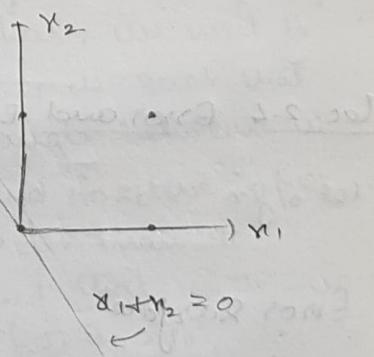
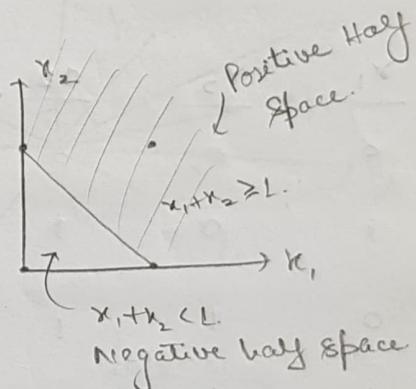
Geometric Interpretation:

OR: $x_1 + x_2 \geq 1 \Rightarrow y = 1$
 $x_1 + x_2 < 1 \Rightarrow y = 0$

* MP Neurons learn a linear decision boundary.

Tautology (always on) $x_1 + x_2 \geq 0 \text{ for } g \geq 0$.

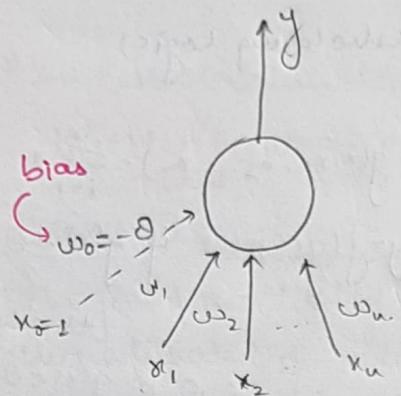
* MP Neurons can represent boolean functions which are linearly separable

Lec 2.3: Perceptions

What about non-boolean (real) inputs?

How to assign more weights to some inputs?

How to deal with non-linearly separable?



Here ~~w₀~~ w_0 is called the bias.

* weighted inputs and possible to learn these weights (not hard coded)

* Real values can be used.

$$y = 1 \quad \text{if} \quad \sum_{i=1}^n w_i x_i \geq 0$$

$$y = 0 \quad \text{if} \quad \sum_{i=1}^n w_i x_i < 0$$

weighted sum.

$$y = 1 \quad \text{if} \quad \sum_{i=1}^n w_i x_i - \theta \geq 0$$

$$\Rightarrow y = 1 \quad \text{if} \quad \sum_{i=0}^n w_i x_i \geq 0$$

$$y = 0 \quad \text{if} \quad \sum_{i=0}^n w_i x_i < 0$$

Considering $x_0 = 1$ and $w_0 = -\theta$

~~lec 2.4 Error and Error surface~~

x_1 x_2 OD.

$$0 \quad 0 \quad 0 \quad \sum_{i=0}^2 w_i x_i < 0 \Rightarrow w_0 < 0$$

$$1 \quad 0 \quad 1 \quad " \geq 0 \Rightarrow w_0 + w_1 \geq 0$$

$$0 \quad 1 \quad 1 \quad " \geq 0 \Rightarrow w_0 + w_2 \geq 0$$

$$1 \quad 1 \quad 1 \quad " \geq 0 \Rightarrow w_0 + w_1 + w_2 \geq 0$$

Lec 2.4 Error and Error Surface:

let's fix $w_0 = -1$

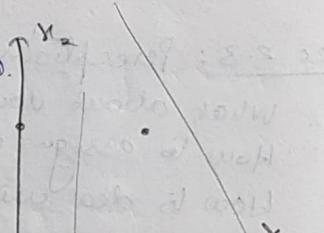
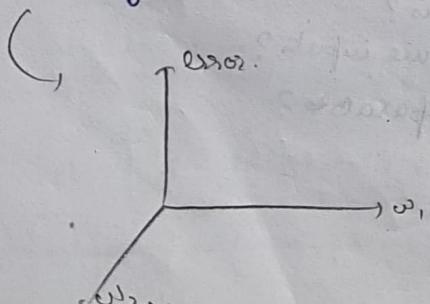
Solve these to get the linear threshold.

choose w_1 and w_2

such that

Error is minimized.

Error Surface



for different values of w_1, w_2

Lecture 2.5: Perceptron Learning Algorithm

* Adjusting the weights of the input of a perceptron.

P: inputs with label 1 (positive inputs)

N: " " " 0 (-ve inputs)

Initialize $w = \{w_0, w_1, \dots, w_n\}$ randomly

while ! convergent do

Pick random $x \in P \cup N$:

if $x \in P$ and $w^T x < 0$, then

$$w = w + x;$$

end

if $x \in N$ and $w^T x > 0$, then

$$w = w - x;$$

end

end.

We need to find $w^T x$ such that it divides the positive and -ve inputs.

for $x \in P$ and $w^T x < 0$

$$\Rightarrow \cos \alpha < 0$$

$$\Rightarrow \alpha > 90^\circ$$

and we want $\alpha < 90^\circ$

So modify w

$$w_{\text{new}} = w + x$$

$$\cos(\alpha_{\text{new}}) < \cos(\alpha_{\text{old}})$$

$$\propto (w+x)^T x$$

$$\propto w^T x + x^T x$$

$$\propto \cos \alpha + \text{a small number}$$

$$\Rightarrow \cos(\alpha_{\text{new}}) > \cos \alpha.$$

$$\Rightarrow \alpha_{\text{new}} < \alpha$$

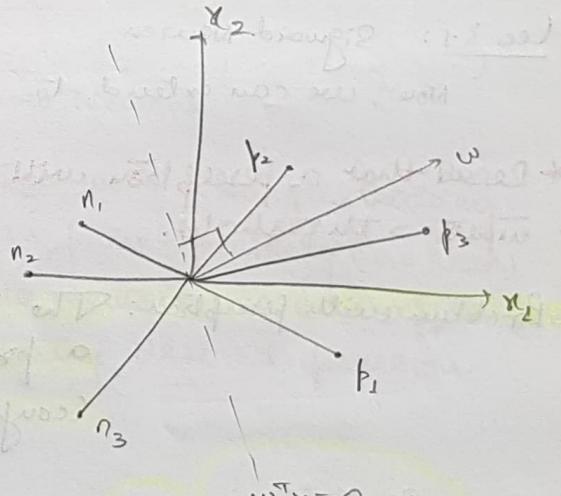
Thus we decrease α

repeatedly.

Similarly, for $x \in N$ and $w^T x > 0$

$$\alpha_{\text{new}} > \alpha.$$

we increase α repeatedly.



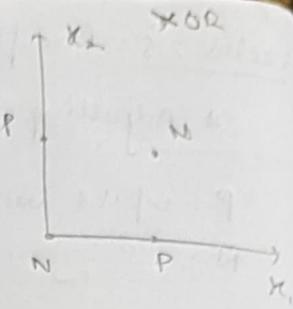
Thus, we need to find w such that angle between w and positive inputs is less than 90° and that with -ve inputs is $> 90^\circ$.

$$\cos \alpha = \frac{w^T x}{\|w\| \|x\|}$$

Lec 2.7 : Linearly Separable Boolean Func:

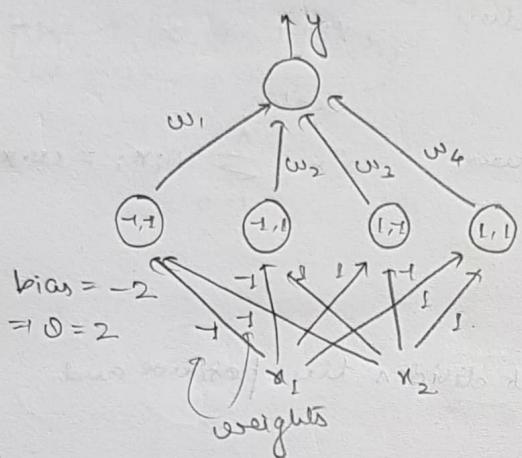
XOR , and XNOR are non linearly separable.

Multiple perceptions are needed to implement non-linearly separable.



Lec 2.8: Network of Perceptrons:

Implementing XOR :



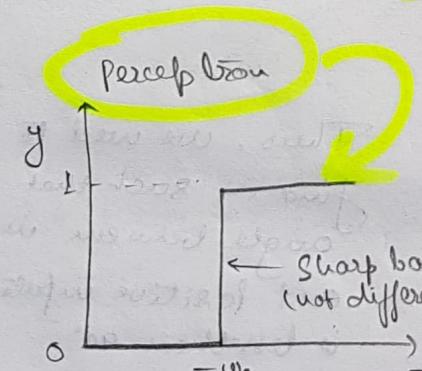
$$\text{For } w_2, w_3 = 1 \\ \text{and } w_1, w_4 = 0$$

Lec 3.1: Sigmoid Neuron:

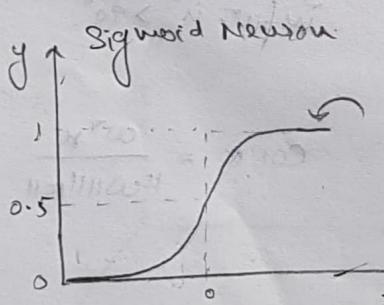
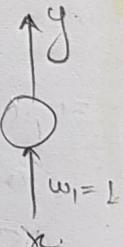
Now, we can extend to $y = f(x)$ $\rightarrow x \in \mathbb{R}^n$ and $y \in \mathbb{R}$ (non-binary).

* Recall that a perception will fire if the weighted sum of inputs $>$ threshold.

Problem with perception: The output is 1 or 0. We don't have a probability for each output (confidence is not known).



Eg: $-w_0 = 0.5$



$$\text{logistic function (differentiable)} \\ y = \frac{1}{1 + e^{-w^T x}}$$

$y = 1 \text{ if } x > 0.5$
 $y = 0 \text{ if } x \leq 0.5$
 $x = 0.49, 0.51$
 are very close but still give very different outputs.

logistic func. output lies between 0 and 1 (analogous to prob.).

Lecture 3.2 Supervised Learning:

$$\hat{y} = \frac{1}{1 + e^{-w^T x}}$$

logistic regression

$$\begin{aligned} \text{Weight vector} & w \\ \hat{y} &= w^T x \\ \hat{y} &= x^T w x \end{aligned}$$

linear reg.

quad. reg.

Given x and \hat{y} , we need to estimate w .

Loss func.: $L(w) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$

Lecture 3.4: Gradient Descent

$$w = w - \eta \nabla w$$

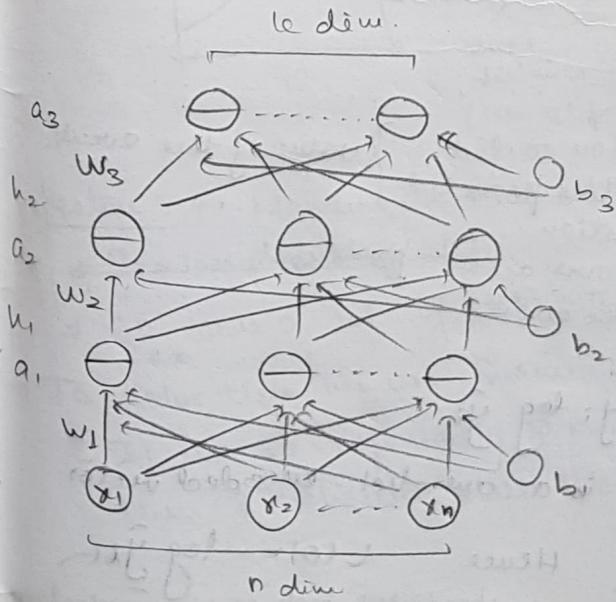
$$b = b - \eta \nabla b$$

gradient of bias
learning rate

Lec 3.5: Multilayer Network of Sigmoid Neurons

* A multi-layer network of perceptrons with a single hidden layer can represent any Boolean func. precisely (no error) whereas for the case of Sigmoid Neurons, we can approximate any func. (continuous) to any desired precision.

Lecture 4.1: Feed forward NN:



Preactivation:

$$a_i(x) = w_{i-1}x + b_i$$

Activation:

$$h_i(x) = g(a_i(x))$$

element wise (can be operation sigmoid, tanh etc)

Activation of output layer:

$$h_o(x) = O(a_o(x))$$

can be softmax, linear etc

$$h_1(x) = g(w_1^T x + b_1)$$

$$h_2(x) = g(w_2^T h_1(x) + b_2)$$

$$h_3(x) = \sigma(w_3^T h_2(x) + b_3)$$

Lecture 4.3: Output func. and loss func.

For regression problems, we want the outputs to be real numbers and unbounded. Then, σ can be a linear func.

$$\hat{y} = h_L = \sigma(a_L) = w_0 a_L + b_0$$

For classification problems, the output is a probability distribution over all the classes. Then, σ ~~should~~ should be softmax. we can't use softmax because the sum of the individual prob. will not equal 1.

$$\hat{y} = [\hat{y}_1, \hat{y}_2, \dots, \hat{y}_k]$$

$$a_L = [a_{L1}, a_{L2}, \dots, a_{Lk}]$$

$$\hat{y}_j = \frac{e^{a_{Lj}}}{e^{a_{L1}} + e^{a_{L2}} + \dots + e^{a_{Lk}}}$$

Loss function for classification cross entropy

$$\text{Information: } I(S) = -\log_2 P(S)$$

$$\text{Entropy: } H(S) = -\sum_i P(S_i) \log P(S_i)$$

Cross entropy:

P → true distribution.

Q → predicted distribution.

↳ tells how close predicted dist. is to the true dist.

$$= -\sum_i P_i \log(Q_i)$$

↳ information content based on the predicted distribution
 ↳ value of the event
 ↳ summing over actual prob. dist.
 ↳ prob. of the event

$$L(\sigma) = -\sum_{i=1}^k y_i \log \hat{y}_i$$

But for classification, y is a one-hot-encoded vector

thus $y_i = 1$ if $i = \text{correct class}$
 $y_i = 0$ otherwise. Hence $L(\sigma) = -\log \hat{y}_{\text{correct}}$

Thus we minimize $-\log p_e$ or maximize $-L(\theta) = \log p_e$

probability that x belongs to class-1
(log likelihood of data)

Minimizing the cross entropy:

minimize $\sum_i -p_i \log q_i$ under the constraint $\sum q_i = 1$

using Lagrangian multipliers:

$$\min_{q_i} \left[-\sum_i p_i \log q_i + \lambda (\sum q_i - 1) \right]$$

$$\text{Diff} \quad \frac{\partial}{\partial q_i} \left(-\sum_i p_i \log q_i + \lambda (\sum q_i - 1) \right) = 0$$

$$\Rightarrow p_i = \lambda q_i$$

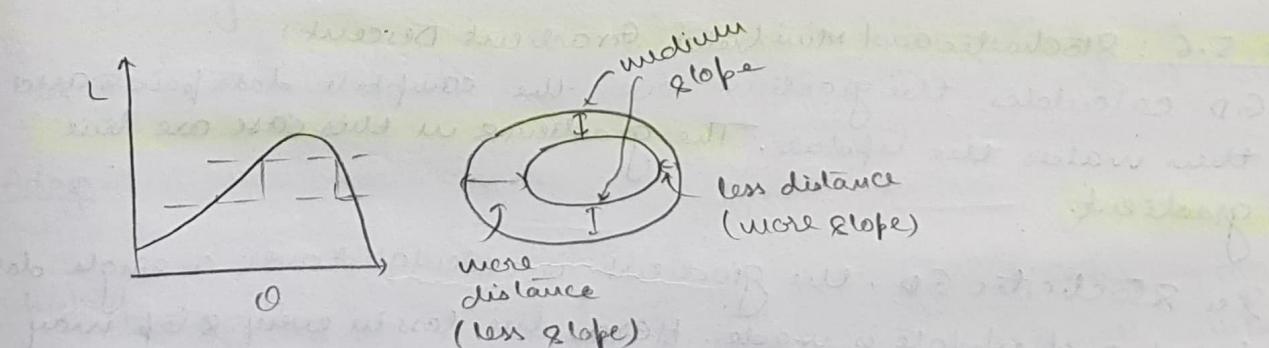
$$\text{Here, } \lambda = 1 \Rightarrow p_i = q_i$$

Lecture 5.1: Gradient Descent.

$$\theta = \theta - \eta \nabla \theta$$

= When gradient (slope) is large update is large and vice versa

Lecture 5.3: Contour Maps



Lecture 5.4: Momentum Based Gradient Descent:

* ~~takes a lot of steps~~

* GD takes a lot of steps to navigate regions of gentle slope.

To solve this we use momentum:

If we are repeatedly taking steps in a particular direction, then we should start taking bigger steps in that direction.

$$\theta_{t+1} = \theta_t - \eta \nabla \theta_t - r \text{ update}_{t-1}$$

$$\theta_{t+1} = \theta_t - \text{update}_t$$

$$\boxed{\text{update}_t = \eta \nabla \theta_t + r \cdot \text{update}_{t-1}}$$
 current update depends upon prev update

$$\text{update}_0 = 0$$

$$\text{update}_1 = \eta \nabla \theta_1$$

$$0 < r < 1$$

$$\text{update}_2 = \eta \nabla \theta_2 + r \cdot \eta \nabla \theta_1$$

Momentum term

$$\text{update}_3 = \eta \nabla \theta_3 + r \eta \nabla \theta_2 + r^2 \eta \nabla \theta_1$$

more weight exponentially weighted avg. less weight

- * Due to build up of momentum, the loss overshoots (problem) but still converges faster than GD.

↓ oscillation

lec 5.5: Nesterov Accelerated Gradient (NAG):

- * Solves the oscillating loss problem.
- * Look before you leap.
- * first move in the intended direction according to the history and check the derivative at this lookahead point and move accordingly. Thus oscillations are reduced.

$$\theta_{\text{lookahead}} = \theta_t - r \cdot \text{update}_{t-1}$$

$$\text{update}_t = r \cdot \text{update}_{t-1} + \eta \cdot \nabla \theta_{\text{lookahead}}$$

$$\theta_{t+1} = \theta_t - \text{update}_t$$

lec 5.6 : Stochastic and Mini-batch Gradient Descent:

GD calculates the gradient over the complete data points and then makes the update. The gradients in this case are true gradient.

In Stochastic GD, the gradient is calculated over a single data point and update is made. Here, the loss in every step may not decrease. We will have oscillations in loss during training.

- * In batch GD, we take a random batch of data points and calculate gradients (close to true gradients) and update is done. It has faster convergence than GD.
- * It has lower oscillations than SGD.

Lec 5.7: Tops for adjusting learning rate and momentum

Learning Rate annealing:

- Step decay
- Half-the learning rate after an epoch if the valid. loss is more than that at the prev. epoch.

Suboptimal Soln:

- To tackle local minima, start again with different init. learning rate and find the optimal loss. Repeat it several times and choose the one with min. loss.

Momentum:

- We start with a small value of γ (discounting more on history). After every step update, γ is increased (having more faith in history)

Lec 5.9: Gradient Descent with adaptive learning Rate

when any of the features is sparse (zero for most data points), the weights corresponding to that input feature will not get sufficient updates compared to dense features. This sparse feature could be very important and must be learnt properly. we need to increase the learning rate for sparse features.

Decay the learning rate according to the parameters update history

Adagrad:

$$v_t = v_{t-1} + (\nabla \theta_t)^2 \quad \text{accumulate the magnitude of gradient}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * \nabla \theta_t$$

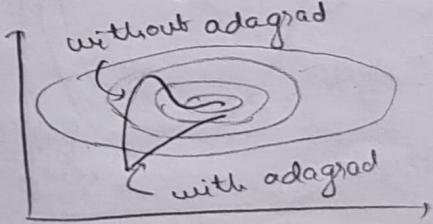
$v_t \uparrow \Rightarrow$ learning rate \downarrow

$\epsilon \rightarrow$ small number

Bg:-

Consider 1000 data point (x, y) pairs with x being 0 80% of the times.

$x \rightarrow$ sparse feature
 $y \rightarrow$ dense feature



without adagrad, initially updates in x -direction were little (sparse). updates occur along y and then slowly along x .

with adagrad, the learning rate for x-direction is large due to small history. Thus, the updates occur towards the optima but due to large no. of updates along y-dir. the history becomes large and hence learning rate becomes very small. Thus we are not able to reach the minima in y-dir.

* Adagrad decays the learning rate very aggressively.

RMS Prop:

Here, we update the history slowly. Thus, the learning rate decay is slow.

exponentially weighted

$$\text{weighted} \rightarrow v_t = \beta v_{t-1} + (1-\beta) (\nabla \theta_t)^2$$

$$0.95 \leq \beta \leq 1$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t + \epsilon}} (\nabla \theta_t)$$

with RMS Prop, we will reach the minima.

Adams Adaptive Moments Estimation

Here, we add momentum to the RMS prop.

$$m_t = \beta_1 m_{t-1} + (1-\beta_1) \nabla \theta_t$$

→ running avg of the gradients (mean) (moment)

$$v_t = \beta_2 v_{t-1} + (1-\beta_2) (\nabla \theta_t)^2$$

update history

* we take the mean of the gradient as it is better to update according to avg. behaviour of gradients.

$$\hat{m}_t = \frac{m_t}{(1-\beta_1)}$$

$$\hat{v}_t = \frac{v_t}{(1-\beta_2)}$$

Bias correction

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} * \hat{m}_t$$

* we want $E[\hat{m}_t]$ to be equal (near) to $E[\nabla \theta_t]$. For that we do the bias correction.

Lec 8.1: Bias and Variance

Lecture 8.1: Bias and Variance

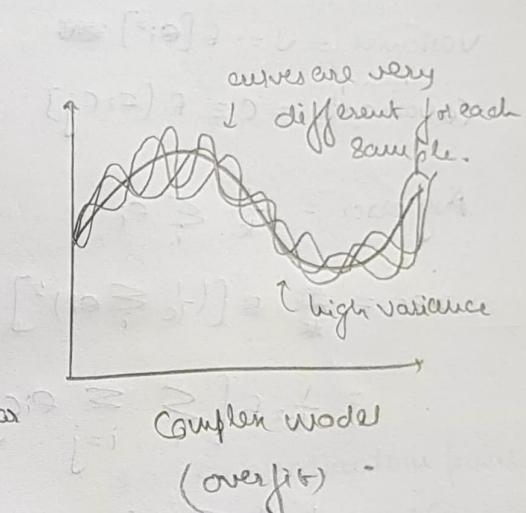
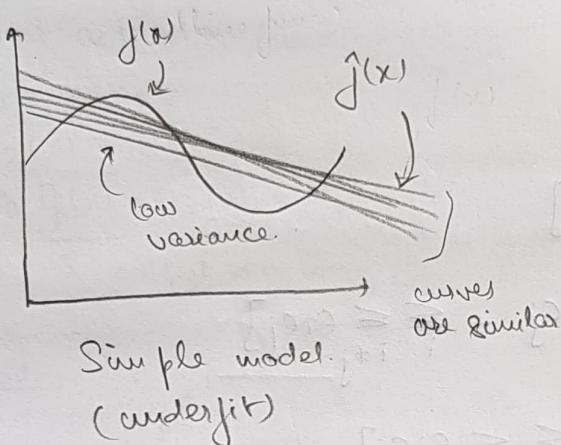
Consider 25 points randomly sampled from a sinusoidal curve is used to train a model.

We have two models:

Degree 1 polynomial

Degree 25 Polynomial.

If we sample 10 times and train we will obtain different curves each time.

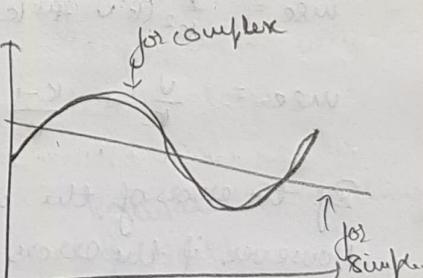


$$\text{Bias} = E[\hat{f}(x)] - f(x)$$

expected value
from approximated
func.

* High for simple model.
low for complex "

$$\text{Variance} = E[\hat{f}(x) - E(\hat{f}(x))]^2$$



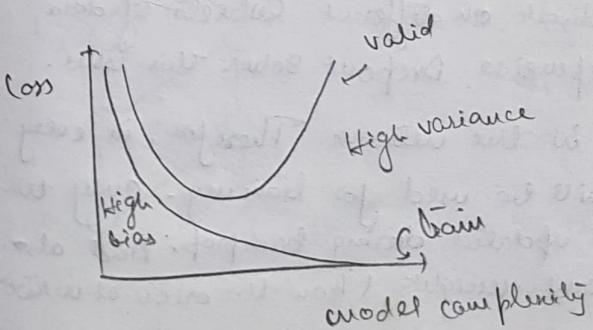
Lecture 8.2: Train loss vs Valid. loss

$$\text{Loss} = E[(y - \hat{f}(x))^2] = \text{Bias}^2 + \text{Variance} + \text{some error}$$

Therefore, for minimum loss both bias and variance should be low.

Lec 8.6: Parameter Sharing & Tying:

* It's used in CNNs



Ensemble Methods:

Combine output of different models to reduce generalization error.

Bagging: ensemble of different instances of some classifier trained on different subsets of training data.

Why ensemble works:

Let's have k logistic regression models.

Each model makes error ϵ_i on the test example.

$$\text{Variance} = V = E[\epsilon_i^2]$$

$$\text{Covariance} = C = E[\epsilon_i \epsilon_j]$$

$$\text{Avg error} = \frac{1}{k} \sum_i \epsilon_i$$

$$\text{MSE} = \frac{1}{k} E\left[\left(\frac{1}{k} \sum_i \epsilon_i\right)^2\right]$$

$$= \frac{1}{k^2} E\left[\sum_i \sum_{i=j} \epsilon_i \epsilon_j + \sum_i \sum_{i \neq j} \epsilon_i \epsilon_j\right]$$

$$= \frac{1}{k^2} \sum_i E[\epsilon_i^2] + \sum_i \sum_{i \neq j} E[\epsilon_i \epsilon_j]$$

$$= \frac{1}{k^2} \sum_i E[\epsilon_i^2] + \sum_i \sum_{i \neq j} E[\epsilon_i \epsilon_j]$$

$$\text{MSE} = \frac{1}{k^2} (kV + k(k-1)C)$$

$$\text{MSE} = \frac{V}{k} + \frac{k-1}{k} C$$

If the errors of the model are perfectly correlated $C=V \Rightarrow \text{MSE}=V$

However, if the errors ϵ_i are uncorrelated $C=0 \Rightarrow \text{MSE}=\frac{V}{k}$

This error decreases with ensemble as different models trained on different subsets of data will have uncorrelated errors.

Dropouts

* Training several instances of network on different subsets of data is expensive. Even testing is expensive. Dropout solves this issue.

Dropouts drop nodes randomly in the network. Therefore, in every iteration, different networks will be used for training. Only the weights for the active nodes are updated during backprop. And also, different networks share the same weights (from the original network).

So, even if ~~a particular~~ network gets sampled only a few times, its weights get updated by the other networks (weight sharing).

Since, each node was present with a prob. p , we scale-the weights by p during testing (we only have p confidence).

Dropout also prevents the neurons from depending on other neurons.

The neighbouring neurons in a layer can be dropped, so the active neurons must learn to improvise. The neurons are prevented from co-adapt.

Lee 9.3: Better Activation Func.:

without activation func.:

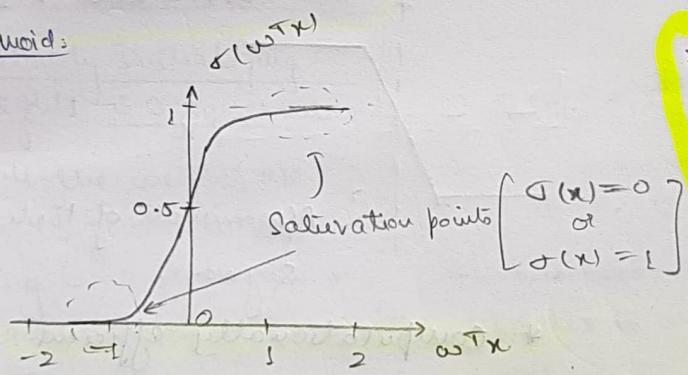
$$h_1 = w_1 x + b_1$$

$$h_2(x) = w_2 h_1 + b_2$$

$$\hat{y}(x) = h_3(x) = w_3 h_2 + b_3$$

∴ output is a linear func. of x .

Sigmoid:



When $\nabla \theta \rightarrow 0$

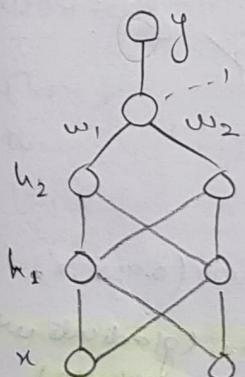
$$\theta_{t+1} = \theta_t - 2 \nabla \theta$$

⇒ no weight update.

Day 2

* At the saturation points, the gradient of $\sigma(x)$ will be close to zero (vanishing gradients). So, if in any layer if the weights become large then $w^T x$ becomes large. ∴ $\sigma(w^T x)$ can be saturated (vanishing gradient).

* Sigmoid is also not zero centred.



$$a_3 = w_2 h_{21} + w_2 h_{22} \quad (\text{common})$$

$$\nabla w_1 = \begin{bmatrix} \frac{\partial L(w)}{\partial y} & \frac{\partial y}{\partial h_{21}} & \frac{\partial h_{21}}{\partial a_3} & \frac{\partial a_3}{\partial w_1} \\ \frac{\partial L(w)}{\partial y} & \frac{\partial y}{\partial h_{22}} & \frac{\partial h_{22}}{\partial a_3} & \frac{\partial a_3}{\partial w_2} \end{bmatrix}$$

w_2 both the

output of $\sigma(x)$ is between [0, 1]

w_2

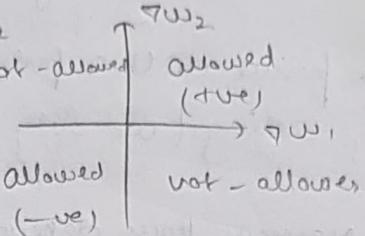
So, If this is true

$$\Rightarrow \nabla w_1, \nabla w_2 \rightarrow \text{true}$$

$$\text{If } -\text{ve} = \nabla w_1, \nabla w_2 \rightarrow -\text{ve}.$$

Then all the gradients of a layer are either +ve or -ve. Restricts update direction (slow convergence).

Eg:- for 2nd order

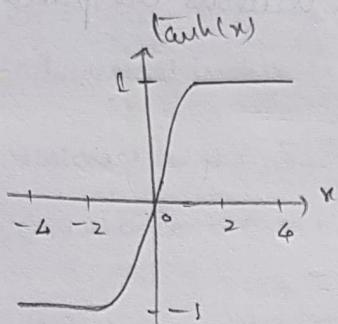


* Sigmoid is computation.

- Only expensive due to $\exp(x)$.

$\exp(x)$.

Tanh:



* zero centred $[-1, 1]$

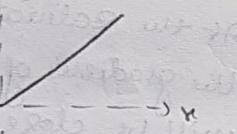
* Vanishing gradient at saturation.

* Computationally expensive: $\exp(x)$

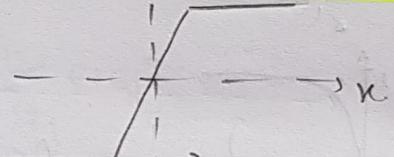
ReLU:

How, is it non-linear?

$$f(x) = \max(0, x)$$



$$\text{ReLU}(x+1) - \text{ReLU}(x-1)$$



Approximate of tanh or sigmoid.

* Computationally efficient

* Does not saturate in the region.

* faster convergence

* can lead to dead neurons.

↳ derivative is 1 for the

* If a large gradient (due to large learning rate) causes the bias to become a large -ve value.

$$\text{ReLU}(w_1x_1 + w_2x_2 + b)$$

↳ normalized

Thus input can be -ve.

This is called dead neuron. (former)

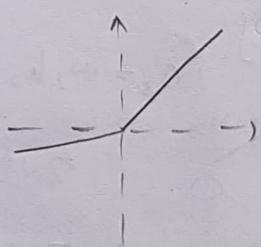
After this, w_1 and w_2 (weights)

will not get updated (0 gradient), thus the neuron will remain dead.

* If the learning rate is high.

$$\frac{\partial \text{ReLU}(x)}{\partial x} = \begin{cases} 0 & ; x < 0 \\ 1 & ; x > 0 \end{cases}$$

Leaky ReLU:

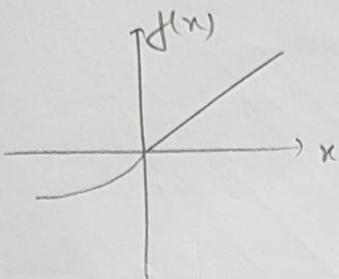


$$f(x) = \max(0.01x, x)$$

* Not satn. (gradients will not die)

* close to zero centred.

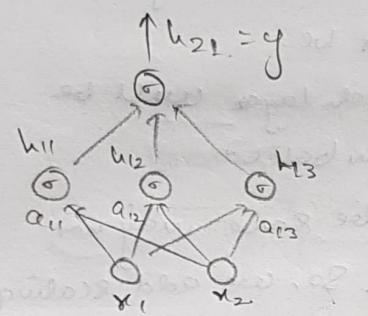
* ReLU: Exponential Linear Unit:



- * Benefits of ReLU
- * Small gradient in -ve side.
- * Expensive: $\exp(x)$
- * Close to zero output

Lec 9.4: Better Initialization Strategies.

Initializing weights to zero:



$$\text{Here: } h_{21} = w_{11}x_1 + w_{12}x_2 + a_{11}$$

$$a_{11} = w_{21}x_1 + w_{22}x_2 + a_{12}$$

If all weights are zero

$$a_{11} = a_{12} = 0 = a_{13}$$

$$\Rightarrow h_{11} = h_{12} = h_{13}$$

$$\text{Now: } \nabla w_{11} = \frac{\partial L(w)}{\partial y} \times \frac{\partial y}{\partial h_{11}} \times \frac{\partial h_{11}}{\partial a_{11}} \cdot x_1$$

$$\nabla w_{12} = \frac{\partial L(w)}{\partial y} \times \frac{\partial y}{\partial h_{12}} \times \frac{\partial h_{12}}{\partial a_{12}} \cdot x_2$$

$$\text{Hence, } \nabla w_{11} = \nabla w_{12} = \nabla w_{13}$$

In the next iteration, ∇ has the updates to all the neurons in a layer will be same because the gradients will be same.

Thus the weight update is restricted.

In general, initializing weights to the same value, we will have symmetry breaking problem.

Randomly Initializing weights to a small value.

$w^T x \rightarrow \text{large} \leftarrow$
If each $w_i \approx 1, \approx w_i x_i$

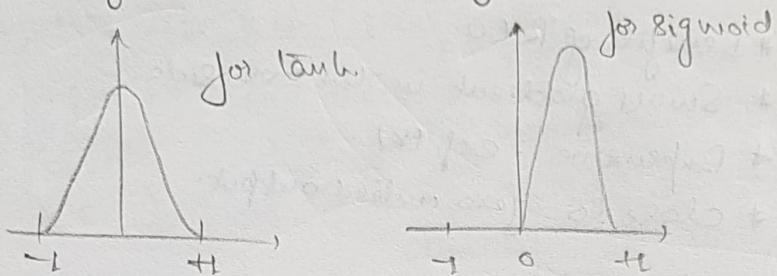
Neuron will saturate
(vanishing gradients)

* Therefore draw the weights from a unit Gaussian and scale them by $1/\sqrt{n}$. This way the variance in the output of a layer does not blow up or shrink.

! Here, $n \rightarrow$ no. of neurons in the given layer.

* weights depend on the no. of neurons in the layer. Since, the weights are divided by n , they can't blow up.

* weights should have a good distribution. (with sufficient variance)



Lec 9.5: Batch Normalization:

In mini-batch GD, the distribution of the output of each layer keeps changing across mini-batches. Training becomes difficult.

Even if the ~~original~~ inputs to the first layer is unit gaussian, the input to subsequent layers may not be.

So, we normalize - the outputs of each layer will be made unit gaussian over the current mini-batch.

But, forcing the distribution to be same may not always be good. (Same variance & bias). So, we add scaling and shifting operation (learnable parameters) for the network to decide what is best.

$$y_{i0} = \gamma_0 z_{i0} + \beta_0$$

(Learnable γ)

Lec 11.1: The conv. operations:

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \rightarrow \text{Blur (averaging)} \quad \begin{bmatrix} 0 & 1 & 0 \\ -1 & 5 & -1 \\ 0 & 1 & 0 \end{bmatrix} \rightarrow \text{Sharpen}$$

(Subtract the neighbours and emphasise the centre pixel)

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{bmatrix} \rightarrow \text{edge detector}$$

(evaluates to zero when neighbour pixels are same as centre pixel)

- * For 3d image, kernel will also be 3D and usually has the same depth.
- Then, the output of conv. will be a 2D image (feature map).
- * we can have k filters for k different feature maps.

Lec 11.2: Input and output size relation:

- * For a filter size F , and stride = 1, output dim. decreases by $(F-1)$

$$\Rightarrow W_{\text{out}} = W_{\text{in}} - (F-1)$$

- * If we want the same output dim., we pad input.

$$W_{\text{out}} = W_{\text{in}} - F + 2P$$

- * For stride S ,

$$W_{\text{out}} = \frac{W_{\text{in}} - F + 2P}{S} + 1$$

Lec 11.3: CNN

SIFT, HOG were dominant at that time.

- * Earlier people used convolution to find edges etc. in an image, flatten them manually and feed them to a DNN for image classifications. The filters were hand-crafted.

Instead, we can learn these filters and have multiple layers of these filters. This is CNN.

- * CNN have sparse connection, weight sharing (as the same kernel is used for the entire image)

* Since the kernel size is small, overfitting is a problem. So, we have multiple kernels.

* Weight sharing results in much less no. of parameters.

* Maxpooling has no parameters. It is just a max operation.

* In CNNs, Dense connections (FC layers) have large no. of parameters.

Eg:- Significance of weight sharing:

$$\text{Input} = 32 \times 32$$

$$\text{Output} = 28 \times 28$$

$$\text{FC-layer: } 22^2 \times 28^2$$

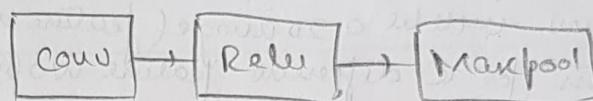
$$\text{CNN: } 5 \times 5 \times 6 = 150$$

$$\begin{bmatrix} P=5 \\ K=6 \\ S=1 \end{bmatrix}$$

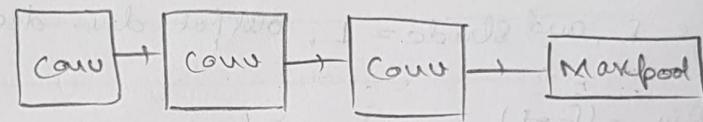
Ques: when we count the no. of layers in a Network, we only count those that have parameters.

Thus, Maxpool is not counted.

Ques: Non-linearity should be done after convolution and before Maxpool.



* we can have back to back conv layer and then a maxpool (VGG style)



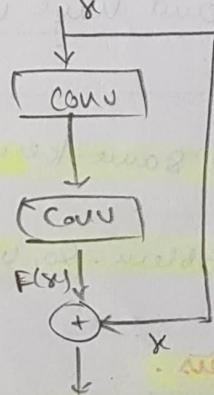
Googlenet:

* uses inception modules.

Inception modules first use 1x1 conv to reduce the no. of channels (depth), then they compute different convolutions with different filter size in parallel and concat the outputs. This lets the network to utilize the good of different filter sizes.

Resnet:

A shallow network works better than a deep network because it is very difficult to learn identity mapping by conv. net. So, it is better to use skip connections. This way a superset of a good working network will perform at least as good as the shallow network.



$$H(x) = F(x) + x$$

Eg:- Significance of weight sharing:

$$\text{Input} = 32 \times 32$$

$$\text{Output} = 28 \times 28$$

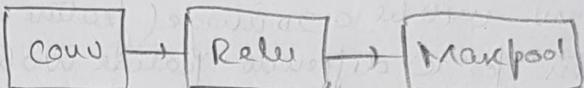
$$\text{FC-layer: } 22^2 \times 28^2$$

$$\text{CNN: } 5 \times 5 \times 6 = 150$$

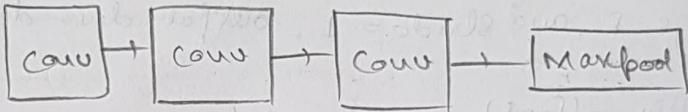
$$\begin{cases} P=5 \\ K=6 \\ S=1 \end{cases}$$

Rule: when we count the no. of layers in a Network, we only count those that have parameters.
Thus, Maxpool is not counted.

Rule: Non-linearity should be done after convolution and before Maxpool.



* we can have back to back Conv layer and then a Maxpool (VGG style)



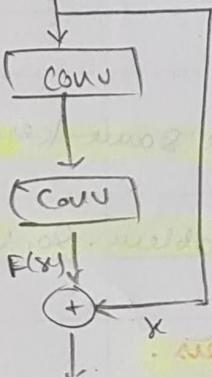
Googlenet:

* uses inception modules.

Inception modules first use 1×1 conv to reduce the no. of channels (depth), then they compute different convolutions with different filter size in parallel and concat the outputs. This lets the network to utilize the good of different filter sizes.

Resnet:

A shallow network works better than a deep network because it is very difficult to learn identity mapping by conv. nets. Soln. is to use skip connections. This way a superset of a good working network will perform at least as good as the shallow network.



$$H(x) = F(x) + x$$

Lec 12-3: Occlusion Expt:

It can be used to check the robustness of the network. The idea is to randomly occlude (gray out) random portions of the image and check for prediction.

For eg:- a good face detector should work even if we cover the nose. (Should not depend too much on nose)

* The network should not rely heavily on a single part of image.

Lec 12-5: Guided Backprop:

This is used to find the influence of input image on any neuron in the network. (Basically finding out the specific portions like eyes, ears etc that the neuron responds to).

We first do a forward pass through a network. We then take a neuron



Selecting this neuron.

The other neurons in this layer are set to zero.

This method is called guided backprop. This gives a more clear influence of input image on the neuron.

Now, backpropagate to get the reconstructed image.

While backpropagating, use ReLU just like we did during forward pass. This will stop the -ve gradients.

Lec 12-6: Optimization over images

Suppose we have a trained network and we want to generate an image that when passed through the network maximizes the output to a certain class.

We can solve this as an optimization problem. Assume the input image to be a trainable parameter that we will update to minimize loss.

- * Start with zero image
- * Set the target (label) to $[0 \ 0 \ \dots \ 1 \ \dots \ 0]$
- * Back propagate and update the pixel values. $\delta = \delta - \eta \nabla \delta$
- * Do a forward pass and calculate cross entropy loss.
- * Repeat till convergence.

one hot encoded

* These generated images are nowhere near to the actual image of the particular class but can be used to fool the network.

Lec 12.7: Create images from embeddings

Take an input image and do a forward pass through the network. Now, the output of each of the layers in this network can be used as embeddings to generate the original image (reconstruction).

To regenerate, put a blank image at the input and the loss will be MSE between the embeddings and output of the chosen layer. Now, use GD to update pixel values of the input image. As we go deeper into the network, the reconstructed image starts to look blurred and weird.

Lec 12.8: Deep Dream:

Basically, we start with some input image (eg:- sky) and we want to maximize the firing of a certain neuron in some layer. So, objective function = $(h_{ij})^2$

to be maximized output of the neuron.

Now, if we update the image to maximize the objective func, the network will draw ~~with~~ certain objects in the input image in correlation to the input image.

Bg:- for sky, birds and castles may be formed and for grasses, certain ~~animal~~ animals may be formed.

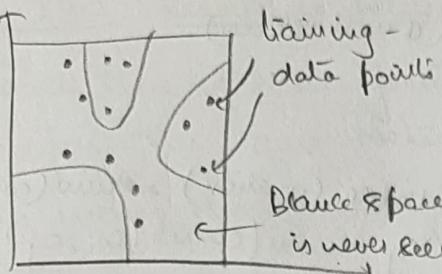
Basically, when a network sees an input image (eg sky), it expects to see birds or castles. If they aren't present, it draws them.

Lec 12.10: Fooling Deep conv nets:

To fool a network, we either take a blank image or an image of one of the classes eg. bus and then update this image slightly using back prop to fit some other class eg. cat.

This image will still look like the bus, but the network will predict it as cat. We can start from blank image and add the updated img. to bus image to obtain the same effect.

Why this happens is explained by (carpacity). Images have very high dimensionality eg (224×224) and in any data set we only see a very small subset of this space. Using the training data, we fit some decision boundary but we leave out the other possible images. We generalize the output (decision) for a large no. of points that we haven't seen.



Blue & face
is never seen by
the network
but we still decided
which class it belongs to.

Since only 1-bit is hot

Lecture 10.1: One-hot representation of words (sparse representation)

Given a corpus (collection of sentences), the set of all the possible words in it is called the vocabulary ' V '. Each word in the vocabulary can be represented by one-hot encoded vectors of length V .

- Disadvantage: One-hot encoding of words does not capture similarities between words.

Eg:- $\text{euclid_dist}(\text{cat}, \text{dog}) = 52$. } same for all pairs.
 $\text{euclid_dist}(\text{dog}, \text{beech}) = 52$

$\text{cosine_sim}(\text{cat}, \text{dog}) = 0$. } orthogonal pairs.
 $\text{cosine_sim}(\text{cat}, \text{beech}) = 0$. pairs.

- As V increases, large space is required.

Lec 10.2: Distributed Representations of words:

Co-occurrence Matrix:

Context

	human	machine	system	...	user
human	0	1	0	...	0
machine	1	0	0	...	0
system	0	0	0	...	2
:	:	:	:	:	:
user	0	0	2	...	0

Column = length of vocab.

Rows = length of vocab.

dimn.

[Words x context]

Co-occurrence matrix captures the no. of times a word appears in context of another word. Context here is defined as a window of 6 words around that word.

* A few columns can be dropped as words like {a, for, the, etc} don't carry much info.

Stop words → their count is very high.

Now, each row can be taken as a vector representation of the word. The representation is still sparse (very few non-zero columns).

Some fixable Problems

Stop words are very frequent \Rightarrow count is large (it will skew)

Soln 1: Ignore high count words.

Soln 2: Use a threshold t , every entry in the table (matrix) = $\min(\text{count}_i, t)$
 $\Rightarrow x_{ij} = \min(\text{count}(w_i, c_j), t)$

Soln 3: Instead of count, use PMD.

$$\text{PMD}(w, c) = \log \left[\frac{\text{count}(w, c) * N}{\text{count}(w) * \text{count}(c)} \right] \quad N \rightarrow \text{total no. of words.}$$

If two words occur frequently independently, but rarely together
 $\Rightarrow \text{PMD}$ is low.

PMD will be high for very frequently co-occurring words.

Another variant,

$$\begin{aligned} \text{Positive PMD (PPMD)} &= \text{PMD}(w, c) \quad \text{if } \text{PMD}(w, c) > 0 \\ &= 0 ; \text{ otherwise} \end{aligned}$$

* we are facing bigger problems:

- very high dimensionality (N)
- very sparse
- grows with the size of vocab.

Soln: use dimensionality reduction (SVD) Singular value Decomposition

* PCA only works for Sq. matrix

* SVD is a generalization of PCA.

Lec 13.1 Sequence Learning Problems:

In feed forward NN and CNN, we had inputs of fixed size and dependency on the previous inputs.

When dealing with sequences eg text completion, the ~~inputs~~ inputs may be correlated and of variable size. Each input occurs at a time step.

Sometimes, we produce output only at the final stage. Eg:-

In video analysis, we can only predict after watching the complete video.

Lecture 13.2: Recurrent Neural Networks

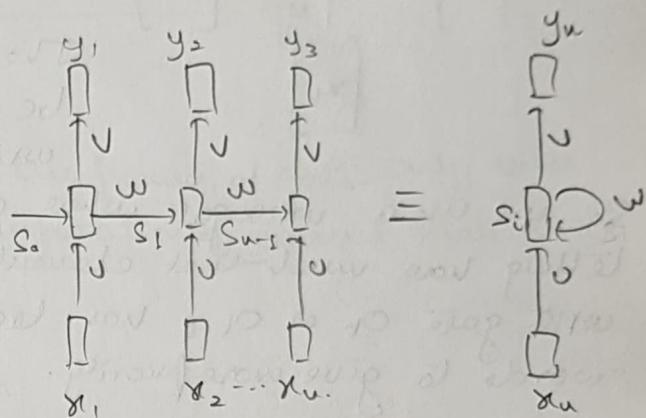
At each time step we feed an input to the NN and it gives an output. Therefore, the same network can be passed each element of the sequence. This way, the length of sequence doesn't affect the network.



$$s_t = \sigma(Ux_t + b)$$

$$y_t = \sigma(Vs_t + c)$$

Add recurrent connection to enable dependency between inputs.



s_t : State of the network at time t
(encodes all the inputs information so far)

(Through s_{t-1} , y_t depends on all the prev. inputs)

$$s_t = \sigma(Ux_t + Ws_{t-1} + b)$$

$$y_t = \sigma(Vs_t + c)$$

Lecture 13.3: Back prop through time:

Lee 13.4: Exploding and Vanishing gradients:

When back propagating through a ~~too~~ long sequence, a small change in the gradient at each time step (> 0 or less than 1) can cause the gradient to explode or vanish.

To avoid this, we backpropagate only till n time steps. (Truncated backprop)
To avoid exploding gradients, we can clip the magnitude.

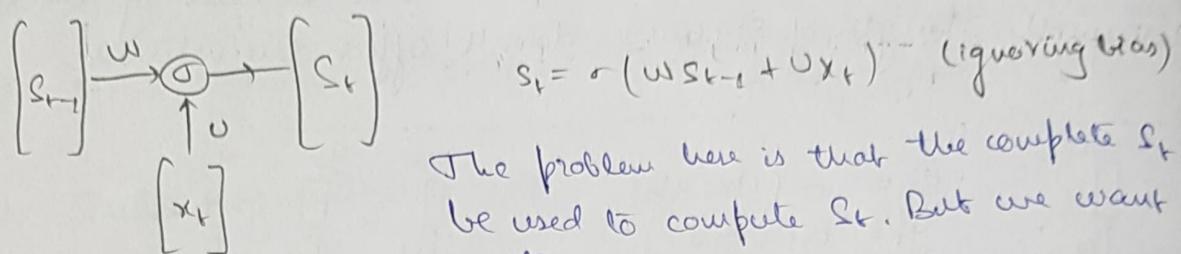
Lee 14.1: Finite state size

The state s_t of an RNN records information from all prev. time steps. Due to fixed ~~size~~ size of s_t , old information gets morphed. So, old information becomes impossible to extract. Thus the morphing should be done in a controlled way.

Lee 14.2: LSTM and GRU:

Consider the task of sentiment analysis of a movie review. ~~The~~ The network (RNN) will read the review from left to right. By the time we reach the end of the review, information about the first few words will be lost (finite state size). So, ideally, our goal is to save important words in the state and forget stop words.

At every time step 't', we want to compute s_t , given s_{t-1} and the input at current time step ' x_t '.



The problem here is that the complete s_{t-1} will be used to compute s_t . But we want selective write.

So, we use a write gate where each element is a value between 0 and 1, telling how much that element of s_{t-1} should be sent to s_t . The write gate o_t or o_{t-1} has learnable parameters. So, it can learn which words to give more priority. [Sigmoid is used to get values between 0 & 1]

$$\underbrace{\begin{bmatrix} s_{t-1} \end{bmatrix} \odot \begin{bmatrix} o_{t-1} \end{bmatrix}}_{\text{Selective write}} = \begin{bmatrix} h_{t-1} \end{bmatrix}$$

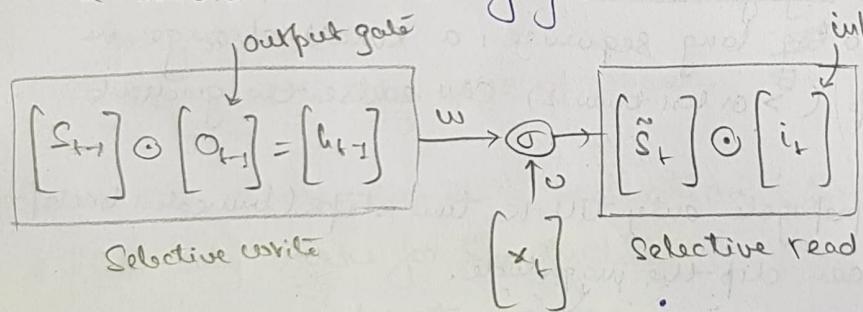
$$o_{t-1} = \sigma(w_{o} h_{t-2} + b_o)$$

Learnable parameters

$$h_{t-1} = o_{t-1} \odot s_{t-1}$$

element wise multiplication

* using output gates, we have created a model which can forget an important words.



input gate

NOW, we calculate \tilde{s}_t which captures complete information from h_{t-1} and x_t . So, we selectively read using $i_t = \sigma(w_i h_{t-1} + b_i)$

Now, we have s_{t-1} and $\tilde{s}_t \odot i_t$, we need to combine them to get s_t . Instead of using s_{t-1} completely, we forget some parts of it (forget gate)

$$\underbrace{\begin{bmatrix} s_{t-1} \end{bmatrix} \odot \begin{bmatrix} o_{t-1} \end{bmatrix} = \begin{bmatrix} h_{t-1} \end{bmatrix}}_{\text{Selective write}} \xrightarrow[w]{\odot} \underbrace{\begin{bmatrix} \tilde{s}_t \end{bmatrix} \odot \begin{bmatrix} i_t \end{bmatrix}}_{\text{Selective read}} + \underbrace{\begin{bmatrix} s_{t-1} \end{bmatrix} \odot \begin{bmatrix} f_t \end{bmatrix}}_{\text{Selective forget}} = \begin{bmatrix} s_t \end{bmatrix}$$

$$f_t = \sigma(w_f h_{t-1} + b_f)$$

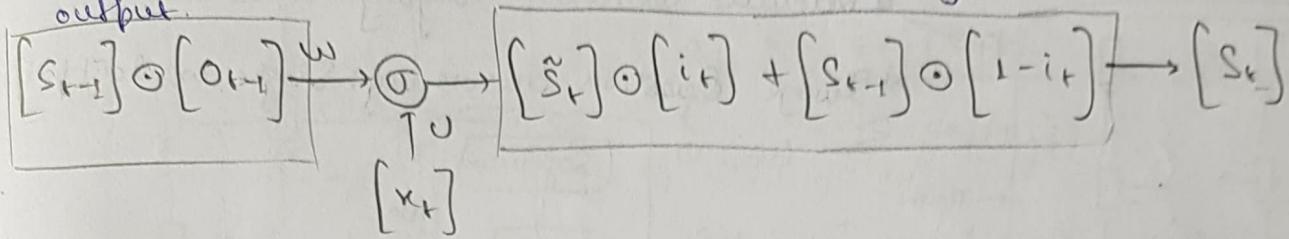
Therefore, in a typical LSTM, we have 2 gates (learnable)

$$s_t = f_t \odot s_{t-1} + \tilde{s}_t \odot i_t$$

* Sometimes s_t is notated as C_t (cell state)

Gated Recurrent units (GRU):

One variant of LSTM is GRU where the input and forget gates are tied together. The idea here is that, we take it from s_{t-1} , then just take $(1-i_t)$ from s_{t-1} as output.



*LSTM prevents vanishing gradients due to the presence of gates. Only those prev. states which contributed to the computation of current state will be used for backprop.

Lec 15.1: Dated to Encoder Decoder Models:

Language Modelling:

Given ' $t-1$ ' words, predict ' t ' word (Text completion) contains information till y_{t-1} (using RNN)

$$y^* = \operatorname{argmax} P(y_t | y_1, y_2, \dots, y_{t-1}) = \operatorname{argmax} P(y_t | s_t)$$

We basically want to find a prob. dist. over the vocab. and find the max.

Compact way of writing:

$$s_t = \text{RNN}(s_{t-1}, x_t)$$

$$s_t = \text{GRU}(s_{t-1}, x_t)$$

$$h_t, s_t = \text{LSTM}(h_{t-1}, s_{t-1}, x_t)$$

Image Captioning:

Here, we are interested to find,

$$P(y_t | y^{t-1}, I) = P(y_t | s_t, f_I(I))$$

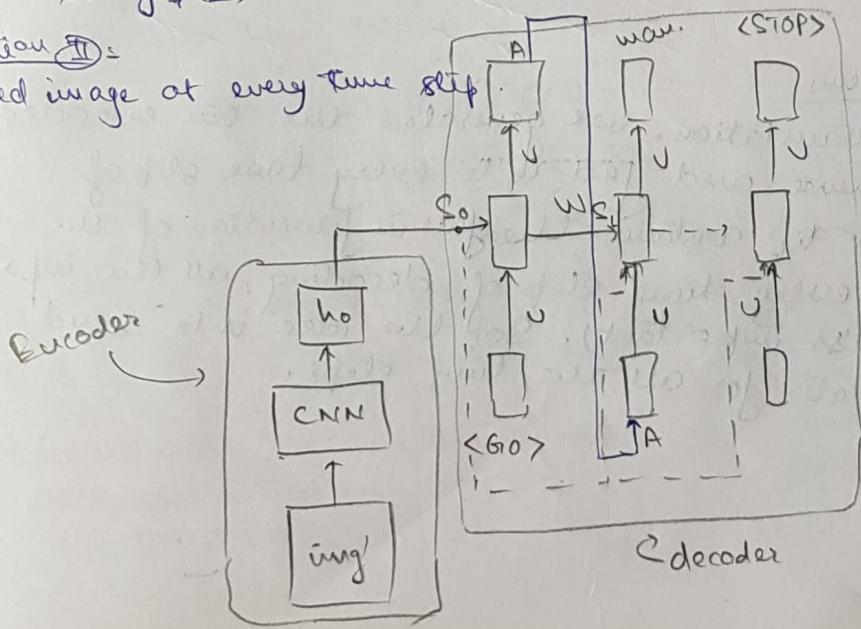
abstract resp. of
image (7th layer)

Option I:

$$\text{Set } s_0 = f_I(I)$$

Option II:

Feed image at every time step



* To train the RNN,

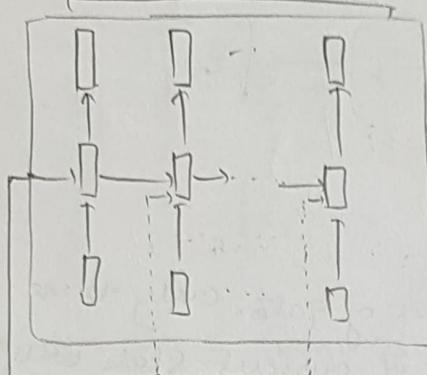
the algorithm used is backpropagation through time (BPTT).

* The loss used is sum of cross-entropy.

Lec 15-2: Applications of Encoder-Decoder Models:

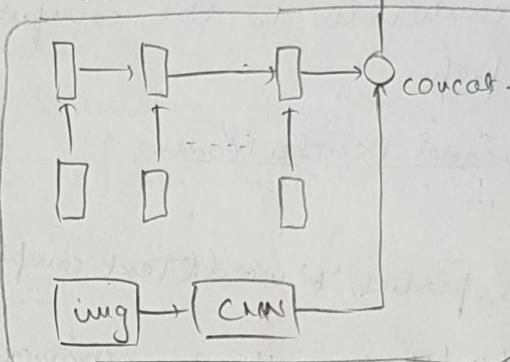
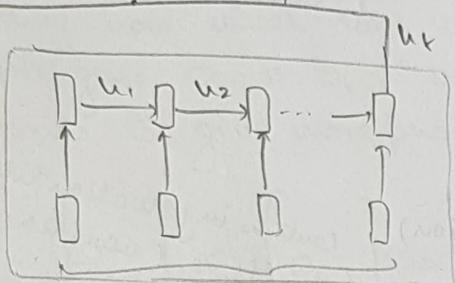
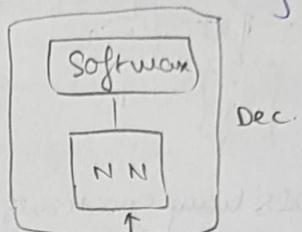
* Textual entailment: Input: Sentence Output: Paraphrasing sentence

Output seq.



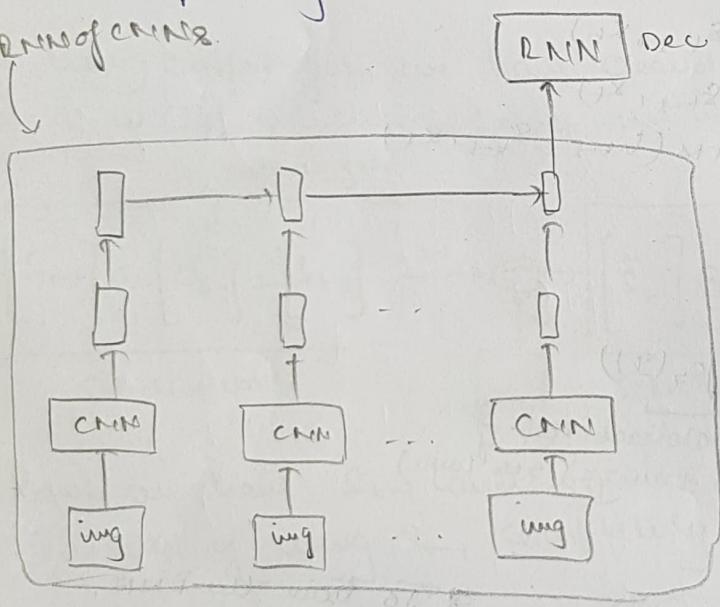
* Machine translation: Language translation

* Image Question answering: Single word output



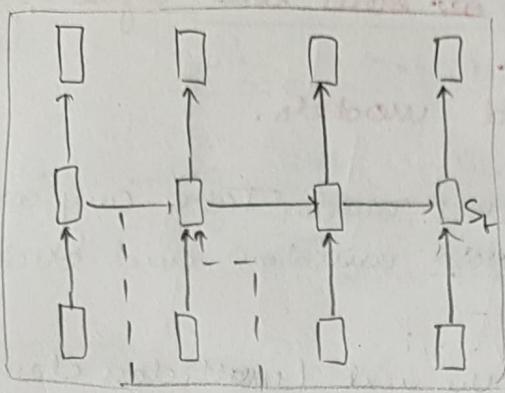
* Video captioning:

enc of CNNs



Lec 15-3: Attention Mechanism:

In the example of machine translation, we generated the ~~the~~ encoded state h_T of the input sequence and fed it to every time step of decoder. This encoded state h_T contains the full information of the input sequence. However, at every time step of decoding, all the input words are not required (less important). So the idea is to send a weighted sum of input states for all the time steps.



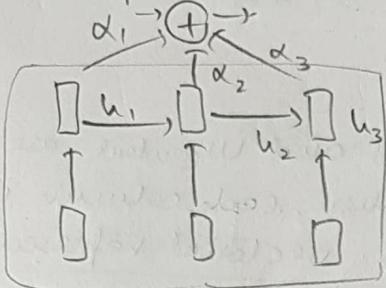
learnable func

$$e_{jt} = f_{ATT}(h_{t-1}, h_t)$$

$$d_{jt} = \frac{\exp(e_{jt})}{\sum_j \exp(e_{jt})} \quad \leftarrow \text{softmax}$$

j^{th} word

$$e_{jt} = \frac{V_{attn}^T \tanh(W_{attn} s_{t-1} + W_{attn} h_t)}{T \text{ scalar}}$$



* When back propagating, these parameters will learn and hence, the attention network will learn too.

Continuation:

10.3: SVD for learning word representations

SVD gives the best rank-\$k\$ approximation of the original rank-\$n\$ matrix. 'k' here is variable, we can decide what rank we want. closer 'k' will be to 'n', more the accuracy.

SVD also brings out latent co-occurrence between the words.

Eg:- user-human has value = 0 in co-occurrence matrix(\$X\$), but since user-install has value = 0 in co-occurrence matrix(\$X\$), but since these two words occur in context with words like "install" etc, in the low rank matrix, their value is non-zero. This hidden co-occurrence becomes visible.

* After doing SVD, the matrix is still of the same size, it is just an approximation (best).

* If we do \$XX^T\$ of the original matrix, each value is a cosine similarity. after SVD this cosine similarity increases (learning better co-occurrences).

From SVD, \$\hat{X} = (U \Sigma V^T) \in \mathbb{R}^{m \times n}\$

$$\text{Now: } X \cdot X^T = (U \Sigma V^T)(U \Sigma V^T)^T$$

$$= U \Sigma \Sigma^T U^T \quad (\because V^T V = I)$$

$$= (U \Sigma)(U \Sigma)^T = W W^T \quad \Rightarrow \boxed{X \cdot X^T = W W^T}$$

$$W = U \Sigma \in \mathbb{R}^{m \times k}_{\text{word}}$$

* Remember, \$k\$ is very small compared to \$m \times n\$. \$k \approx 100\$ or 200.

Thus, cosine similarity is maintained and \$W_{\text{word}}\$ is the representation of \$n\$-words in the vocab.

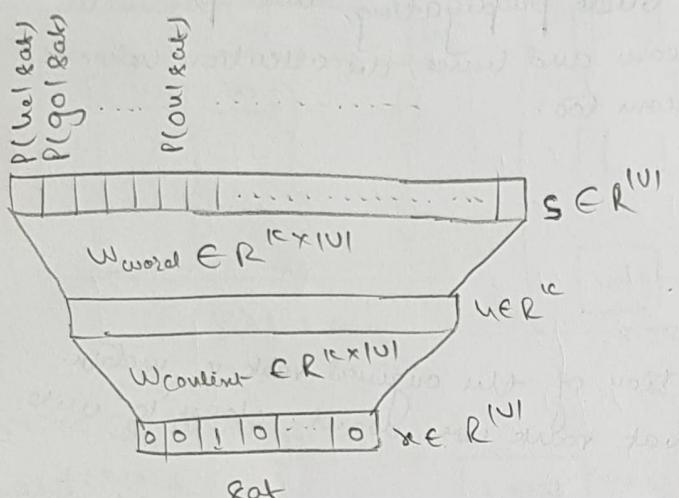
Also, \$W_{\text{content}} = V \in \mathbb{R}^{n \times k}\$

\$\leftarrow\$ representation of context words

Lee 10.4: Continuous Bag of words model:
 * Till now, whatever word representations we had was before deep learning approaches. (count based models).
 Now, we will learn prediction based models.

Consider the task of predicting n th word, given $n-1$ words. Take any written material eg:- stories, blog etc, make a n -word window and slide over it for training.

Let $n=2 \Rightarrow$ given one word predict the next (multiclass classification).



Here, w_{word} and w_{content} are $R^{n \times n}$ matrix where, each column is a n -dimensional vectorial representation of the word. These will be learnt.

$$u = w_{\text{content}} \cdot x$$

$$= \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} a_{12} \\ a_{22} \\ a_{32} \end{bmatrix}$$

Thus, if input is j th word, then j th column of the w_{content} is selected.

Now, to get $p(\text{word}=i|\text{eat})$

$$\begin{aligned} &= \frac{e^{(w_{\text{word}}, u)[i]}}{\sum_j e^{(w_{\text{word}}, u)[j]}} \end{aligned}$$

Thus, $p(\text{word}=i|\text{eat})$ depends

on j th column of w_{content} and

i th column of w_{word} . Thus, i th column

of w_{word} is the vectorial representation

of word 'i'.

Let's see how this will learn good word representations:

Let's denote context word (eat) by c and the collect output word (ou) by index w .

Output func: Softmax

Loss func: Cross Entropy

$$l(w) = -\log \hat{y}_w$$

Maximize w th entry in \hat{y}

$$u_c = u = w_{\text{content}} \cdot x_c$$

The vectorial rep. of context word.

$$\hat{y}_w = \frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})}$$

vectorial rep. of word. (Output)

$$\Rightarrow l(\theta) = -\log \frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})} = -(u_c \cdot v_w - \log \sum_{w' \in V} \exp(u_c \cdot v_{w'}))$$

$$\text{Now, } \nabla_{v_w} = \frac{\partial l(\theta)}{\partial v_w} = -u_c \left(1 - \frac{\exp(u_c \cdot v_w)}{\sum_{w' \in V} \exp(u_c \cdot v_{w'})} \right)$$

$$\nabla_{v_w} = -u_c (1 - \hat{y}_w)$$

$$\text{Update rule, } v_w = v_w - \eta \nabla_{v_w}$$

$$v_w = v_w + \eta u_c (1 - \hat{y}_w)$$

Now, $\hat{y}_w = 1 \Rightarrow$ prediction is correct

$\hat{y}_w = 0 \Rightarrow$ wrong

$$\Rightarrow v_w = v_w + \eta u_c$$

* we are basically adding a fraction of context vector to word vector.

Thus, the word representations will become ~~closed~~ closer to context rep. Thus, cos-sim between v_w and u_c is maximized.

Thus, if we take two words like dog and cat, both of these words will be similar in context word like eats, sleeps etc but a breed will not. This however does not mean that dog and sleep, eat are similar. We only care about the word representations.

Similarly $u=d$,

for this, just concatenate the word 1-hot vector, so, the input will be a vector of dim $|V|$ ~~with~~ with d -elements

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} a_{12} + a_{13} \\ a_{22} + a_{23} \\ a_{32} + a_{33} \end{bmatrix}$$

$\underbrace{\quad}_{\text{sum of columns}}$

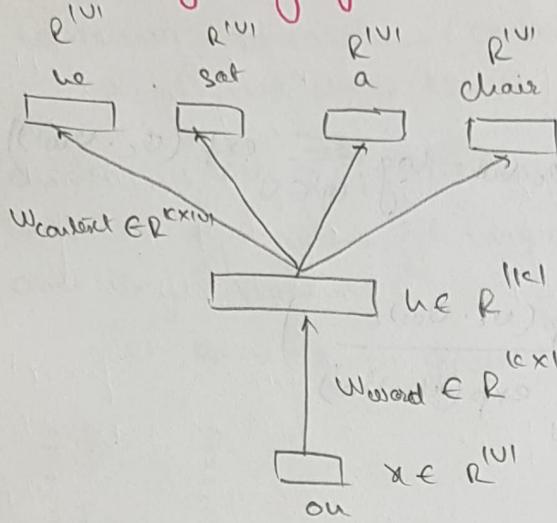
Here, the order of words in input does not matter. (bag of words)

Problem: Softmax over $|V|$ is computationally expensive.

* Bag of words model predicts the output given the context.

Lec 10.5: Skip gram model

Reverse of bag of words: Given a word, predict all the context words.

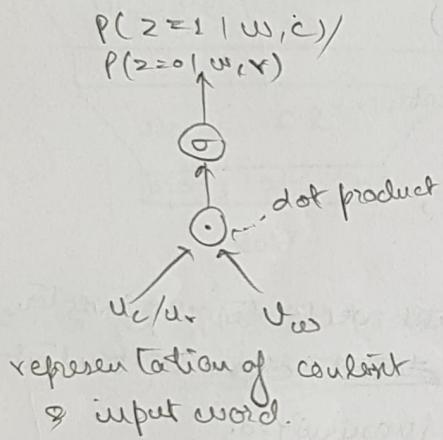


The sum of crossentropy will be the loss.

It has the same problem as bag of words
(computationally expensive softmax)

- Sols: ① Negative Sampling
- ② Contrastive estimation
- ③ Hierarchical Softmax

Negative Sampling:



Make a set D of correct (w, c) pair and incorrect (w, r) pairs D' . $r \rightarrow$ random word.
The task now is to maximize the sigmoid operation of the dot product between input and context words. At the same time we need to make the output ≈ 0 for incorrect pairs.

The objective therefore is to maximize.

$$\underset{\mathcal{D}}{\text{maximize}} \prod_{(w,c) \in D} P(z=1 | w, c) \underset{(w,r) \in D'}{\text{maximize}} P(z=0 | w, r)$$

* we have ' \prod ' because we want to maximize this for all (w, c) and (w, r) pairs.

* The idea here is that we are trying to bring context words close to input words. Thus, words in the same context will come closer. Additionally, we are also making the dot product (cosine similarity) between unrelated words very small. Thus unrelated words are also pushed far apart. (improvement from bag of words)

* In general size of D' is k times that of D because in natural language incorrect pairs are more common than true pairs.

* Now, instead of sampling r using uniform dist, we can do it based on the freq of occurrence.

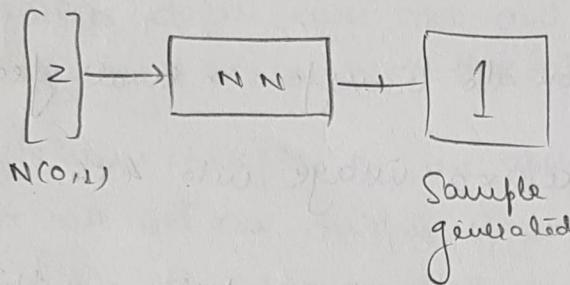
$$p(r) \sim \frac{\text{count}(r)}{N}$$

(3/4) ↴ hyper parameters

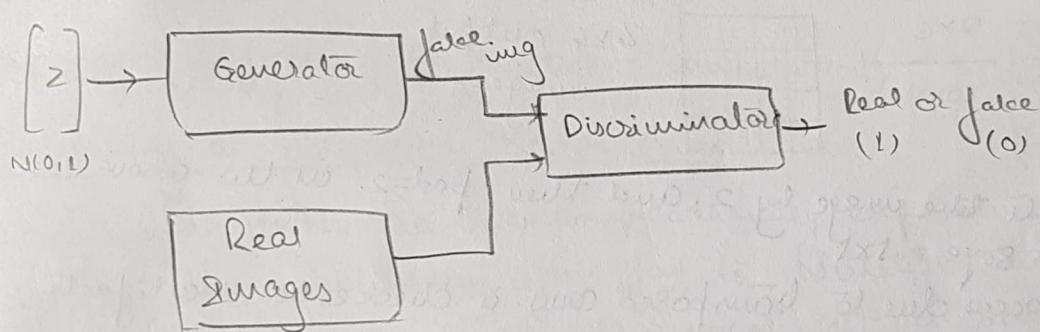
Total no. of words in the corpus.

Lee 22.1: Generative Adversarial Networks (GANs):

Given some dataset, the goal is to generate more data that looks like they have come from the distribution. Data sets belong to a very complex high dimensional dist. It is very complex (difficult) to sample from this very high dimensional dist. So, instead we can sample from a simple normal dist. and learn a transformation from this to the training data dist.



Architecture:



There are two networks, generator and discriminator. Generator's job is to generate better and better images that look like they have come from the same dist. (training dist.). The discriminator's job is to identify whether the images are real or fake (generated). The discriminator outputs 1 for real images and 0 for fake. The generator keeps trying to fool the discriminator.

$$D(x) = 1 \quad \} \text{ for ideal disc.}$$

$$D(G(z)) = 0$$

Discriminator's job is easier, so it will learn quicker and disc. loss will initially be lesser than generator's loss. Then, as generator starts to learn its loss will decrease and $D(G(z))$ will start to increase and $D(x)$ will decrease. This indicates that discriminator is having a hard time identifying (generated images are really good).

Special types of convolution

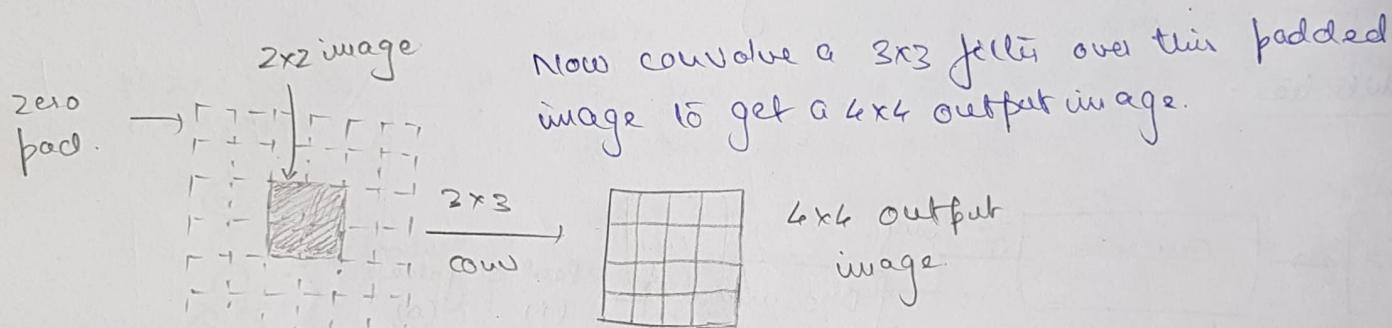
1x1 conv:

- It is used to reduce the depth of feature maps without the expensive 3×3 and 5×5 convs.
- Dimensionality reduction results in a small but dense feature maps (feature pooling)
- 1×1 conv along with ReLU can be used to learn complex functions.

Transposed Convolution:

The CNN uses the interpolation to transform a small feature map (image) into a larger one.

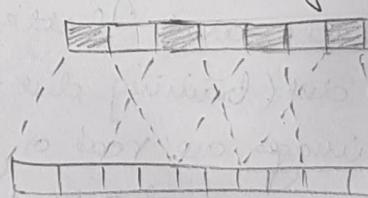
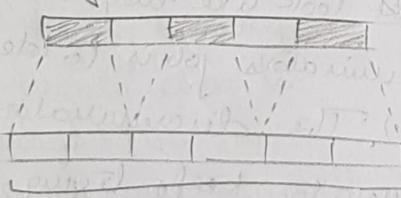
Bg:- We want to convert a 2×2 image into 4×4 .



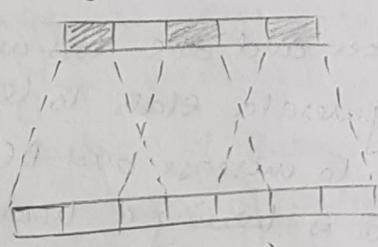
* If we dilate the image by 2, and then pad=2. in the above case.
output size = 5×5 .

A problem that can occur due to transposed conv. is checkerboard artifacts.

Size=2, Stride=2 Stride=2 Size=4.



Size = 3, & stride = 2



Every pixel in the output receives different amount of info. compared to its adjacent pixel.

To prevent, filter size must be divisible by stride.

Dilated convolution (Atrous conv):

They are used to increase the receptive field of the kernel without increasing its size. When $d = 1$, it behaves like a normal conv.

A no. of dilated conv. can be stacked with increasing receptive field. The output of these can be used to aggregate multi-scale contextual information without losing resolution.

Depthwise Separable Conv:

It consists of 2 steps depth wise conv. and 1×1 conv. They are commonly used in deep learning (MobileNet & Xception).

To a volume of depth D , we apply depth wise conv. (filter of depth 1). So, we get an output consisting of depth D . Note that D different kernels convolve on the respective channel of the input.

The output vol. of depth D is converted to depth 1 using a 1×1 conv.

Advantages: Efficiency (a fraction of computation)

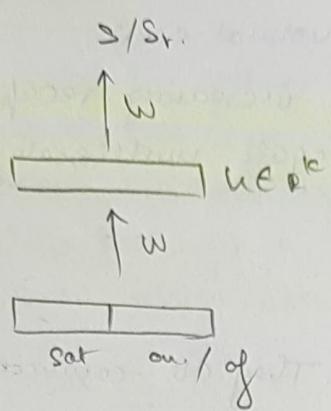
Deep learning: ai Sequence Models Course:

* character level RNN's are more computationally expensive because of the long sequences and they often fail to utilize the correlation due to longer sequences. However, the vocab. size is small and there are no `<unk>` tokens needed.

* language modelling is basically learning the probability of a sequence of words. To determine which sentence is more proper.

* RNN's suffer from vanishing and exploding gradients. Exploding gradients can simply be solved using gradient clipping. Vanishing gradients are harder to solve.

Skip-gram model (contrastive estimation)



Here, we use a NN where we pass input and context word pair of or incorrect pair. The network will output a score. (single value)

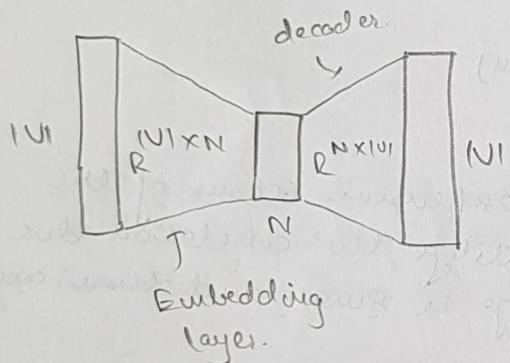
We want the score for correct pair should be at least ' m ' (margin) greater than S_r .

Thus we want to maximize the ...

$$\text{maximize } [\max(0, S - (S_r + m))]$$

Word2Vec (Skip gram model)

In Skip gram model we, we use central words in the corpus to find the context word and this way we learn the embeddings. we use correct pairs only to train the ~~the~~ network.



We input the words as 1 hot vector into the embedding layer. The output to the hidden layer is just the i th row of embedding weights if i th index is hot.

The task ~~of~~ of the decoder is to map the embedded vector to context word. Thus words that are similar in context come close to each other.

Beam Search

When an RNN generates a sequence of words or characters, the output of every time step may be optimal considering the part of the seq. that it has already written but may not be optimal considering the full sentence.

So, instead of using greedy search at each step, we use beam width eg:- 3. Let's say we have a vocabulary of 1000 words. To choose first word, we select top 3 words.

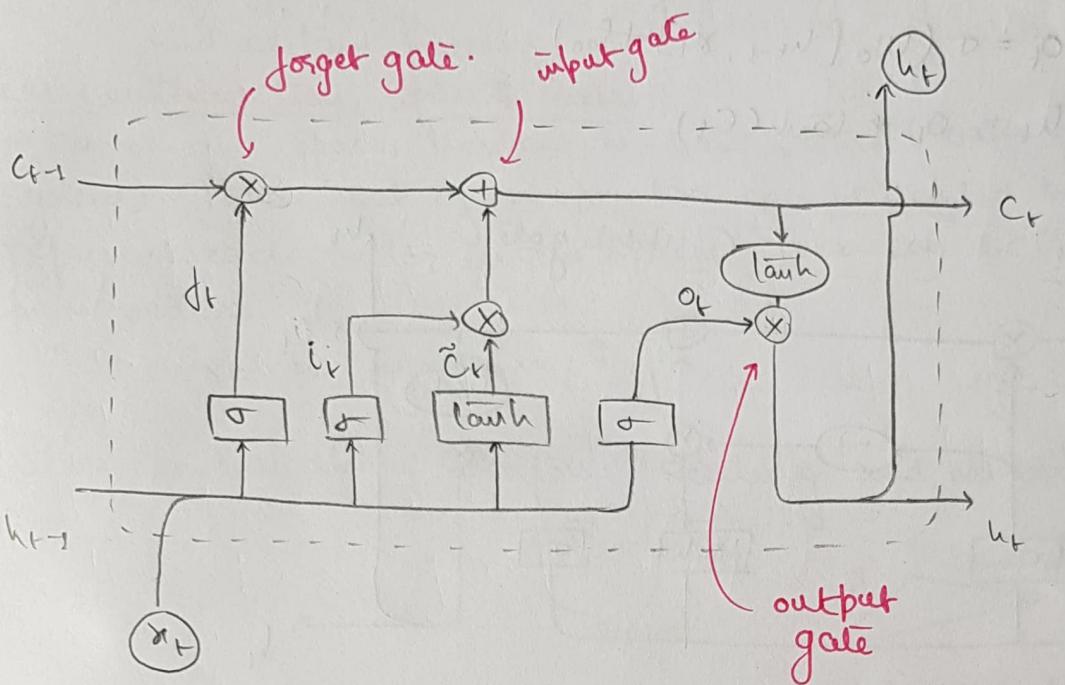
Now, input each of these to the next step, for each we get 1000 outputs (softmax), so total 3000.

Now, choose top 3. We now have top 3 most probable pairs (1st and 2nd word).

Now, input 2nd word into the network, again choose top 3 from 3000 outputs.

This is the top 3 combination of 1st, 2nd and 3rd word.
Keep repeating.

Christopher Olah's Blog (LSTM):



[] NN layers
 ○ point wise operation
 → concatenation
 ↗ copy

The key to LSTM is the cell state. It runs through the entire chain (time step). Each LSTM cell along the time step has the option to selectively add or remove information from this cell state. This selective operation is done through gates.

First, it decides what information to throw from the previous cell state (forget gate).

$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

Next, we decide what new information we want to add. For this sigmoid layer (NN) generates attention map. Then, a tanh layer (NN) makes some useful transformations before adding. This new representation is then weighted by the input gate and added to the cell state.

$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{c}_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

Thus, the total update to \$c_t\$ (cell state):

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

forget information add info.

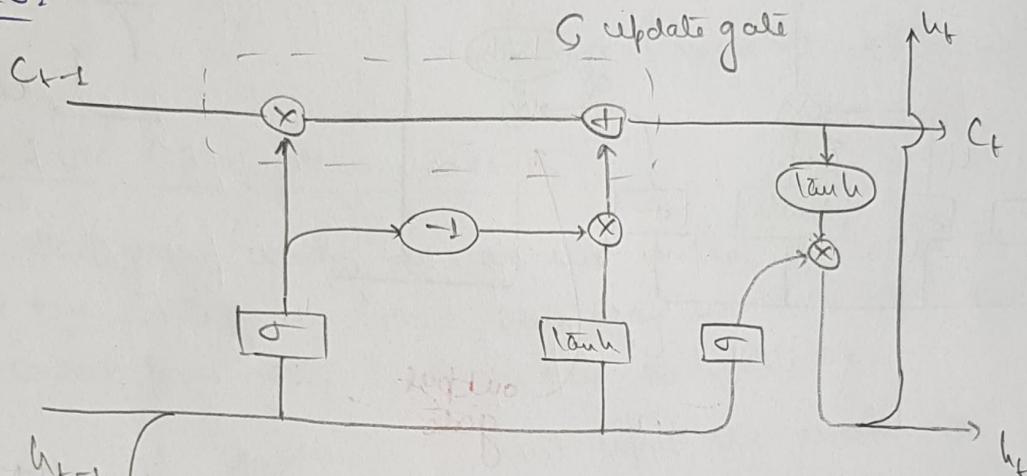
Now, the cell state captures all the essential information about the seq. over all the time steps, but we don't need all the information at each

time step to predict the output. Thus, we filter some information from the cell state using x_t and h_{t-1} to get the hidden state h_t . This h_t is then passed to a NN to output the element for that time step. This filtering is done using output gate.

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

GRU:



In GRUs, the input and forget gates are combined into one update gate. Basically what we allowed to pass in the forget gate, the opposite of that must be added (info. that is forgotten from prev. step should be added in the current step)

$$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$$

These GRUs have only 2 gates (update and output)

`if __name__ == '__main__':
 main()`

This statement is used to check if the current .py file is imported or being run directly.

Bg:- If the file is func.py and we import it as `import func`.

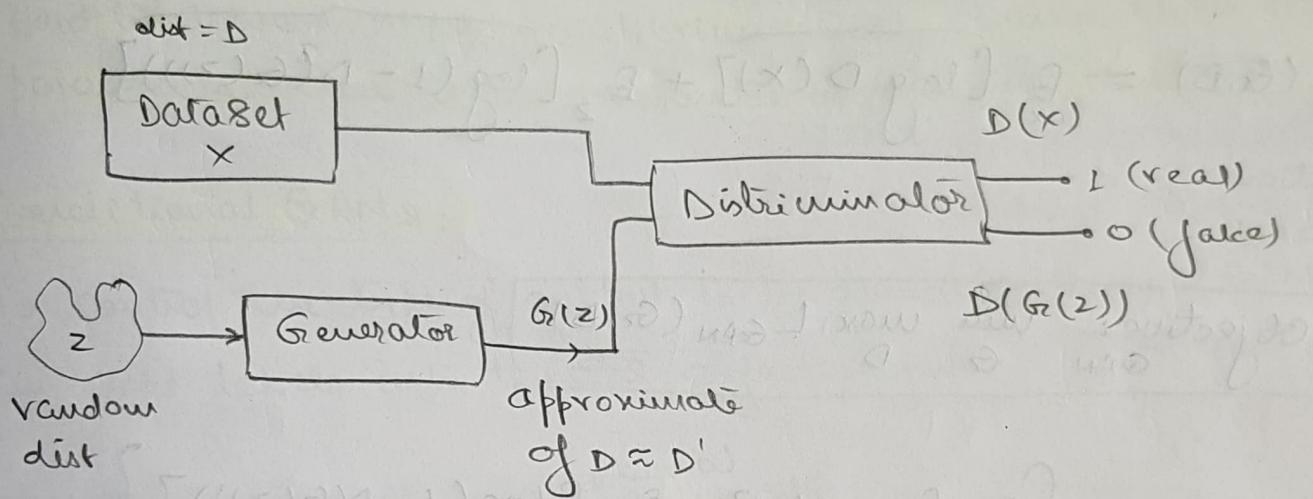
Now, running `print(__name__)` will output 'func'. However, if we are running func.py directly, the output will be '__main__'.

Generative Adversarial Networks:

Objective:

Given, some dataset having distribution D . We don't know this distribution, so we cannot sample new data from this distribution.

However, we can sample data from a random distribution and using this try to ~~to~~ train a generator to approximate D .



Loss Functions:

$$BCE(\hat{y}, y) = [y \log \hat{y} + (1-y) \log (1-\hat{y})]$$

Now,

$$BCE(D(x), 1) = \log D(x) \quad \text{--- (1)}$$

$$BCE(D(G(z)), 0) = \log (1 - D(G(z))) \quad \text{--- (2)}$$

Now, the objective of Discriminator is to classify correctly.

$$\Rightarrow D(x) = 1 \quad \text{and} \quad D(G(z)) = 0$$

Since, $D(x) \in [0, 1]$, when $D(x) = 1$, $BCE(D(x), 1) = \text{minimum}$.

Similarly,

when $D(G(z)) = 0 \Rightarrow BCE(D(G(z)), 0) = \text{maximum}$.

Therefore, the goal of the discriminator is to maximize ① and ②
 Now, let's consider the generator. It's objective is to fool the discriminator.

$$\Rightarrow D(x) = 0 \quad \text{and} \quad D(G(z)) = 1$$

* we don't want the discriminator to misclassify the true examples.

~~when $D(x) = 0$, $\text{BCE}(D(x), 1) \rightarrow \text{minimum}$~~
 when $D(G(z)) = 1$, $\text{BCE}(D(G(z)), 0) \rightarrow \text{minimum}$.

Therefore, the goal of the generator is to minimize ① and ②.

Now,

$$\text{BCE}(D(x), 1)$$

$$\text{BCE}(D(G(z)), 0)$$

$$L_{\text{GAN}}(G, D) = B_x \underbrace{[\log D(x)]}_{\substack{\text{loss function} \\ \text{for the GAN}}} + B_z [\log (1 - D(G(z)))]$$

$$\boxed{\text{Objective: } \min_{\text{GAN}} \max_{G} \max_{D} L_{\text{GAN}}(G, D)}$$

$$\text{Objective}_D = \max_D \left\{ B_x [\log D(x)] + B_z [\log (1 - D(G(z)))] \right\}$$

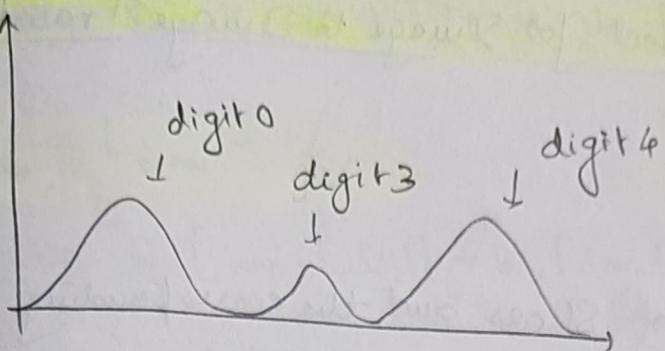
$$\text{Objective}_G = \min_G \left\{ B_z [\log (1 - D(G(z)))] \right\}$$

Limitations of GANs:

* Vanishing Gradient: When starting to train the GAN, the discriminator is trained first. It becomes reasonably good and easily identifies the fake images as false $\Rightarrow D(G(z)) \rightarrow 0$. The derivative of the loss func. becomes small.

* Mode Collapse: The GAN may start to produce identical images after some epochs. Most of the real world data has multimodal dist.

Besides MNIST: digits belonging to each class may be drawn from a normal dist and each digit group may have a different dist.



* The generator is supposed to learn ~~this~~ this multimodal dist, but sometimes it learns only a single mode (mode collapse).

* This happens because the only requirement for the generator is to fool the discriminator. It is much easier for the generator to learn one mode and produce really good samples from that mode.

* Hard to achieve Nash Equilibrium: Both change their own parameters. There is no guarantee of reaching the optima.

Conditional GANs:

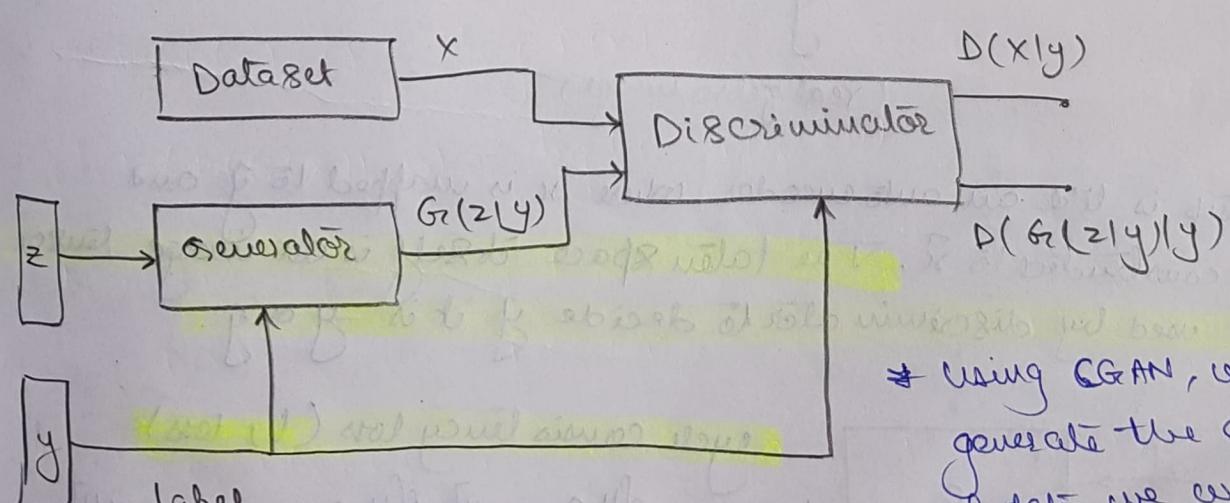
To control the type of output to be generated through GANs, we can add labels into the ~~the~~ objective function.

~~Objective~~

$$L_{CGAN} = E_x [\log D(x|y)] + E_z [\log (1 - D(G(z|y)|y))]$$

conditional probability

$$\text{Objective}_{CGAN} = \min_{G} \max_{D} L_{CGAN}(G, D)$$



* Using CGAN, we can generate the exact type of data we want to generate.

Pix2Pix vs Cycle GAN (used for Image to Image Translation)

used for paired data

for unpaired data

Paired Data: Silhouette of shoes and the corresponding shoes.

Unpaired Data: Collection of Scenery and collection of paintings.

Cycle GAN:

Let's convert horse images into zebra images.

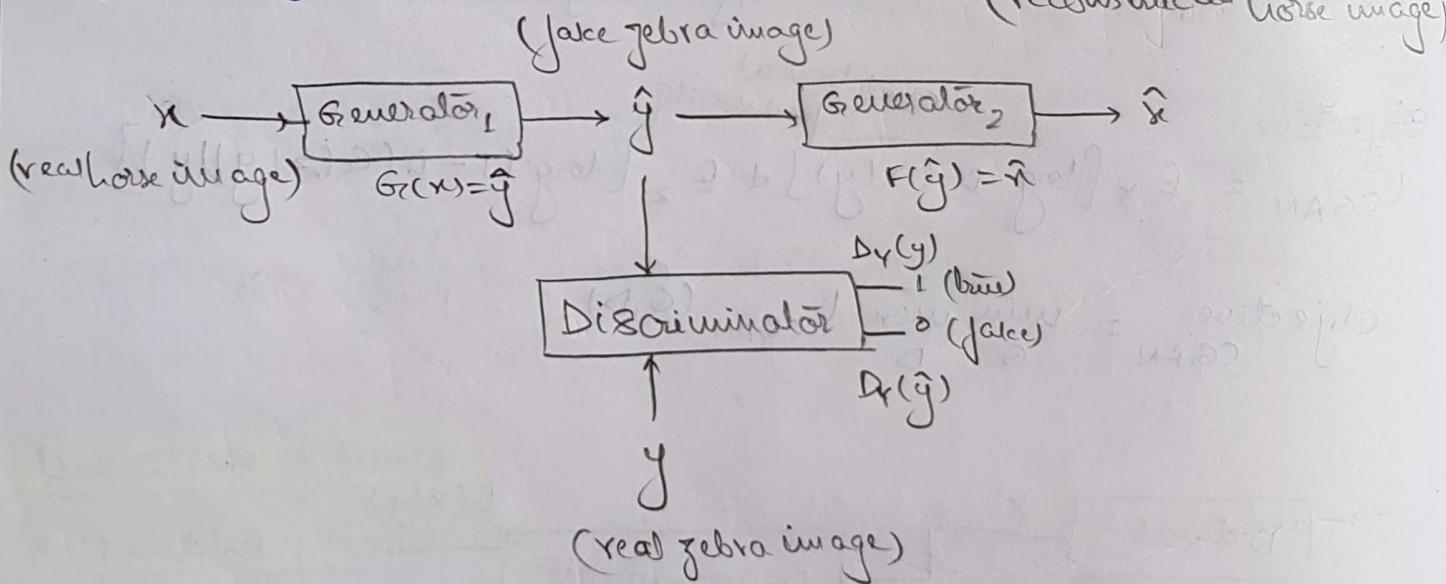
X : Collection of Horse images.

x : A single horse image taken from X .

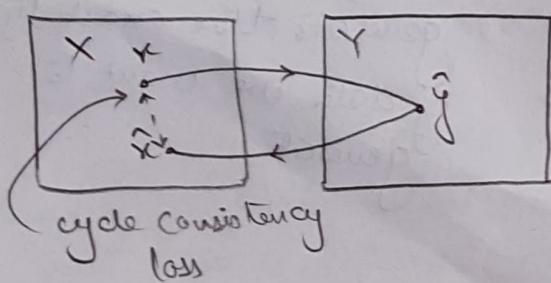
Y : Collection of zebra images

y : A single zebra image taken from Y .

(reconstructed horse image)



* The setup is like an autoencoder where x is mapped to y and then reconstructed to x . The latent space itself is an image tensor which is used by discriminator to decide if it is y or $y-hat$.



cycle consistency loss (L_1 loss)

$$L_{\text{cycle}} = \mathbb{E}_x [\| \hat{x} - x \|] + \mathbb{E}_y [\| F(G_2(y)) - y \|]$$

$$L_{\text{cycle}} = \mathbb{E}_x [\| F(G_2(x)) - x \|]$$

Similarly, $Y \rightarrow X$ generation is also possible.

$$\begin{aligned} L_{\text{cycleGAN}} &= L_{\text{GAN}} + L_{\text{cyc}} \\ &= E_y [\log D_Y(y)] + E_x [\log (1 - D_Y(G(x)))] \cancel{+ E_x [||R(G(x)) - x||]} \\ &\quad + E_x [||R(G(x)) - x||] \end{aligned}$$

DICE GANs

* In DICE GANs, the cycle consistency loss has L_2 Norm.