

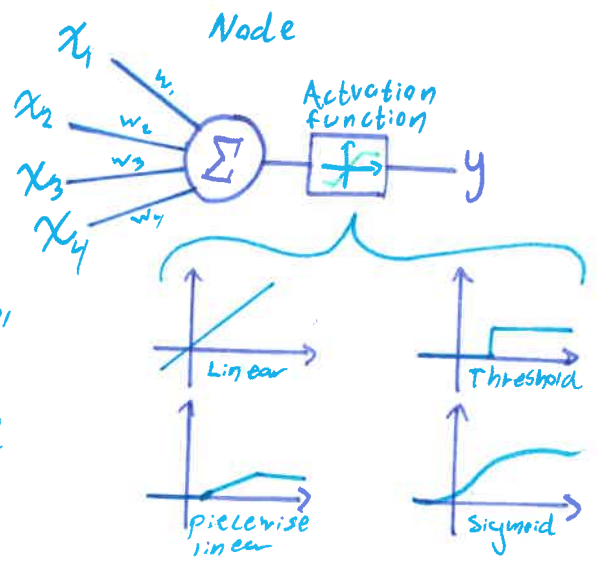
Fundamentals

Node - Smallest computational Unit.

Activation function - When should a Node "fire"

Learning Rule - How should success and fails be handled

Topology/Architecture - How the network is structured



Learning Principles

Supervised - Input with labels

Unsupervised - Input without labels.

Reinforcement - Reward given on success.

Single layer



Multiple layer

Recurrent

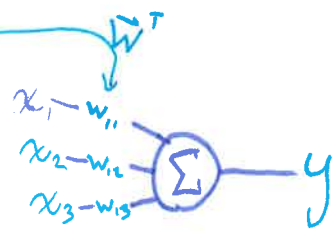


Linear Networks

Weight matrix

$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}$$

sender W_{xy}
receiver



$$\vec{W}^T \cdot \vec{x} = y$$

$$W \cdot \vec{x} = \vec{y}$$

Note: Chained single layer networks with weights w_1, w_2, w_3 have a total weight:

$$\vec{y} = w_3(w_2(w_1 \vec{x})) = (w_3 w_2 w_1) \vec{x}$$

in short, this is just W another single layer network.

Hebb's Learning Hypothesis

Neurons which activate together should strengthen their connection.

$$\Delta W_{ij} = x_i y_j$$

if sign of $x_i y_j$ is the same their product is positive, else it is negative.

Same as

$$\Delta W_{ij} = (x_i - \bar{x}) (y_j - \bar{y})$$

"Fire Together Wire Together"

Storing Mappings

$$W = \sum_{p=1}^n \vec{y}^{(p)} \cdot \vec{x}^{(p)T}$$

$$y_{out} = W \vec{x}^{(k)} = \sum_{p=1}^n \vec{y}^{(p)} (\vec{x}^{(p)T} \vec{x}^{(k)})$$

$$= y^{(k)} (\vec{x}^{(k)T} \vec{x}^{(k)}) + \underbrace{\sum_{p \neq k} \vec{y}^{(p)} (\vec{x}^{(p)T} \vec{x}^{(k)})}_{\text{close to zero if } \vec{x}^p \text{ and } \vec{x}^k \text{ are orthogonal.}}$$

$$\approx \alpha y^k$$

close to zero if \vec{x}^p and \vec{x}^k are orthogonal.

Preception Learning

1 Data is correct; Do nothing.

2 Data is false positive $\Delta W = -\vec{x}$

3 Data is false negative $\Delta W = +\vec{x}$

Convergence

If solution exists it will be found in finite number of steps.

Delta Rule (Widrow Hoff)

While preception learning will always terminate with linearly seperable data, it is not always the most intuitive solution. Delta rule will always converge on an intuitive solution by minimizing an error function.

1 symmetric targeting values $\{-1, 1\}$

2 Error is measured before threshold
 $e = t - \vec{w}^T \vec{x}$

3 Minimize weights according to
 $E = \frac{e^2}{2}$ (steepest decent)

$$\frac{\partial E}{\partial \vec{w}} = -e \vec{x}$$

$$\Delta \vec{w} = \eta \frac{\partial E}{\partial \vec{w}} = \eta e \vec{x}$$

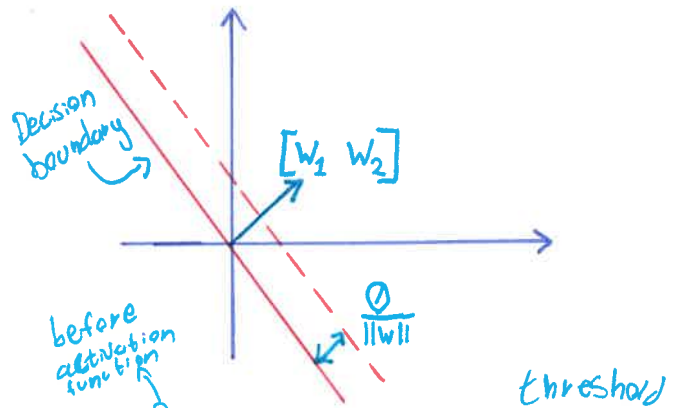
learning rate

Threshold Logic Unit

McCulloch pitts neuron (1942)

Introduction of the threshold function.

Add function to output to map some values to 0 and some to 1. Also creates a decision boundary.



$$y = w_1 x_1 + w_2 x_2 \dots w_n x_n - \theta$$

$$0 = w_1 x_1 + w_2 x_2 \dots w_n x_n + y' \quad \text{Plane equation}$$

θ becomes an additional input to the node and acts as a bias.

Multilayer Networks

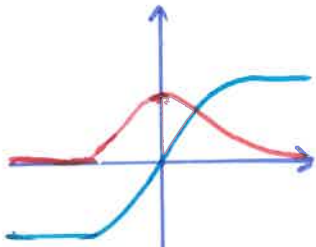
Dilemma:

- > Thresholding destroys information needed for delta rule.
- > No thresholding loses functionality of multiple layers

Solution:

Use threshold-like but differentiable function.

$$\varphi \quad \varphi'$$



$$\varphi(x) = \frac{1 - e^{-x}}{1 + e^{-x}}$$

$$\varphi'(x) = \frac{1 - \varphi(x)^2}{2}$$

Sequential & Batch → stochastic gradient descent

Either error is accumulated and then updated or after every sample.

Sequential:

- faster but needs higher learning rate.
- less likely to get stuck in local minima.
- Is not true gradient descent

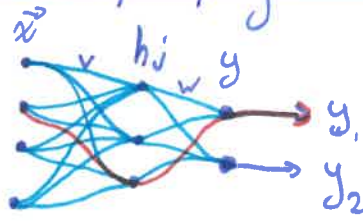
weight should be randomly initialized

$$N(0, \sigma)$$

$$\sigma = \frac{1}{\sqrt{fan-in}}$$

← number of inputs for node

Backpropagation



1 forward pass; compute h_j & y_k

$$h_j = \varphi\left(\sum_i v_{ji} \vec{x}_i\right) \quad y_k = \varphi\left(\sum_j \vec{w}_{kj} h_j\right)$$

2 Backward pass; compute δ_k & δ_j

$$\delta_k = (t_k - y_k) \cdot \varphi'(y_k^{in})$$

$$\delta_j = \sum_k \delta_k \cdot w_{kj} \cdot \varphi'(h_j^{in})$$

Convenient if φ' can be expressed in terms of φ

3 update; Δw & Δv

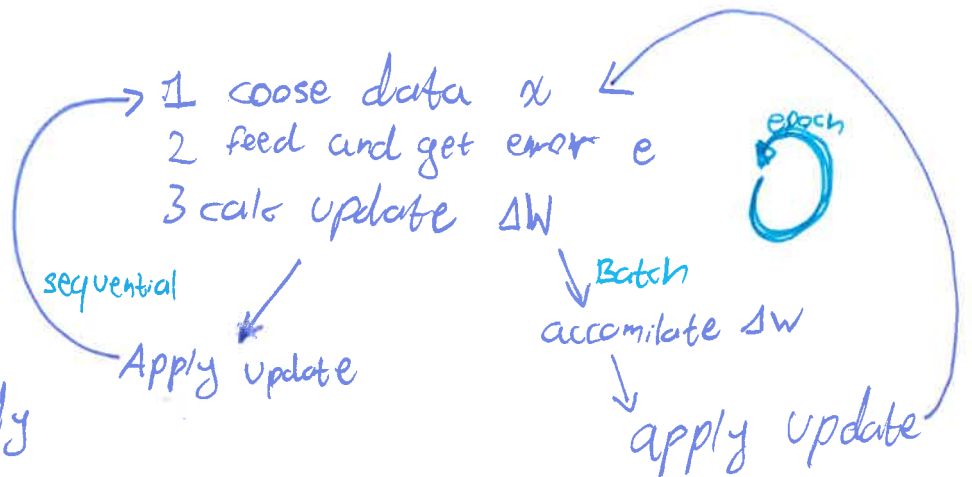
$$\Delta w_{kj} = \eta \delta_k h_j \quad \Delta v_{ji} = \eta \delta_j \vec{x}_i$$

Problems:

Convergence is slow
getting stuck in local minima.

Many parameters must be tuned

Challenging scaling
Biologically unrealistic.



Learning Rate η

too low \rightarrow learning slow

too large \rightarrow overshoot convergence

Adaptive $\rightarrow \eta(t) = \eta(1)/t$ or $\eta(t) = \eta(1)/(1+t/\alpha)$

\hookrightarrow depending of sign of consecutive iterations!

Momentum $\Delta \vec{w} = \beta \Delta \vec{w} - (1-\beta) \frac{\partial e}{\partial \vec{w}}$

\hookrightarrow previous update affects new update
keeps old changes for a while.

overrepresenting hard data.

Adding noise during training.

Hyperparameter selection.

Generalization

Our model should work on data never seen before.

Key factors:

\hookrightarrow Quality & size of training data

\hookrightarrow Complexity of model

\hookrightarrow Complexity of problem

Classification RBF

Training for MLP involves changing weights and biases.

Training for RBF means changing locations and radii.

Competitive Learning

Winner takes all.

node closest to the input pattern learns by the rule.

$$\Delta w = p \vec{x}$$

Pos of RBF \nearrow or \searrow

$$\Delta w_a = \eta(\vec{x} - w_a)$$



Vector Quantization

Represent data in terms of a limited number of typical data vectors.

- > Compression
- > noise Reduction.

Spills-over Learning

When learning, nodes connected to winning nodes (neighbour) also learn a bit.

$$\Delta w = \eta h(\vec{x} - w)$$

\uparrow
spill rate depending on neighbourness



Important!

neighbours are determined by the output space NOT input space

(winner is still in the input space)

Cover's Theorem

Projecting a pattern to a higher dimension will cause problems to be solved easily.

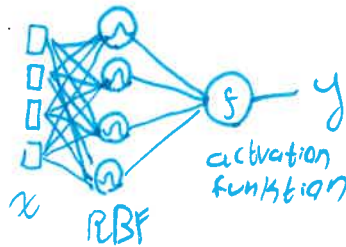
RBF

Use non-linear mapping function.

$$\varphi_i(\|\vec{x} - \vec{x}_i\|)$$

\searrow RBF centre

$$F(x) = \sum_{i=1}^N w_i \varphi(\|x - x_i\|)$$



Learning Algorithm.

fundamental principle:

Update winning node to make it more specialized.

Algorithm.

- Initialize clusters randomly

w_1, \dots, w_c

- For each input vector, find winner weight. (in the input space)
- update weight using rule.

Result

clusters in data is found and protected by nodes



Problems

Dead or unuseful nodes.

Self Organizing Maps (SOM)

- | | |
|------------------|----------------|
| 1 competition | } output space |
| 2 cooperation | |
| 3 weight updates | } input space |

\rightarrow Topological ordering

\rightarrow high learning rate

\rightarrow large neighborhood

} global order

\rightarrow Convergence

\rightarrow low learning rate

\rightarrow small neighborhood

} Local fit

Learning Vector Quantization

Supervised version of VQ
when classes are known.

$$\Delta \vec{w}_i = \eta (\vec{x} - \vec{w}_i)$$

when classified
correctly

$$\Delta \vec{w}_i = -\eta (\vec{x} - \vec{w}_i)$$

when classified
incorrectly

Hopfield Network (Hobbin)

All nodes are connected to every
other node.

$$\vec{x} = \text{sgn}(Wx + \theta)$$

$$E(\vec{x}) = -\frac{1}{2} \sum_i \sum_j w_{ij} x_i x_j + \sum_i \theta_i x_i$$

Asynchronous update:

$$x_i(t+1) = \text{sgn} \left(\sum_j w_{ij} x_j(t) \right)$$

i is random

$$\Delta E_{x_i \rightarrow x_i^*} = -\frac{1}{2} \left(\sum_j w_{ij} x_i x_j^* - \sum_j w_{ij} x_i x_j \right)$$

$$= -\frac{1}{2} (x_i^* - x_i) \sum_j w_{ij} x_j \leq 0$$

Synchronics can flip

Memory capacity

Memory capacity M nodes
is proportional to $0.138n$

$$M \leq \frac{n}{4 \ln(n)}$$

Sparsely patterns $n \times \log(n)$

Associative Pattern Recognition

Hetero Associative: boat \rightarrow water

Auto Associative: boat \rightarrow

P.R. networks: boat 1 \rightarrow boat 2

Works when \vec{x} are orthogonal

Auto associative case:

$$W = X X^T$$

\vec{x} is eigenvector of W

$$\text{sgn}(W\vec{x}) = \vec{x} \rightarrow$$

W describes non-orthogonal projection of
subspace spanned by \vec{x}

Resonance, Energy (BAM)

$$\vec{y}_0 = W \vec{x}_0 \quad \vec{e} = W^T \vec{y}_0$$

how far is e from x ?

$$E = -\vec{x}_0^T \vec{e} = -\vec{x}_0^T W^T \vec{y}_0 = -\vec{y}_0^T W \vec{x}_0$$

Converges if W is symmetric

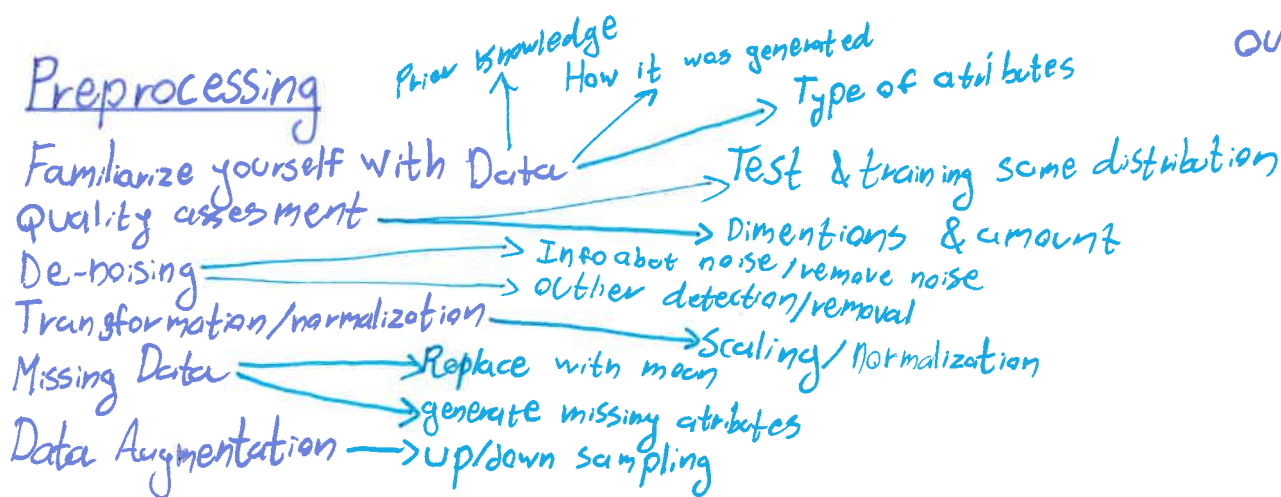
Memory state are fixed point attractor

Batch mode is faster

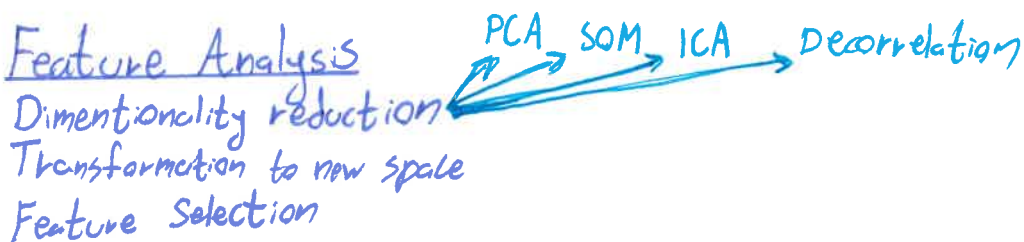
Pattern Recognition Pipeline

Data → Preprocessing → Feature Analysis → Classification / Regression → Postprocessing
 ↓
 output

Preprocessing



Feature Analysis



Classification / Regression

Choosing correct error and loss function

- ↳ Depends on: problem type (regression/classification)
- ↳ Expected value in output layer (function / num outputs)
- ↳ learning algorithm

relu sigmoid

Maximum likelihood

↳ How close the distributions of predictions by model match with data

↳ Cross-Entropy between empirical distributions defined by training set, and Probability distribution of model.

↳ MSE is the same for regression

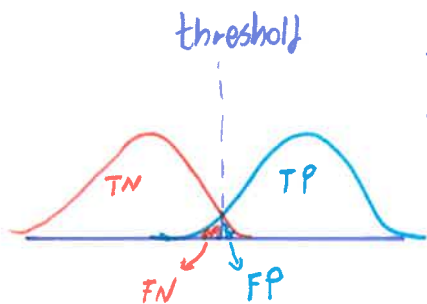
Cross Entropy

↳ sigmoid activation of output

Binary classification → $L(y, t) = \sum_{i=1}^N t^i \log y^{(i)} + (1-t^i) \log(1-y^i)$

Multiple classes → $L(y, t) = \sum_{i=1}^N \sum_{c=1}^C t_c^i \log(y_c^i)$

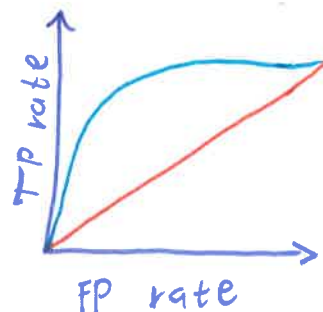
Receiver Operator Characteristics



Threshold can be moved to decrease FP or FN, but increase the other.

$$\text{Specificity} = \frac{TN}{TN+FP}$$

$$\text{Sensitivity} = \frac{TP}{TP+FN}$$



Why deep learning

Deep learning performs better with more data while other learning plateaus eventually.

Depth - longest path from x to y

Layer-Wise Training

Only Empirical proof



Train one layer at the time while keeping others locked. Unsupervised.

Then tune with backprop.

Pretraining minimizes variance
Controls complexity
implicit penalization (Regularization)

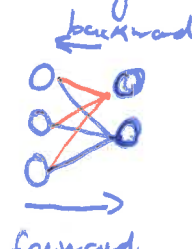
2010

2007

Pretraining has better initial condition training of entire architecture

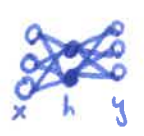
RBM

Two layer autoencoder



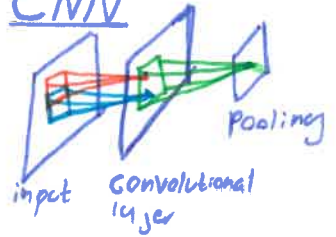
features are extracted in hidden layer and then reconstructed in backward pass

Autoencoder



Unsupervised (still gradient descent)
 x should equal y despite dimension reduction.

CNN



Multiple layer allow for hierarchical feature extraction.
Multiple sampling.

Problems

- Vanishing gradients in backprop algorithm with more layers added.
- Unstable gradients
- local minima
- overfitting

2006

↳ pre-train deep architectures with layer-wise unsupervised learning.

Classical Architectures

Restricted Boltzmann Machine

↳ Contrastive divergence for pre-training

Autoencoder

↳ gradient descent based alg for pre-training

Convolutional Neural Network

↳ No need for pretraining

↳ often uses transfer learning

Representation Learning

Many levels of representation allow multitask-learning

Distributed representation

$P(\text{data} | \text{latent})$ - generation

↳ learn features

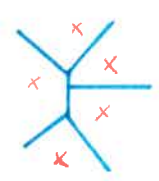
↳ layers of abstraction

↳ multitask/transfer learning

$P(\text{latent} | \text{data})$ - recognition

↳ non-local generalization (sparse code) (multi clustering)

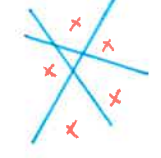
Local cluster



Hard boundaries

Class is not determined by a single feature.

Multi cluster



Soft (many) boundaries

Residual activity prevent overfitting

Filters convolve (slide) across input to analyse said input and sub sample

Distributed vs local representation



Distributed
 Information is spread across multiple units
 Represents similarity between objects

More easily separate (orthogonal)

- implications
- ↳ Transfer learning
 - ↳ Multitask
 - ↳ zero shot

Supervised

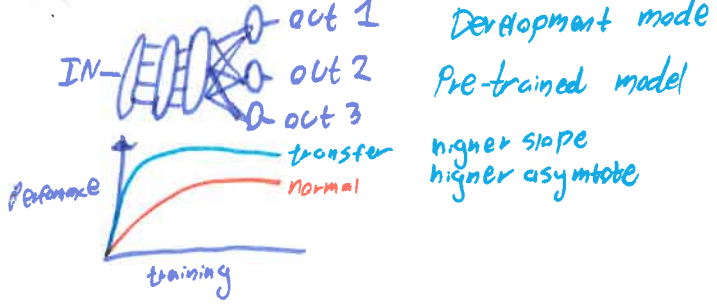
Very useful at solving specific tasks
 Relies on labeled data
 Information bottleneck
 Lacking "common sense"

Unsupervised

Learning world before tasks
 Less susceptible to information loss

Transfer learning

A model developed for one task is reused as the starting point of another



Multitask learning

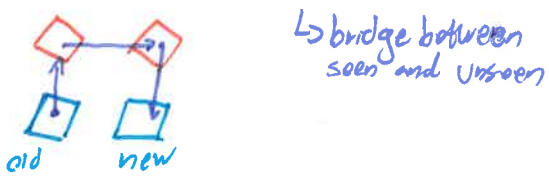
Learn multiple tasks simultaneously
 Improves generalization by using domain specific information contained in training signals



Zero-Shot learning

Recognize objects from unseen classes.

- Available:
- ↳ seen data (labels)
 - ↳ Unseen data (no labels)
 - ↳ Auxiliary information from seen and unseen classes



Autoencoders

$$x \begin{matrix} \diagdown \\ \diagup \end{matrix} \begin{matrix} w \\ w \end{matrix} x \quad \text{Loss}(x, g(f(x)))$$

Encoder $x' = f(x) = w'h$
 Decoder $h = g(f(x)) = wx$

Undercomplete

↳ hidden layer < input

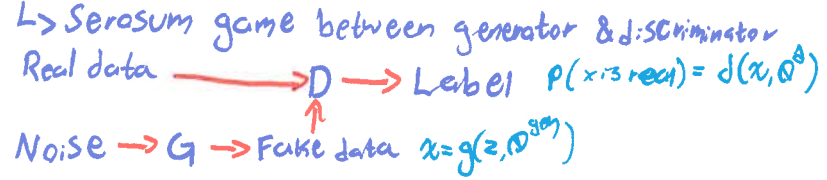
Overcomplete

- ↳ hidden layer > input
 - ↳ needs regularization
 - ↳ sparse & denoising
 - ↳ Deep autoencoder
- Fixed by adding latent variable with a prior and maximize likelihood

Generative modeling

Objective is to capture hidden structure in data

Generative Adversarial network



Variational Autoencoders

an autoencoder with generator capabilities

$x \rightarrow z \rightarrow x'$ last layer of q define mean and variance of latent space

Vector arithmetic in latent space

Linear Networks.

Supervised learning
 $y = \vec{w}^T \cdot \vec{x}$

Stacking multiple does not work (becomes new linear mapping)

Storing Mappings

↳ Program Resides in weights

↳ Learning Adapts weights

Hebbian Learning

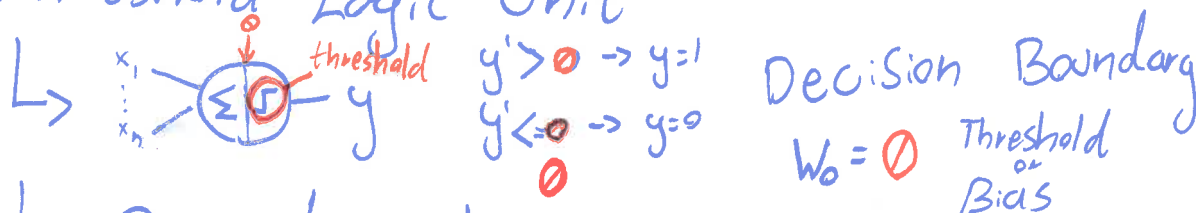
↳ $\Delta w_{ij} = x_j y_i$

↳ Fire together wire together.

↳ $\Delta W = \vec{y} \times \vec{x} = \vec{y} \vec{x}^T = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} = \begin{bmatrix} x_1 y_1 & \dots & x_3 y_1 & \dots \\ x_1 y_2 & \dots & x_3 y_2 & \dots \end{bmatrix} = \begin{bmatrix} \Delta w_{11} & \dots & \Delta w_{15} \\ \Delta w_{21} & \dots & \Delta w_{25} \end{bmatrix}$

↳ Remember orthogonal patterns

Threshold Logic Unit



↳ Perceptron learning

↳ Update weights depending on classification correctness

↳ Result=0 label=1 $\Delta W = p \vec{x}$ Result=1 label=0 $\Delta W = -p \vec{x}$ Correct classification do nothing

↳ Always converges for linearly separable data.

↳ Delta Rule

↳ Results $\in \{-1, 1\}$, error = $t - \vec{w}^T \vec{x}$ minimize $E = \frac{e^2}{2}$

↳ Steepest descent move down in gradient $\Delta \vec{w} = -\eta \frac{\partial E}{\partial \vec{w}} = \frac{\partial}{\partial w} E(e(w)) = -e \vec{x}$

↳ finds lowest error (not necessarily correct classification)

$$\frac{\partial}{\partial w} = -e \vec{x}$$

$$\Delta w = \eta e \vec{x}$$

↳ Cannot solve xor problem

Hopfield Networks

Learn and remember patterns. Autoassociative

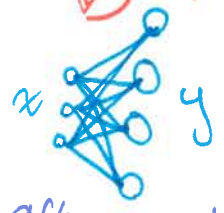
~~$Y = WX$~~ $\vec{x}^{(i+1)} = Wx^{(i)}$
recurrent network

Bipolar encoding $x \in \{1, -1\}$

$$y = Wx_p = \sum_k (y_k x_k^T) x_p = \sum_p (x_p^T x_p) y_p + \sum_k \underbrace{y_k (x_p^T x_k)}_{\text{if patterns are orthogonal}}$$

Hebbian learning

loop many times



$W = X X^T$
↳ matrix of patterns

$\text{sgn}(W\vec{x}) = \vec{x}$ $\text{sgn}(WX) = X$
↳ \vec{x} are eigenvectors

after many iterations x will converge to an "attractor"

$\forall w_{ij} = 0$
 $E(\vec{x}) = -\frac{1}{2} \sum_{i,j} w_{ij} x_i x_j + \text{const}$

Bidirectional Associative memory

Will always converge since energy decreases (when w is symmetric)

$\Delta E(x_i \rightarrow x_i^*) = -\frac{1}{2} (x_i^* - x_i) \sum_j w_{ij} x_j \geq 0$

opposites of attractors are stable

$[-1]$ & $[1]$

Synchronous & asynchronous

given patterns x_1, x_2, \dots, x_n

$W_{ij} = \begin{cases} \frac{1}{n} \sum_k x_{ki} \cdot x_{kj} & \text{if } i \neq j \\ 0 & \text{if } i = j \end{cases}$
Hebb learning

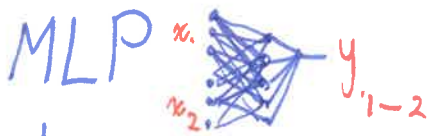
Patterns must be orthogonal for best performance

Patterns $M \leq 0.138n$ nodes $n \log(n)$
one-shot learning sparse patterns

$M \leq \frac{n}{4 \ln n}$

Temporal Processing

Discrete & Continuous. Time domain data is important

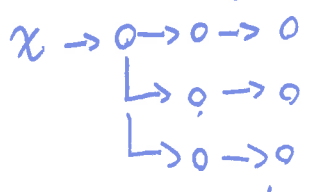


Time lagged feed forward network

- ↳ Multiple different times in the input
- ↳ only some exact times are allowed
- ↳ gives ability to add static nodes

Recurrent Architecture Backprop through time

↳ Multiple layers talking to each other



for training the entire time series is needed since you cannot know future samples.

Kalman filters

vanishing gradient (very bad) since w is multiplied many times and eigenvalues λ are usually less than 1 so $\lambda^n \rightarrow 0$ for large n .

Echo State Network

↳ Reservoir computing



Reservoir is Very large 10k nodes
 sparse 2-20% of nonzero elements
 spectral radius $\rho < 1$ forgetting speed

Nodes have leaky memory modulated by α $h(t) = \alpha h(t-1) + (1-\alpha)x(t)$

one shot learning with least squares (LMS)

fill and activate reservoir with $x(n)$, compute $h(t)$, fit y with LMS

↳ Long Short-term Memory

very good, hard to understand
 backprop is possible because everything is differentiable

Multi-layer Networks

$$\varphi(x) = \frac{1 - e^{-x}}{1 + e^{-x}} \quad \varphi'(x) = \frac{1 - \varphi^2}{2}$$

$$\varphi(x) = \frac{1}{1 + e^{-x}} \quad \varphi'(x) = \varphi(x)(1 - \varphi(x))$$

$$\varphi = \tanh(x) \quad \varphi'(x) = 1 - \varphi^2$$

$$\varphi = \arctan(x) \quad \varphi' = \frac{1}{1 + x^2}$$

Instead of thresholding use differentiable function

1 Define error function $E = \frac{1}{2} \|\vec{t} - \vec{y}\|^2 = \sum_k (t_k - y_k)^2$

2 Minimize E with steepest descent (towards minima)

$$\Delta w = -\eta \frac{\partial E(\vec{x}, \vec{w})}{\partial \vec{w}} \quad \Delta w_{kj} = -\eta \frac{\partial E}{\partial w_{kj}}$$

$$\Delta v_{ji} = -\eta \frac{\partial E}{\partial v_{ji}}$$

Error of previous layer depends on the next, however it diminishes with each layer (vanishing gradient)

$$\delta_k = (t_k - y_k) \cdot \varphi'(y_k^{in}) \quad \delta_j = \sum_k \delta_k \cdot w_{kj} \cdot \varphi'(h_j^{in})$$

$$\Delta w_{kj} = \eta \delta_k h_j \quad \Delta v_{ji} = \eta \delta_j x_i$$

1 Forward pass: Compute each hidden layer and output

2 Backward pass: Compute δ_k & δ_j

3 Update weights: $\Delta w_{kj} = \eta \delta_k h_j$ $\Delta v_{ji} = \eta \delta_j x_i$

Sequential

- ↳ update w after each sample
- ↳ faster but needs lower learning rate
- ↳ Less likely to get stuck in local minima
- ↳ not true gradient descent, but rather stochastic gradient descent

Adaptive Learning Rate

↳ num iterations $n(t) = \frac{n(1)}{t}$ $n(t) = \frac{n(1)}{1 + t/\eta}$

↳ Momentum $\Delta w = \beta \Delta w - (1 - \beta) \frac{\partial E}{\partial w}$

↳ antisymmetric activation function old w Small change converge faster.

Radial Basis Function

$$\mathbb{R}^M \rightarrow \mathbb{R}^N \rightarrow \mathbb{R}^1 \quad N > M$$

Separate mappings that are nonlinear after mapping them to a high dimension using a linear separability

$$\varphi_i(\|x - x_i\|)$$

φ_i → gauss
 $\|x - x_i\|$ → gauss of distance from x_i
 x_i → RBF centre

Linear operation in N dimensional space: $F(x) = \sum_i^N w_i \cdot (\text{RBF})$

Unsupervised learning, mainly for clustering

Number of RBF need to be less than num samples but larger than input size. \Rightarrow more samples than dimension is needed

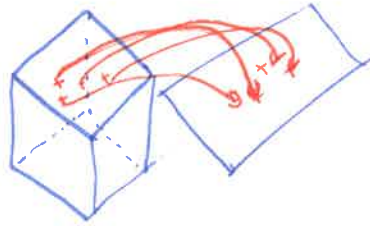
Winner takes all, winner is the closest RBF to the input data $\|x - w\|$
 Dead units are units that never win. (Fixed by reinitialize, or give a node a point to start at)

Self Organizing Map

hyperparameter: neighborhood size
 learning rate

Vector Quantization - let single vector represent a large area
 \hookrightarrow compression / noise reduction.

When a node wins it will spill over to its neighbours in the node space
 This means nodes close in input will remain close in output despite (output space) the dimension up or down scale.



Learning Vector quantization

$$\Delta w = +\eta(\tilde{x} - \tilde{w}) \quad \Delta w = -\eta(\tilde{x} - \tilde{w})$$

when correctly classified when incorrectly classified

Decreasing learning rate and neighbourhood size

Deep Belief Networks

Create a network that maximises $p(x)$ not $p(y|x)$

Make inferences instead of learning training data.

Stacked RBM \rightarrow Energy based

Train layer by layer using unsupervised methods

Tune after, using supervised learning and labels.

\hookrightarrow Works because pretraining minimizes variance

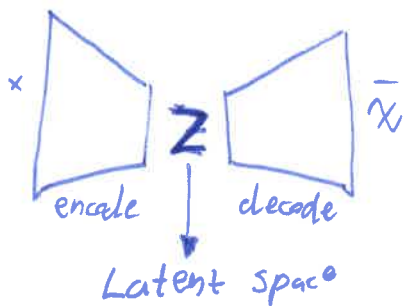
\hookrightarrow Controls complexity

\hookrightarrow implicit penalty and acts as regularization (less susceptible to overfit)

\hookrightarrow Better initial position in terms of local minima

Contrastive Divergence

Variational autoencoders



generation capabilities comes from having probability distribution in the latent space $P(x|y)$

Vectors in latent space can be decoded into the real space

Undercomplete - Dense

overcomplete - sparse

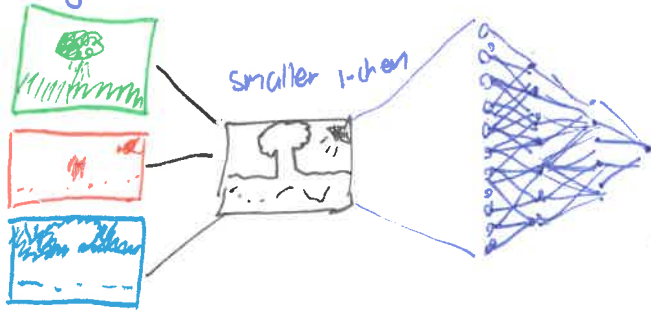
\hookrightarrow avoid trivial connection

\hookrightarrow L1 regularization

Convolutional Neural Networks

Take original input (image) and apply filters across. (Convolution)

Big 3-channel



Representation Hierarchy

Generative Adversarial Network

Real image \rightarrow Discriminator \rightarrow label

Generator \rightarrow Fake image \uparrow

Probabilistic

Regularisation & Bayesian Techniques

Generalization is the ability to apply learned knowledge in unseen situation

↳ Training size, complexity NN, complexity of problem

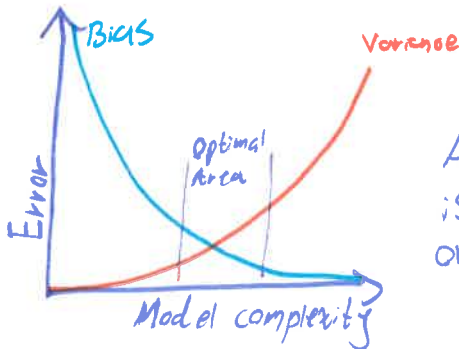
↳ Underfitting / overfitting

Bias: Difference between average prediction and the correct value.
High bias implies the model oversimplifies the problem and underfits training results giving high error in training

Variance: The variability of model prediction on a given data point.
High variance is very good with the training data but does not generalise well on unseen data overfit

True Relationship $Y = f(x) + e$ ↓ gaussian error Modelled relationship $Y = \hat{f}(x) + e$

$$\text{Error}(x) \quad E[(Y - \hat{f}(x))^2] = \underbrace{(E[\hat{f}(x)] - f(x))^2}_{\text{Bias}} + \underbrace{E[(\hat{f}(x) - E[\hat{f}(x)])^2]}_{\text{Variance}} + \underbrace{\sigma_e^2}_{\text{irreducible error}}$$

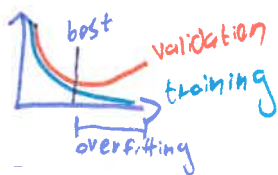


Adding noise in every sample right before it is seen minimises the variance since network cannot overfit to data seen once.

Early stopping

Create an additional dataset called Validation set

Train network on training data until validation error is minimum



Early stopping can also be based on problem-specific limits.

Network Growing and pruning

Growing - Add more units or layers when specification is not fulfilled

Pruning - remove parts of networks based on saliency

Penalised learning
Add a penalty term to error function

$$\tilde{E} = E + \lambda \Omega$$

Penalty
factor

L2 Regularisation

$$\Omega = \|W\|^2 = \sum_j w_j^2 \Rightarrow \tilde{E} = \frac{1}{N} \sum_i (t_i - y_i)^2 + \lambda \sum_j w_j^2$$

gradient: $\frac{\partial \tilde{E}}{\partial w_j} = \frac{\partial E}{\partial w_j} + 2\lambda w_j$

L1 regularisation

$$\Omega = \|W\| = \sum_j |w_j|$$

$$\tilde{E} = \frac{1}{N} \sum_i (t_i - y_i)^2 + \lambda \sum_j |w_j|$$

gradient: $\frac{\partial \tilde{E}}{\partial w_j} = \frac{\partial E}{\partial w_j} + 2\lambda \text{sgn}(w_j)$

Both L2 & L1 Penalise large weights

L2 penalty is proportional to the weights size

L1 is constant penalty

L1 leaves few large connections, encourages sparsity

Bayesian Regularisation

Bayesian belief is to model underlying (multiple) probability distribution

Frequentist belief there is a single true model and data is just realized

Bayesian: probability of hypothesis given data

Frequentist: probability of data given hypothesis

i can assume a distribution
only real data is relevant

$$P(H/D) = \frac{P(H) P(D|H)}{P(D)}$$

Ensemble Learning

Dropout - approximative bagging, drop some nodes randomly

n-fold cross-validation -> group data into n-subsets, train on n-1 samples and use the nth sample for validation, repeat n times and use avg test error

Ensemble learning -> combine multiple networks, assumes networks are independent

↳ Boosting -> Treat previously misclassified samples more seriously

↳ on different weak learners

↓
Decrease
Variance

Bagging reduces variability among weak learners

Network	Algorithm	Properties	Application
Perceptron	Perceptron learning	always converges for linearly separable, premature termination	linear classification
Multilayer Perceptron	Backprop with Delta rule	Local minima vanishing gradient, supervised global classification	
Radial Basis Function	Position RBF with winner takes all, Delta rule train.	Clustering, classification	
Self Organizing Maps	Competitive cooperative learning, with map in output space	neighborhood preserving not distance preserving, vector quantization	
Hopfield Networks	Hebbian learning unsupervised	one-shot learning orthogonal vectors energy based	Denoising
Restricted Boltzmann Machine	Contrastive Divergence, Gibbs sampling	Energy based Stochastic	
Deep Belief Network	greedy layerwise training, sleep-wake fine tuning	Generative, Discriminative	
Convolutional Neural Network	shared weights Convolutional, Pooling, backprop	Translational invariance, weight sharing	Image classification
Generative Adversarial Networks	0-sum game train.	Probability is modeled implicitly	Image generation
Autoencoder	backprop with original as label (not supervised kinda) (MSE)	over/undercomplete sparse base	Compression Denoise
Variational Autoencoder		Stochastic latent space generative vector arithmetic in latent space	

Classification

needs many samples

images

CNN

sigmoidal outputs
(transfer)

input space = size of image

Labeled data
Backprop

regression

reduce image size
when few data samples

generalization loss

Cross entropy

Approximative capabilities
Cluster

Competitive learning
Delta rule
for linear W

RBF

feed forward.
Sigmoid output
for classification

hyperparameter:

- Num RBF
- width (maybe)
- learning rate

Regression

10-fold cross validation

MLP

L2 regularization
dropout

Hopfield

Hebbian Memory
Pattern

sampling without label

Representations
(no images)

SOM

$P(\text{out}|\text{in})$ ← discriminative

Probabilistic (input+output) ← generative
generate data

probabilistic
Label in output

DBN

one-hot encoding

Contrastive divergence

wake-sleep tuning

underlying characteristics

Probabilistic model over data

Gibbs sampling

backprop

MSE Backprop

AE

denoising

Probabilistic

Probabilistic

VAE

Unsupervised

latent variables

underlying characteristics

Probabilistic (input)

GAN

$P(x, y)$

Unsupervised

Underlying characteristics

multiple time scales

LSTM

recurrent

Backprop

Mean sq error

L2 regularization

Temporal prediction
& Regression

Size depends on persistence

Time Delayed fN

MLP

Backprop

Cyclic pattern

Stochastic g.d.

ESN

Train & validation

RNN

Backprop through
time