

solution:-

Forward Pass

(a) h

$$\text{net } h = w_1 x = 0.15 \times 5 = 0.75$$

$$\text{out } h = \frac{1}{1 + e^{-\text{net } h}} = \frac{1}{1 + e^{-0.75}} = 0.679$$

error f^n is calculated
based on predn & actual target

(b) o

$$\text{net } o = w_2 \text{out } h = 0.4 \times 0.679 = 0.272$$

measures how far netw is from making correct prediction

$$\text{out } o = \frac{1}{1 + e^{-\text{net } o}} = 0.568$$

$$E = \frac{1}{2} [y - \hat{y}]^2 = \frac{1}{2} [1 - 0.568]^2 = 0.081 \quad \text{error as } y = 1$$

netw error

- Backpropagation: \rightarrow error is propagated backward through the netw.
Starting from output layer & moving to input layer.

(a)

$$w_2(\text{new}) = w_2(\text{old}) - \eta \frac{\partial E}{\partial w_2}$$

During this, partial deriv. of error wrt wts & bias is calculated using chain rule. (eff. computation of gradient for all parameters in network)

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial \text{out}_o} \times \frac{\partial \text{out}_o}{\partial \text{net}_o} \times \frac{\partial \text{net}_o}{\partial w_2}$$

$$\frac{\partial E}{\partial \text{out}_o} = \frac{\partial [\frac{1}{2} (y - \hat{y})^2]}{\partial \text{out}_o} = \frac{y - \text{out}_o}{\text{out}_o} \times -1 = 1 - 0.568 \times -1 = -0.438$$

$$\frac{\partial \text{out}_o}{\partial \text{net}_o} = \frac{\partial \left(\frac{1}{1 + e^{-\text{net}_o}} \right)}{\partial \text{net}_o}$$

${}^1 \text{ derivative of } f^n \Rightarrow A \cdot F(1 - A) \cdot F'$

$$= \text{out}_o (1 - \text{out}_o) = 0.568 (1 - 0.568) = 0.245$$

$$\frac{\partial \text{net}_o}{\partial w_2} = \frac{\partial w_2 \text{out}_h}{\partial w_2} = \text{out}_h = 0.679$$

$$\frac{\partial E}{\partial w_2} = -0.432 \times 0.245 \times 0.679 = -0.072$$

$$w_2 \text{ new} = 0.4 - \eta \frac{\partial E}{\partial w_2}$$

$$= 0.4 - 0.5(-0.072)$$

$$= 0.436$$

Once gradients are calculated, an optimization algo was then to update

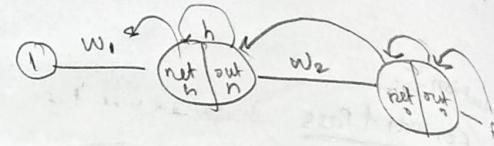
algo was then to update network's wts & biases

'Iteratively reducing the error & improving the network's performance.'

forward:

a) h :-
net h
out h

b) o :-



$$(b) \boxed{w_1(\text{new}) = w_1(\text{old}) - \eta \frac{\partial E}{\partial w_1}}$$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial \text{out}_h} \times \frac{\partial \text{out}_h}{\partial \text{net}_h} \times \frac{\partial \text{net}_h}{\partial w_1}$$

already done

$$\frac{\partial \text{net}_h}{\partial w_1} = \frac{\partial (w_1 \text{out}_h)}{\partial \text{out}_h} = w_1 = 0.4$$

$$\frac{\partial \text{out}_h}{\partial \text{net}_h} = \frac{\partial (\frac{1}{1+e^{-\text{net}_h}})}{\partial \text{net}_h} = \text{out}_h(1-\text{out}_h)$$

$$= 0.679(1-0.679) = 0.218$$

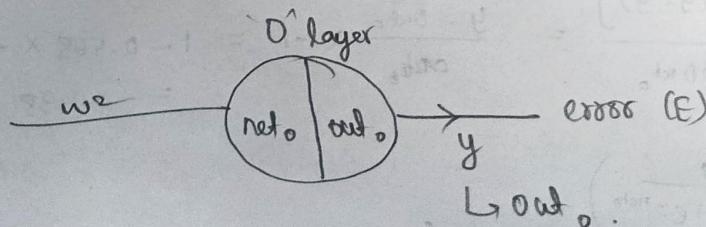
$$\frac{\partial \text{net}_h}{\partial w_1} = \frac{\partial w_1 \times x}{\partial w_1} = x = 5$$

$$\frac{\partial E}{\partial w_1} = -0.432 \times 0.245 \times 0.4 \times 0.218 \times 5 = -0.046$$

$$w_1(\text{new}) = w_1(\text{old}) - \eta \frac{\partial E}{\partial w_1} = 0.15 - 0.5 \times -0.046$$

$$= 0.15 + 0.023 = 0.173$$

$$\left. \begin{aligned} w_1(\text{new}) &= 0.173 \\ w_2(\text{new}) &= 0.436 \end{aligned} \right\} \text{new values}$$



$$E = \frac{1}{2}$$

Backpropagation

$$\frac{\partial E}{\partial w_2}$$

$$\frac{\partial E}{\partial w_1}$$

\Rightarrow

* Error function :- (Loss function)

discrepancy

- discrepancy b/w the predicted output of NN & actual target output for a given input
- provide a measure of how well the network is performing.

Common error fns include :-

→ MSE : used in regression problems

- calculate avg of sq. differences b/w predicted & actual values.

→ Cross-entropy ^{loss} ~~values~~:

- classification problems (Multi-class)

- diff b/w predicted prob. distribution & true distribution

→ backward propagation of errors :

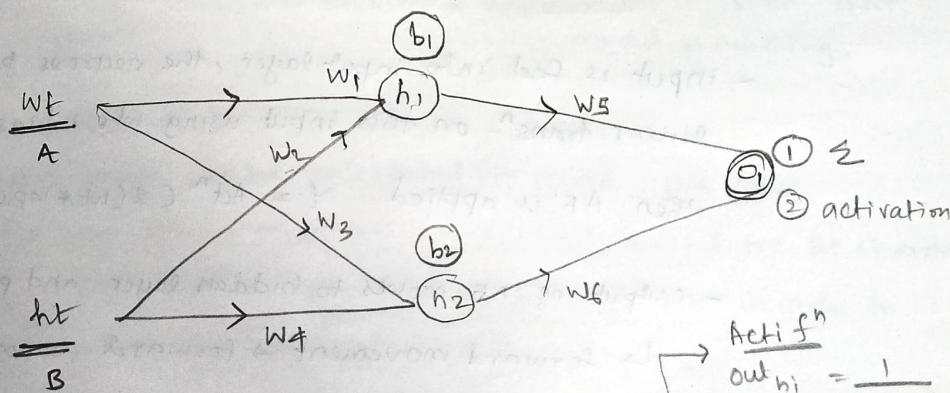
- used to eff. compute the gradients of error fn w.r.t wts & biases of NN

- these gradients indicate the dir & mag by which the wts & biases shld be adjusted to minimize the error.

Input layer

Hidden layer

Output layer



$$\text{net } h_1 = w_1A + w_2B + \delta_1$$

$$\text{net } h_2 = w_3A + w_4B + \delta_2$$

$$\begin{aligned} \text{Actf}^n \\ \text{out}_{h_i} &= \frac{1}{1+e^{-x}} \\ &= \frac{1}{1+e^{-\text{net } h_i}} \end{aligned}$$

$$\begin{aligned} \downarrow \\ \frac{1}{1+e^{-\text{net } h_2}} \end{aligned}$$

$$\text{net } o_1 = w_3 \text{ out } h_1 + w_6 \text{ out } h_2 + \delta_3$$

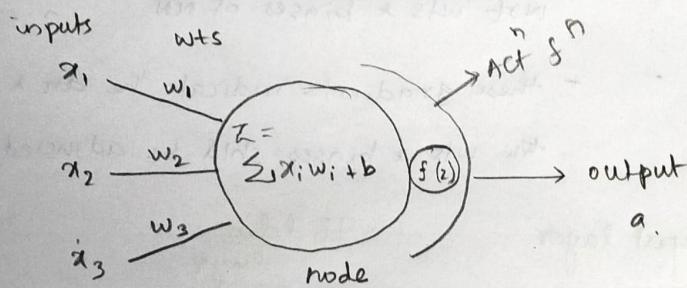
* Activation fns (AF) - linear, softmax, Tanh, ReLU,

- An AF decides whether a neuron shld be activated / not.
this mean that it'll decide whether the neuron's input to
the network is imp/not in the process of predn using
math- operations.
 - to transform the summed wt input from node into
an output value to be fed to the next hidden layer
or as output.

NN \approx brain \Rightarrow node = neuron (that receives a set of input signals - external stimuli)

⇒ brain based on \int processes them & decides whether neuron shld be act^d/fired or not.

- AF also called as transfer S^m



- input is fed into input layer, the neurons perform a linear transⁿ on this input using wts & biases.
then AF is applied $y = \text{Act}^n(\sum(wt * \text{input}) + \text{bias})$
 - output of AF moves to hidden layer and process is repeated
 - ↳ forward movement \rightarrow forward propagation.
 - what if output generated is far away from actual value
 - ↳ backward prop.

Using this error is calc., based on this ~~with~~ ^{with & biases of}
~~noises are unrelated~~

- A neuron w/o AF \rightarrow will perform only linear function on the w/o weighted input)
- Linear transⁿ \rightarrow simpler but net would be less powerful & not be able to learn ^{under} patterns from the data.
- NN w/o AF in DL is essential just as linear Reg. Model
 i.e. we use non-linearity fⁿ for non-linear transⁿ w/ the help of AF
- AF in NN \rightarrow wt sum of input is transformed into an output from a node or nodes in a layer of the network.
- If output range of AF is int⁺ \rightarrow squashing fⁿ
 Many AFs are non-linear \rightarrow non-linearity is the layer of network design
- choice of AF \rightarrow large impact on capability & performance of NN & diff AF fⁿs may be used in diff parts of model
- may use within/after internal processing of each node in network
- Network layers
 - Input = raw input from domain
 - hidden = input from another layer & pass it ^{to out}
 - output = makes predⁿ.

use AF \hookrightarrow diff AF & depends on the predⁿ that model is making.
- AF \Rightarrow differentiable
 - \Rightarrow 1st order can be calculated for given input value.
 - \Rightarrow req b/c given that NN are typically trained for BP of error algo that req. the derivate of predⁿ error in order to update the wts of the model.

* Types of AF

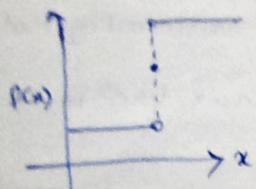
- Binary Step AF
- Linear AF
- Non-linear AF
 - $\tan-h.$
 - Sigmoid
 - ReLU
 - Softmax

① Binary Step fⁿ

- depends on x value that decides whether a neuron should be activated/not.

- input - fed to AF - compared w/ λ if input $> \lambda \rightarrow \checkmark$

$$\begin{cases} < \lambda \\ \geq \lambda \end{cases}$$



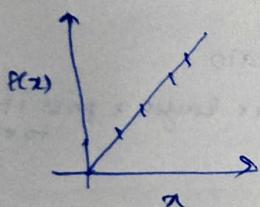
$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Its output isn't passed on to next HL

Limitations :-

- can't provide multi-value outputs (Multi-class Classification problem x)
- gradient of step $f^n = 0$
which causes a hindrance in BP pro.

② Linear/Identity AF



x is input

$f(x) \rightarrow$ output

$$f(x) = x \quad [\text{if } x=2 \text{ then } y=2]$$

- also known as "no activation fⁿ" or "identity fⁿ"

- AF \propto input

- f^n doesn't do anything, just simply splits out the value it was given.

Limitations :-

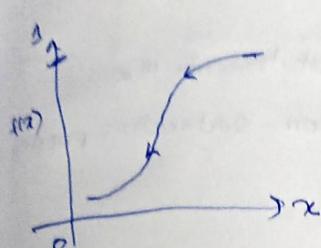
- not possible to use bp as the derivative of f^n is a const & has no relation to the input x.
- All layers of NN will collapse into one if a linear AF f^n is used.

will be a linear fn of the target. So, Linear AF turns NN to a layer (just).

② Non-linear Activation f^n :

- They solve linear AF problems.
- Allow us to now derive f^n w.r.t input & its possible to go back and understand which wts in input neurons can provide a better predn.
- Multi-layers of neurons as output would be non-linear comb'n of inputs passed thru multi-layers.
- any output can be represented as f^n computation in a NN.

(i) Sigmoid



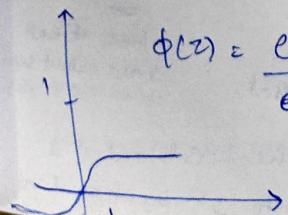
$$\phi(z) = \frac{1}{1+e^{-z}}$$

- commonly used for models where we've to predict the prob. as an output.
- Since prob exists only b/w 0 & 1, sigmoid is the ryt choice.
- f^n is diff'ble & provides a smooth gradient i.e., preventing jumps in output values.
 ↳ rep by S-shape of sigmoid f^n .

Limitations:

- δ value $\rightarrow 0$, network ceases to learn & suffers from vanishing δ prob.
- output of logistic f^n isn't symmetric arnd 0, so the output of all neurons will be of same sign (0 to 1) \Rightarrow this makes training unstable and difficult.

(ii) tan-h

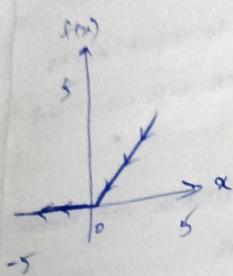


$$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

- same S-shape but values [-1, 1]
- it is zero centered, so we can easily map the outputs strongly -ve, neutral, strongly +ve
- used in HL as values [-1, 1]

! the mean for HL comes out to be 0.5
close to it.
— centering data makes learning much easier

(iii) ReLU (Rectified Linear Unit)



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

i.e., $x > 0 \rightarrow x$
 $x \leq 0 \rightarrow 0$

- ReLU has a derivative of 1
& allows bp with while simultaneously making the computation efficient.

- doesn't activate all neurons at the same time.

- neurons will only fire if the output of Linear trans < 0

Adv:
 - only certain is act^d,
 ReLU is far more computationally
 activated eff compared to Sigmoid,
 $\tan-h$

- accelerates overall convergence of S descent towards the global min of f^n due to its linear x non-saturating property

Disadv: - The dying ReLU problem.

(iv) Softmax AF

$$\text{Softmax}(z_j) = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, 2 \dots K$$

- able to handle multiple classes.
- Sigmoid's prob by a sigmoid are independent \Rightarrow
 \therefore not constrained to sum to 1
- Softmax input: $-0.5, 1.2, -0.1, 2.4$
 output: $0.04, 0.21, 0.05, 0.70$
 \approx
 why? coz they look at each raw output separately
- Whereas, Softmax's output are interrelated.

the sum will be always $\underline{=}$ \Rightarrow \uparrow^s the likelihood of 1

class, other has to \downarrow^s by an equal amt

* Optimization :- used for reducing the loss fn.

Algo :-

BGD

RMSprop

* Gradient Descent :-

- generalization of the derivative to the multivariate fn,
- descent \Rightarrow process of moving downward in fn's value i.e., to find min of a cost fn,
- local slope of f^n .

Batch GD :- computes S^n 's and updates model parameters using the entire dataset at each step. \rightarrow smooth / precise convergence but req significant computational power and memory.

Mini-Batch GD :- balances the trade-offs b/w BGD and Stochastic GD (SGD)

- updates model parameters using small subsets of data
- which in turn \Rightarrow faster training & better use of computational resources than BGD - (more stable convergence but less S^n -variance than SGD)
- hybrid approach \rightarrow ^{hy} eff & widely used tech / esp w DL for large datasets.

BGP

- entire dataset for update
- cost f^n reduces smoothly
- computation cost is very high

SGD

- single observation for update
- lots of variations in cost f^n
- computation time is more

MBGD

- subset of data for update
- smoother cost f^n as compared to SGD
- both are less time \rightarrow SGD cost \rightarrow BGP.

* RMS Prop :-

- we use adaptive learning rates.

- instead of a fixed η , learning rate changes based on past S^n 's.

* Bias-Variance Trade off

bias :- error from overly simple assumptions.

(model too simple \rightarrow underfitting)

var :- error from model being too sensitive to training data

Noise :

(model too complex \rightarrow overfitting)

Trade-off :- low bias \rightarrow flexible model, \uparrow risk of var.

low var \rightarrow rigid model, risk of \uparrow bias

So, we need to find a "sweet spot" where both are balanced \rightarrow best generalization.

Overfitting :- model learns training data too well, including noise & irrelevant patterns.

- \uparrow training accuracy but poor test accuracy.
- coz of too many parameters, insufficient data and lack of regularization.

* Regularization :-

- techniques to prevent overfitting by penalizing model complexity

ex :- L1 regularization (Lasso) \rightarrow adds penalty $\propto |\text{wts}|$
encourages sparsity scattered

L2 " (Ridge) \rightarrow " $\propto |\text{wts}|^2$
reduces large coeffs.

elastic net \rightarrow combⁿ of L1 + L2

- effect : shrinks wts \rightarrow reduces var at the cost of slightly \uparrow er bias
 \rightarrow improves generalization

* Generalization :-

- ability to perform well on unseen data

- opp of overfitting

- Achieved by balancing bias-var, using regⁿ, dropout & choosing good inductive bias

* Inductive Bias :-

- set of assumptions a learning algo uses to make pred's on unseen data
- w/o IB, learning from finite data is impossible.

ex:- LR assumes linear reln

Decision trees assumes data can be split hierarchically.

NN → complex patterns can be rep by layers of fn.

* Dropout :-

- A reg'ly tech used in NN ; randomly drops out (sets to 0) a subset of neurons during training
- prevents co-adaptation of neurons → reduces var
- acts like an ensemble of many smaller networks
↳ general'ly improves.

* Probabilistic Neural Network (PNN)

- used for classification & pattern recognition tasks
- Bayesian Stats & estimates PDFs for each classes, rather than learning wts thru iterative bp.
- A layer FNN (one dir w/o any cycles) In → out

Layers :- (2 feature : Apple & orange)

that needs to be classified

(i) Input layer :- just receives feature vector & feed to pattern layer

- no calculation

if fruit w. $w_b = 150\text{g}$, Color = 0.8

→ distn. the input feature to every neuron

↳ give to PL

in the next layer dim' n layer = 2

(ii) Pattern layer / Hidden layer :-

- memorizes the entire dataset & measure how similar the new inputs is to every single training example.

1 neuron rep ^
specific apple from T.D. - there is 1 neuron for each observation 100 fruits = 100 neurons
← each neuron is associated w. both a specific training vector & its known class

- calculates euclidean dist b/w

input vector (X) & its own stored training vector (x)

- then feeds the dist into Radial Basis fn (RBF) or Gaussian fn

- this fn outputs a value that rep the 1st act?

- closer the input is to neuron's stored vector, ↑ the output

e.g.: input [150, 0.8] compared w. neuron w [150, 0.75] apple ↑ value
orange ↓ value.

(iii) Summation layer :-

- avg the votes from PL for each class
- 1 neuron for each class & collects the output from all pattern neurons that belong to its class
- sum these outputs & avg them by total no. of patterns
- avoids unfair due to more training samples.

(iv) Output layer :-

- make final classification decision
- Choose the largest summation layer neuron.

1. Apple = 0.85
2. Orange = 0.12

→ Parzen Window Est" → PNN approx's the parent PDF for each class



adds the value
and approx to
a prob hill

(non-parametric)
tool
→ doesn't assume
anything

→ don't use bp for training

Numpy → python lib for prototyping & building NN

numpy.algorithms.PNN

check front

RNN → Recurrent NN

- process sequential data (text, speech / time series data)
- by maintaining an internal "memory" of past info.
- unlike FNN, this uses past info from 1 step to next to understand the context & dependencies within a sequence

* Hopfield Network

- RNN network invented by John Hopfield.
- both completes & finds errors.
- uses CAPM (Content Addressable Memory), stores & retrieve pattern based on partial or noisy inputs.
- single layer of interconnected neurons w/ symm wts, where the network's dynamics (driven by an energy fn) cause it to converge to stored mem. states.
- binary neurons (either -1 or 1 / 0 or 1)

brain can identify
somethings, like this do

every neuron connected
to others

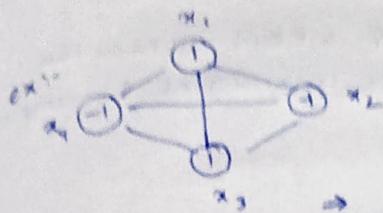
$$\text{net} = \sum_j (w_{ij} \times s_j) \quad \text{sum of all states} \times \text{connection wts.}$$

$$E = -\frac{1}{2} \sum_i \sum_j w_{ij} s_i s_j$$

$$\text{if net}_i \geq 0 \rightarrow S = +1$$

$$\text{else net} < 0 \rightarrow S = -1$$

set $w_{tj} \Rightarrow w_t > 0 \rightarrow$ same state
 $w_t < 0 \rightarrow$ diff.



Given M target N.S = $x_1^{(1)} \dots x_L^{(M)}$

$$w_{1j} = \frac{1}{M} \sum_{k=1}^M x_k^{(k)} X_j$$

$$\begin{bmatrix} 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{bmatrix} \quad w_{12}$$

$$\frac{1}{2} [1-1] = 0$$

Boltzmann Machine :-

- stochastic RNN that learns prob. dist. over its inputs.
- energy-based generative model that uses prob. approach to learn patterns in data by min energy fn
- consists of interconnected visible & hidden neurons learns PDFs & last energy state.
- due to its slow learning process, a restricted version called RBM limits connections b/w layers was developed to improve eff.
- neurons constitute a recurrent structure & they operate in binary manner \Rightarrow On-State $\Rightarrow +1$
Off-State $\Rightarrow -1$

Visible nodes :- receive input data from external world

Hidden nodes :- capture internal features & representations

All-to-all conn's \Rightarrow every neuron is connected to every other neuron

Boltzmann Distri :- $w_{VH} \propto e^{-E}$

$$V-V, V-H, H-H$$

- energy fn $\Rightarrow E = -\frac{1}{2} \sum_i \sum_j w_{ij} s_i s_j$ ∇S ∇E ∇s_i ∇s_j ∇w_{ij} ∇w_{VH} ∇w_{HV}
- chooses a random neuron 'k'
flips the state from x_k to x_{-k} \Rightarrow check the temperature T

$$P(x_k \rightarrow x_{-k}) = \frac{1}{1 + e^{-\Delta E_k / T}} \quad (\text{Sigmoid})$$

$\Delta E_k \rightarrow$ energy change

$T \rightarrow$ pseudo temp

\uparrow temp \rightarrow large explores widely

\downarrow temp \rightarrow \downarrow

$P(\text{flip}) = 1$ if $\Delta E_k > 0 \Rightarrow$ flip

$P(\text{flip}) \rightarrow 0$ if $\Delta E_k < 0 \Rightarrow$ don't

shed search
 \Rightarrow thermal eq

- energy based Model (\downarrow energy = more stable)
- sym. connections $\Rightarrow w_{12} = w_{21}$
- stochastic decisions \Rightarrow prob. decisions to be ON/OFF
- learning thru energy minimization \Rightarrow learns by adjusting its internal wts & biases $\rightarrow \downarrow$ overall energy.

* Deep BM

- deep NN based on BM
- stacked RBM, where output of 1 RBM serves as the input for the next, allowing for learning of highly complex data representations.

* AutoEncoders in DL

- unsupervised NN that learn to reconstruct their input data
- compression algo, learning eff data representations by encoding and decoding the data.
- ignore noise in data.

Encoder :-

- a module that compresses the train validated test input data into a encoded rep. i.e., several orders of mag smaller than the input data.

ex:- Img \rightarrow convolutional layers + pooling layers. 1st dim \rightarrow 1st dim

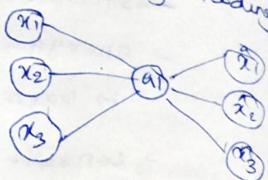
- set of convolutional blocks + pooling layers that compress the input to model into compact section called bottleneck

Bottleneck :-

- latent space
- contains compressed data (knowledge rep)
- captures only most imp features of data to allow successful reconstruction.
- a bottleneck as a compressed rep of an input further prevents the NN from memorizing the input & overfitting on the data.
- Smaller the BN \rightarrow lower the risk of overfitting.
- however, v. small BNs would restrict the amt of info storabl, which \Rightarrow the chances of imp info slipping out thru the pooling layers of encoder.

Decoder :-

- helps the network decompress the knowledge rep & reconstructs the data from its encoded form.
- output is then compared w/ ground truth
- set of upsampling & convolutional blocks that reconstruct



APPENDS :-

- Dim Redⁿ : easy analysis & visualization
 - Feature Extraction
 - Noise Removal : remove noise from data.
 - Anomaly detection : reconstruct → can identify anomalies as input that they can't reconstruct accurately.
 - Generative Models : generate new data

Hyperparameters b4 training autoencoder:

- Code size / BN size : most imp
it decides how much data to be compressed
Small : Strong ; Large : Weak
risk of losing imp info risk of overfitting

- No. of Layers : depth of encoder = decoder
per depth yes model complexity
↓ or → faster to process

- No. of nodes per layer:- wts we use per layer
no. of nodes los w subsequent layer
as input becomes smaller

- Reconstruction loss :- loss fn we use \rightarrow dependent on type of input & output we want the auto to adapt to.

image \Rightarrow MSE loss + L1 loss

Output is range [0, 1] \Rightarrow MNIST,

Binary cross entropy.

Types :-

- Std AE :- basic AE that learns a compressed rep of input data.
 - De-noising AE :- Trained to reconstruct the original data from noisy data.
 - Variational AE :- learns prob-rep of data
 - Convolutional AE : imgs & other spatial data.

* Probabilistic Graphical Models :-

linear chain :- each node depends on its ^{neighbour} ~~know one~~ → sequence

Cyber :- sys logs → user act → packets / @ time

for anomaly detection.

open → read → write → close → Normal

open → read → exec → delete → Malicious behav.

Partition f^n :- \star normalization const used in prob. models
(2) prob sum = 1

Cyber :-

eval prob of given attack pattern, partition f^n helps in normalize over all possible sys states
helps in ranking diff attack hypotheses.

Malware detection model $\Rightarrow \underline{1}$

Markov random field :- nodes & their dependencies w/o assuming spatial/relational dependency direction

Cyber :- can detect inter-dependencies among devices/hosts in a network

useful in DIDS (distributed Intrusion Detection System)

Combined attacks :-

- A → abnormal traffic
- B → logs unknown outbound connection
- C → shows privilege escalation.

Belief Propagation :- (BP)

Algorithm used in PGMs to compute marginal prob.

Cyber :-

helps in propagating across the network of events, one event is likely to happen.

ex :- phishing email leading to cred compromise triggering lateral movement

^A

^B

BP can propagate the belief from A to B,
↑ detection confidence.

Training Conditional Random Fields :-

~~discriminative~~ probabilistic model structured in linear chain / graph

- involves max. the conditional likelihood of labeled sequence.

Cyber:

CRFs used to label seq. network traffic flows/
they learn dependences b/w features & labels.
Open \rightarrow chmod \rightarrow sus.

Hidden Markov Model (HMM):- hidden states (normal, attack) and observable outputs (log entries, pkt features)

Cyber: used for behav modelling, intrusion detection.

HS :- {normal, exploit, data exfiltration}

Obs :- {port scanning, unusual login, large file transfers}.

Entropy :- measures uncertainty/randomness of data.

Cyber:-
+ entropy in files \Rightarrow likely encrypted/malware
+ in Net traffic \Rightarrow normal Comm?

Sudden change \Rightarrow potential intrusion.