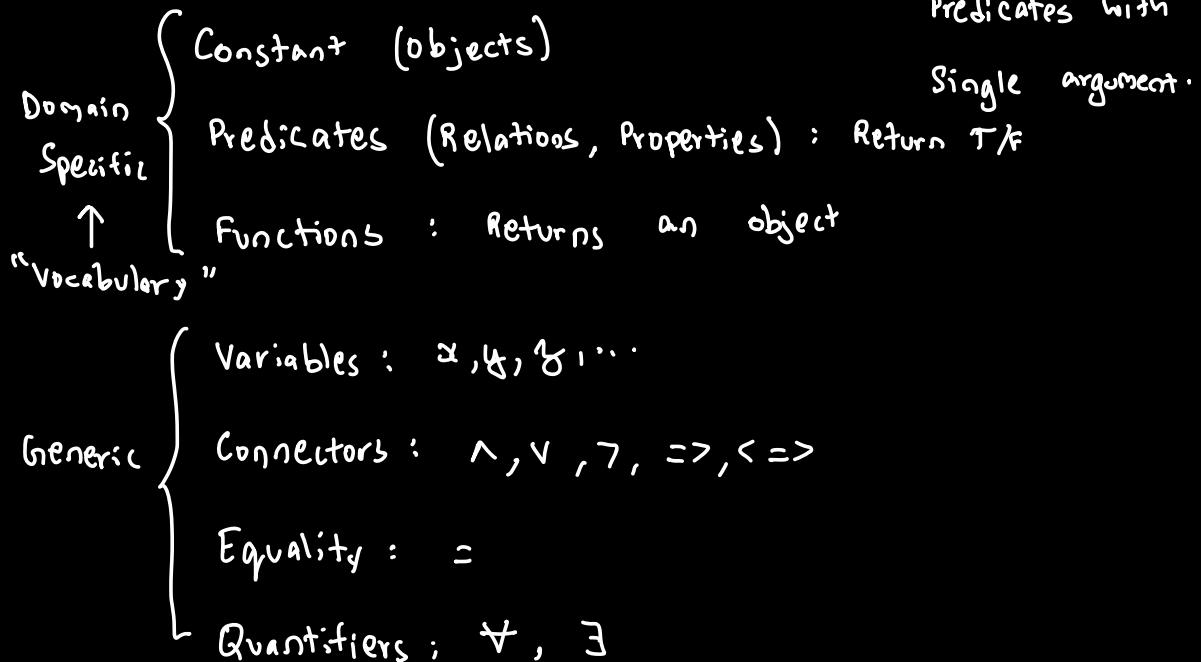


FIRST ORDER LOGIC:

- * more expressive

- * more succinct

SYNTAX:



Atomic Sentence:

Predicate (term₁, ..., term_n)

↳ constant / variable / function (term₁, ..., term_n)

\forall : \Rightarrow

\exists : \wedge

PROPERTIES:

$$\forall x \forall y \equiv \forall y \forall x$$

$$\exists x \exists y \equiv \exists y \exists x$$

$$\exists x \forall y \not\equiv \forall y \exists x$$

$$\neg \exists x \text{ likes}(A, B) \equiv \forall x \neg \text{ likes}(A, B)$$

$\exists ! \text{ king}(x)$: There is exactly one king.

e.g.: Spot has exactly 2 sisters

$$\exists x \exists y \text{ sister}(x, \text{spot}) \wedge \text{ sister}(y, \text{spot}) \wedge \neg(x = y) \wedge$$

$$[\forall z \text{ sister}(z, \text{spot}) \Rightarrow (z = x \vee z = y)]$$

↓
same as $\underbrace{[\neg(\exists z \text{ sister}(z, \text{spot}) \wedge \neg(z = x) \wedge \neg(z = y))]}_{\neg \alpha \vee \beta} \quad \alpha = \neg \beta$

$$\forall z \neg \text{ sister}(z, \text{spot}) \vee ((z = x) \vee (z = y))$$

$$\neg \alpha \vee \beta \quad \alpha = \neg \beta$$

$$\forall z \text{ sister}(z, \text{spot}) \Rightarrow ((z = x) \vee (z = y))$$

Well-formed formula: No free variables
(without quantization)

BINDINGS LIST:

Answers that return True

Query : $\exists x \text{ Mammal}(x)$

{x/spot , x/Tuna}

ENCODE PROBLEM:

1. Vocabulary : Constants, Functions, Predicates
2. KB :
 - i. Independent of case (Domain)
 - ii. Specific case (Instance)
3. Queries

INFERENCE:

i. Reducing FOL inference to propositional inference

ii. Universal Instantiation (UI):

$$\frac{\forall x \ \alpha}{\text{substitute } (\{x/g\}, \alpha)}$$

g: ground term : constants

function (constants)

\Rightarrow If g finite \Rightarrow Equivalence is preserved.

iii. Existential Instantiation (EI):

$$\frac{\exists x \ \alpha}{\text{substitute } (\{x/k\}, \alpha)}$$

k: skolem constant; constant symbol

that does not appear anywhere else

in the KB.

\Rightarrow Satisfiability is preserved. (not equivalence)

Entailment in FOL is semi-decidable:

If α is entailed by FOL KB, it is entailed by finite subset of propositional KB.

For $n=0$ to ∞ do

- * Create propositional KB by instantiating depth- n terms.
- * See if $\text{KB} \models \alpha$

FOL KB $\rightarrow \Delta_0$ (only depth 0, i.e. constants)

Δ Check $\Delta_0 \models \alpha$: Yes $\Delta \models \alpha$

No: Δ_1 (only depth-1 functions + constants)

;

So if $\Delta \models \alpha$, some Δ_i will satisfy $\Delta_i \models \alpha$

if $\Delta \not\models \alpha \Rightarrow$ infinite search.

DEFINITE CLAUSE: Only one positive literal.

2. FORWARD CHAINING:

Start with ground terms and apply the KB,
till we reach the answer.

"Datalog" ; Definite Clause + No function symbols.

3. BACKWARD CHAINING:

Start with result and go back.

→ DFS

→ Logic Programming - Prolog.

4. RESOLUTION:

UNIFICATION:

$$\alpha = \text{Knows}(\text{John}, x)$$

$$\beta = \text{Knows}(\text{John}, \text{Jane})$$

$$\text{Unifier : } \theta = \{x / \text{Jane}\} \quad \text{unify } (\alpha, \beta) = \theta$$

$$\alpha\theta = \beta\theta.$$

$$\Delta \wedge \neg \alpha \rightarrow \text{unsat}$$

$$\text{then } \Delta \models \alpha$$

CONVERSION TO CNF:

1. Eliminate \Rightarrow , \Leftrightarrow

2. Move \neg inwards

$$\neg \forall x \alpha \rightarrow \exists x \neg \alpha$$

$$\neg \exists x \alpha \rightarrow \forall x \neg \alpha$$

3. Standardize variables: each quantifier should use a different one.

4. Skolemize

If \exists is inside $\forall x$

then replace variable with $F(x)$

If \exists is on its own, use some Skolem constant, c_1 .

5. Drop Universal Quantifiers

6. Distribute

$$(\alpha \wedge \beta) \vee \gamma \rightarrow (\alpha \vee \gamma) \wedge (\beta \vee \gamma).$$

PROBABILITY AND BELIEFS:

Classical logic is "monotonic".

→ What is believed, remains to be believed even with new information.

$$\Delta \models \alpha \text{ then } \Delta \wedge \beta \models \alpha$$

Non-classical logic: $\Delta + \text{assumptions}$

"Degrees of belief": α not true/false
but $\text{belief}(\alpha) \in [0,1]$
(PROBABILITY).

PROBABILITY AS A BASIS FOR REPRESENTING BELIEFS:

Propositional logic: $\Delta \models \alpha$

↳ Yes

↳ No

↳ Neither

$$\text{Now: } \Pr(\alpha) = \sum_{\omega \models \alpha} \Pr(\omega)$$

PROPERTIES:

* $\Pr(\alpha \vee \beta) = \Pr(\alpha) + \Pr(\beta) - \Pr(\alpha \wedge \beta)$

* if α, β are mutually exclusive

$$\Pr(\alpha \wedge \beta) = 0$$

* BELIEF CHANGE: given evidence β

$$\Pr(w|\beta) = \begin{cases} 0 & \text{if } w \not\models \beta \\ \frac{\Pr(w)}{\Pr(\beta)} & \text{if } w \models \beta. \end{cases}$$

BAYES CONDITIONING:

$$\Pr(\alpha|\beta) = \frac{\Pr(\alpha \wedge \beta)}{\Pr(\beta)}$$

* INDEPENDENCE:

$\rightarrow \alpha$ is independent of β if

$$* \Pr(\alpha|\beta) = \Pr(\alpha)$$

$$\Rightarrow \frac{\Pr(\alpha, \beta)}{\Pr(\beta)} = \Pr(\alpha) \Rightarrow \Pr(\beta) = \frac{\Pr(\alpha, \beta)}{\Pr(\alpha)}$$

$$\Rightarrow \Pr(\beta|\alpha) = \Pr(\beta).$$

$$\star \Pr(\alpha \cap \beta) = \Pr(\alpha) \cdot \Pr(\beta)$$

$$\star \Pr(\alpha | \beta, \gamma) \text{ need not be } \Pr(\alpha | \gamma)$$

$\rightarrow \alpha$ is conditionally independent of β given γ if

$$\star \Pr(\alpha | \beta \cap \gamma) = \Pr(\alpha | \gamma).$$

$$\star \Pr(\alpha \cap \beta | \gamma) = \Pr(\alpha | \gamma) \cdot \Pr(\beta | \gamma).$$

* **CHAIN RULE:**

$$\Pr(\alpha_1 \cap \alpha_2 \cap \dots \cap \alpha_n) = \Pr(\alpha_1 | \alpha_2 \cap \alpha_3 \cap \dots \cap \alpha_n) \cdot \\ \Pr(\alpha_2 | \alpha_3 \cap \dots \cap \alpha_n) \dots \Pr(\alpha_n).$$

* **CASE ANALYSIS (LAW OF TOTAL PROBABILITY):**

$$\Pr(\alpha) = \sum_{i=1}^n \Pr(\alpha \cap \beta_i) = \sum_{i=1}^n \Pr(\alpha | \beta_i) \Pr(\beta_i)$$

where $\beta_1, \beta_2, \dots, \beta_n$ are mutually exclusive

and exhaustive!

* **BAYES RULE:**

$$\Pr(\alpha | \beta) = \frac{\Pr(\beta | \alpha) \cdot \Pr(\alpha)}{\Pr(\beta)}.$$

$$\text{FALSE POSITIVE} = \Pr(T | \neg D)$$

$$\text{FALSE NEGATIVE} = \Pr(\neg T | D).$$

BAYESIAN NETWORKS:

SYNTAX

↳ * Directed Acyclic Graph (DAG) : Causality

↳ ↳ Probabilities.

VARIABLE INDEPENDENCE:

Variable sets X, Y, Z

X and Y are independent given Z

$$I(X, Z, Y)$$

$$\text{i.e. } \Pr(X|Y, Z) = \Pr(X|Z)$$

$$\text{eg: } X = \{A, B\} \quad Y = \{C\} \quad Z = \{D, E\}$$

4 states

2

4

$$I(X, Z, Y)$$

32
Statements

$$\left\{ \begin{array}{lll} \hookrightarrow A \wedge B \text{ independent of } C \text{ given } D \wedge E \\ A \wedge \neg B \quad \dots & C \dots & D \wedge E \\ \vdots & & \\ \neg A \wedge \neg B \quad \dots & \neg C \dots & \neg D \wedge \neg E \end{array} \right.$$

CAPTURING INDEPENDENCE GRAPHICALLY:

Non-Descendants(v) : All nodes \sim Descendants(v)
 $- v -$ Parents(v)

MARCOVIAN ASSUMPTIONS OF A BAYESIAN NETWORK:

- * $I(v, \text{Parents}(v), \text{Non-Descendants}(v))$
- Markov(b) \rightarrow Apply this for all nodes(v)
 in b .

PARAMETERIZING THE STRUCTURE:

CONDITIONAL PROBABILITY TABLE (CPT):

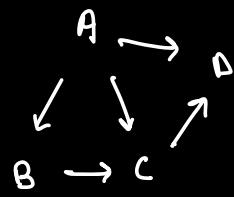
CPT for variable c : $\Pr(c | \text{Parents}(c))$

If c has 2 parents A, B

A	B	C	$\Pr(c a,b)$
			$\left. \begin{array}{l} 8 \text{ parameters} \\ (\text{if independent}) \end{array} \right\}$

$$\theta_{\bar{c}|\bar{a}\bar{b}} = \Pr(\bar{c}|\bar{a}, \bar{b}) = \Pr(c = \text{false} | A = \text{false}, B = \text{true}).$$

$$\Pr(a\bar{b}c\bar{d}) = \Pr(a) \cdot \Pr(\bar{b}|a) \cdot \Pr(c|a\bar{b}) \cdot \Pr(\bar{d}|ac)$$



d- SEPARATION:

We know z

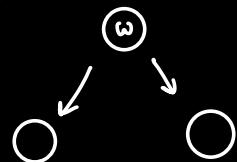
SEQUENTIAL VALVE:



Closed iff

$$w \in z$$

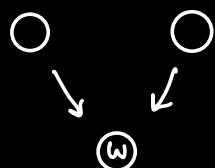
DIVERGENT VALVE:



Closed iff

$$w \in z$$

CONVERGENT VALVE:



Closed iff

neither w or any of its descendants are
in z .

x and y are d -separated given z

$$dsep(x, z, y)$$

iff every path between x, y is blocked by z .



Path is blocked by z iff at least one value on the path is closed given z .

INFERENCE:

- * Queries
- * Algorithms

QUERIES:

1. PRIOR MARGINALS:

Probability of variable without any condition.

$$\Pr(c) \rightarrow c = \text{yes} \quad 3.2\%$$

$$c = \text{no} \quad 96.8\%$$

note: Deterministic $\rightarrow \Pr(x) = 0 / 1$.

2. POSTERIOR MARGINALS:

$$\Pr(c | \beta)$$

3. MPE : Most Probable Explanation:

Say 5 variables

Evidence : One variable given

The most probable instantiation of other variables given evidences.

4. MAP : Maximum a Posteriori Hypothesis:

Most probable instantiation of a subset of variables given evidences.

MPE can be done using MAP.

COMPLEXITY OF INFERENCE:

- Variable elimination
- conditioning

Complexity depends on

↳ Topology

↳ Treewidth (connectedness)

n : #variables

d : #values

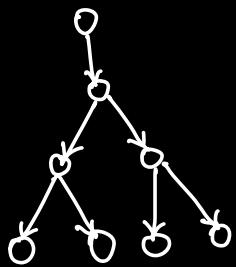
Marginals computed in

w : treewidth

$O(n \cdot d^w)$

TOPOLOGY:

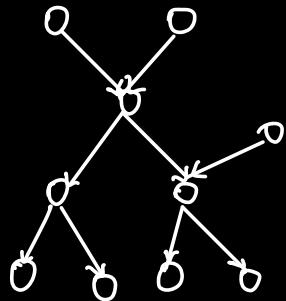
* TREE :



* only 1 parent

* $w = 1$

* POLYTREE: (Semi - CONNECTED NETWORK)



* $w = k$

↳ maximum number
of parents a
node can have.

* Only single path between
2 nodes.

* MULTIPLY - CONNECTED (DMC).

WEIGHTED MODEL COUNTING (WMC):

#SAT: model counting

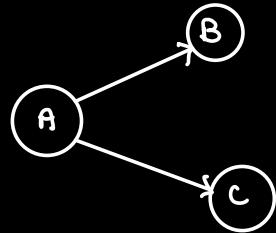
$$WMC = \sum_{w \models \Delta} w \cdot \text{wt}(w) \quad \text{if } \Delta = (A \vee B) \wedge C$$

↓
weight of w

Idea:

$\Delta \xrightarrow{\text{Compile}} d\text{-DNNF Circuit} \longrightarrow WMC \text{ in linear time.}$
+ Smooth

REDUCING PROBABILISTIC INFERENCE TO WNC:



Probabilities: $\theta_A, \theta_{\bar{A}}$

$\theta_{B|A}, \theta_{\bar{B}|A}, \theta_{B|\bar{A}}, \theta_{\bar{B}|\bar{A}}$

$\theta_{C|A}, \theta_{\bar{C}|A}, \theta_{C|\bar{A}}, \theta_{\bar{C}|\bar{A}}$

Step 1 : Take each node in the Bayesian network as a variable. Create a variable for each probability.

A, B, C , $\gamma_A, \gamma_B, \gamma_C \rightarrow$ "Indicator" variables
(A,B,C)

P₁ θ_A

P₂ $\theta_{\bar{A}}$

P₃ $\theta_{B|A}$

P₇ $\theta_{C|A}$

P₄ $\theta_{\bar{B}|A}$

P₈ $\theta_{\bar{C}|A}$

P₅ $\theta_{B|\bar{A}}$

P₉ $\theta_{C|\bar{A}}$

P₆ $\theta_{\bar{B}|\bar{A}}$

P₁₀ $\theta_{\bar{C}|\bar{A}}$

↓ ↓
"Parameter" Variables (P_i)

Step 2 : Create KB

$$\Delta \vdash A \Leftrightarrow P_1,$$

$$\neg A \Leftrightarrow P_2$$

$$A \wedge B \Leftrightarrow P_3$$

$$A \wedge \neg B \Leftrightarrow P_4$$

$$\neg A \wedge B \Leftrightarrow P_5$$

$$\neg A \wedge \neg B \Leftrightarrow P_6$$

$$A \wedge C \Leftrightarrow P_7$$

$$A \wedge \neg C \Leftrightarrow P_8$$

$$\neg A \wedge C \Leftrightarrow P_9$$

$$\neg A \wedge \neg C \Leftrightarrow P_{10}$$

Step 3 : Assign weights

$$w(A) = w(\neg A) = w(B) = w(\neg B) = w(C) = w(\neg C) = 1$$

$$w(P_i) = \theta_i \quad \text{i.e. } w(P_1) = \theta_A, \quad w(P_7) = \theta_{C1A}$$

$$w(\neg P_i) = 1$$

e.g. One Model of Δ :

$$A = T, B = F, C = T \Rightarrow A \wedge B \subset P_1, P_4, P_7, \neg P_2, \dots, \neg P_{10}$$

$$w(\text{model}) = 1 \times 1 \times 1 \times \theta_A \times \theta_{\neg B1A} \times \theta_{C1A} \times 1 \dots \times 1$$

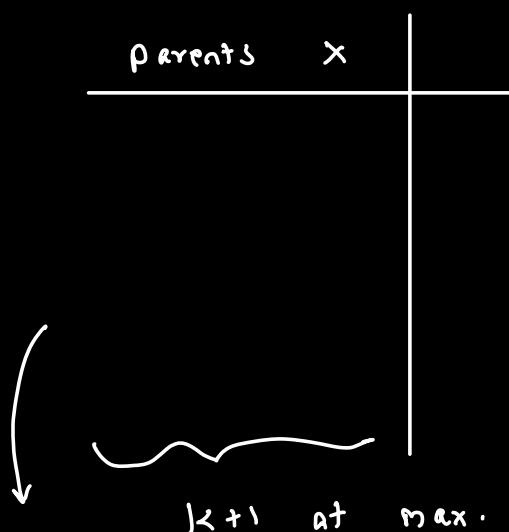
Step 4: $\Delta \xrightarrow{\text{Compose}} d\text{-DNNF Circuit} \longrightarrow \text{WMC in linear time.}$
+ Smooth

Step 5: Given Query α

$$\Pr(\alpha) = \text{WMC}(\Delta \wedge \alpha).$$

ADVANTAGES: Not dependent on treewidth!

SIZE OF BAYESIAN NETWORK:



n : variables

k : maximum # parents
per node.

d : maximum # values
per variable

Size of CPT $\leq 2^{k+1}$ if 2 value per variable

$$\leq d^{k+1}$$

There are n CPTs

$$\text{Size of BN: } O(n \cdot d^{k+1}). \quad \ll d^n$$

Joint Probability
Table Size:

MODELLING BAYESIAN NETWORKS:

Step 1:
Variable, Values

Step 2:
Edge

Step 3:
CPT
→ Problem Statement
→ Subjective Beliefs
→ Learning from data.

DATA:

Complete: All data present
Closed Form (Efficient)

Incomplete: Some data missing.

EM: Expectation Maximization.

Say we have 2 CPTs - A, B
which to choose.

Run the data through and get $\Pr(\text{Data})$

$\hookrightarrow p_1, p_2 \dots p_n \quad P_i = \Pr(\text{example } i)$

\hookrightarrow this is the score

MAXIMUM LIKELIHOOD PRINCIPLE : $\arg \max_{\text{CPT}_S} \Pr(\text{Data} | \text{CPT}_S)$

SENSITIVITY ANALYSIS:

Given we want posterior marginal to be some value, what changes can help achieve this.

Evidence : E

$$\Pr(R|E) = 10\%$$

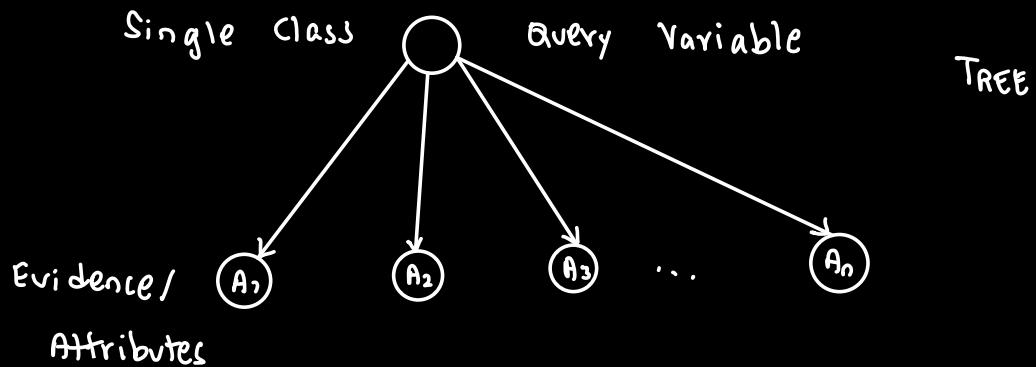
$$\text{We want } \Pr(R|E) \leq 5\%.$$

Run Sensitivity Analysis

e.g. → False Negative for x

reduce from 10% to 4%.

NAIVE BAYES STRUCTURE:



Here we assume final query variable is a

single value. e.g.: Disease query can be

Flu or Cold not both.

SINGLE FAULT ASSUMPTION!

LEARNING:

- * Parameters

- * Structure

PARAMETER LEARNING:

Structure given.

1. Complete Dataset:

1. Take the dataset \Rightarrow generate empirical distribution (probability of each world).

A	B	C	$\Pr_D(\cdot)$	Dataset
T	T	T	2/16	
			:	

F F F 1/16

2. Compute Probabilities:

$$\theta_{\bar{a}|b} = \Pr(\bar{a}|b) = \frac{\Pr_D(\bar{a}, b)}{\Pr_D(b)}$$

These are the

MLE of them.

One ITERATION!

2. INCOMPLETE DATASET:

EM [EXPECTATION MAXIMIZATION]

Generate random CPT_1 .

* $CPT_1 \rightarrow BN_1 \rightarrow \Pr_1(\cdot)$

Replace incomplete eg: $\begin{array}{c} A \\ T \\ F \end{array} \quad \begin{array}{c} B \\ F \\ ? \end{array} \quad \begin{array}{c} C \\ ? \end{array}$

entries

$$\begin{array}{ccc} x & \tau & f & \tau \\ \diagup & & & \diagdown \\ y & \tau & f & f \end{array}$$

$\begin{array}{ccc} \tau & f & ? \end{array}$

\downarrow

$$x = \Pr_1(C=\tau | A=\tau, B=f)$$

$\begin{array}{cccc} \tau & f & \tau & x \end{array}$

$$y = \Pr_1(C=f | A=\tau, B=f)$$

$\begin{array}{cccc} \tau & f & f & y \end{array}$

Now: $\begin{array}{ccccc} \tau & \tau & \tau & | & \\ \tau & \tau & f & , & \end{array} \quad \left. \begin{array}{c} \\ \\ \end{array} \right\} \text{complete}$

$\begin{array}{cccc} \tau & f & \tau & x \end{array}$

$\begin{array}{cccc} \tau & f & f & y \end{array}$

* Use this to get CPT_2

* $CPT_2 \rightarrow BN_2 \rightarrow \Pr_2(\cdot)$

:

converges.

[Local Search Algorithm]

$CPT_1 \dots \xrightarrow{\text{Inference}} CPT_2 \xrightarrow{\text{Inference}} \dots$

If n loops $\Rightarrow n$ inferences).

LEARNING STRUCTURE:

1. LOCAL SEARCH METHODS:

- * Choose one structure.
- * Change and move to a neighboring structure with better score.
 - add, remove, reverse an edge

2. SYSTEMATIC SEARCH METHODS:

A* search.

WHAT IF WE USE MLE?

We end up with complete structure.

Overfitting. More parameters → Higher likelihood.

↪ Won't generalize.

Score : Likelihood + Penalty Term

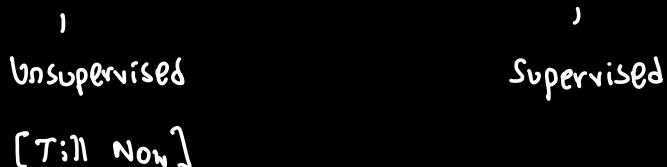
e.g.: MDL Score

↪ number of independent parameters

and

size of data

Model - ORIENTED VS QUERY - ORIENTED LEARNING:



SUPERVISED LEARNING:

$B_N + Q_{\text{query}} \Rightarrow \text{Boolean Formula} + \begin{matrix} \uparrow \\ \text{weights} \end{matrix} \approx NNF \text{ Circuit}$

compilation time
 $O(n \cdot d^{\omega})$

↓
 Convert NNF
 Circuit to AC

Parameters:

* $\theta_i = \Pr(i)$

For evidence: ↘ Indicators

$$A \quad \lambda_a \quad \lambda_{\bar{a}}$$

$$T \quad 1 \quad 0$$

$$F \quad 0 \quad 1$$

$$? \quad 1 \quad 1$$

LOSS FUNCTION:

For each input, find $\Pr(B)$ using AC, say P

True Output is :
 (Label)

B
 $x \quad 0$
 $y \quad 1$
 $z \quad 0$ → one-hot representation.
 Say Q

$$\text{Loss} = \frac{1}{n} \leq \text{Cross-Entropy}(P(x), Q(x))$$

Optimized using Gradient Descent.

CROSS-ENTROPY:

$P(x)$: Prediction

$Q(x)$: Label

$$CE = \sum_x Q(x) \log_2 P(x)$$

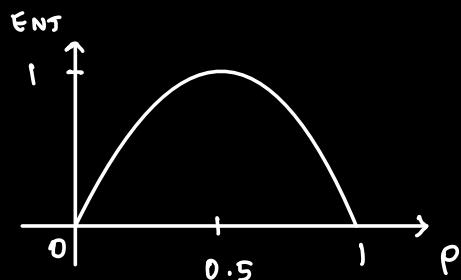
BACKGROUND KNOWLEDGE:

Can be added directly to BN.

DECISION TREES AND RANDOM FORESTS:

ENTROPY:

$$H(x) = - \sum_x p(x) \log_2 p(x)$$



$$H(x|y) = - \sum_x p(x|y) \log_2 p(x|y)$$

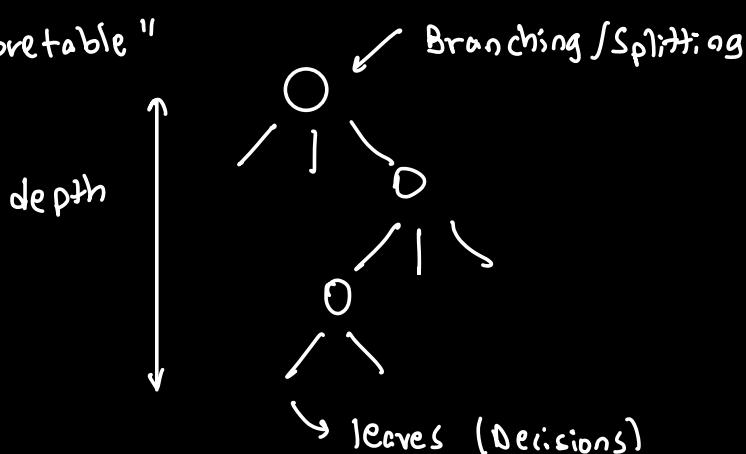
$$H(x|y) \leq \sum_y p(y) H(x|y)$$

$$H(x|y) \leq H(x)$$

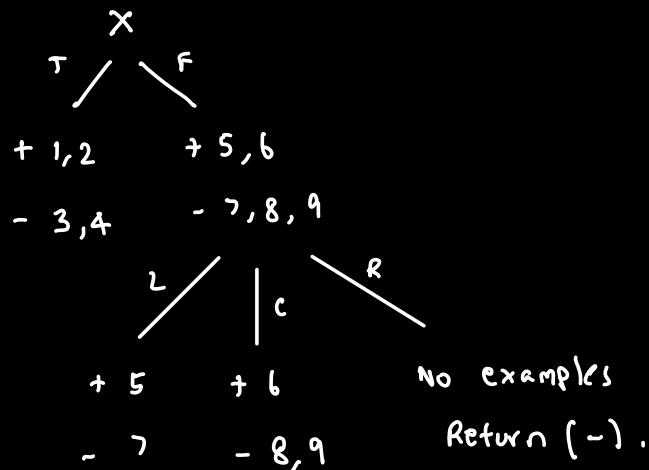
when independent

$$H(x|y) = H(x)$$

"Interpretable"



- * If no attributes left, return the PLURALITY of left example (majority vote).
- * No examples left, return majority of parent examples.

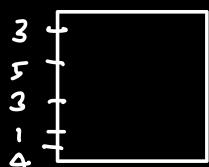


"Interpretable" : Decision Trees > Bayesian Networks > Neural Networks.

RANDOM FORESTS :

1. Ensemble learning methods: Majority vote from multiple Decision trees.
2. Bootstrapped Dataset: Randomly pick with replacement.

So, each Decision tree has different Dataset.

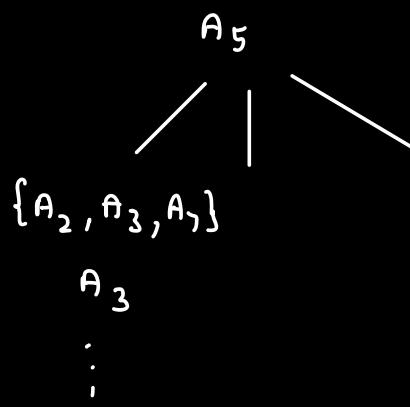


2 is "Out-of-bag" example.

3. Only a subset of attributes is considered at each point. At each point choose random 'm' attributes and select best of the m. Then next Step choose random m again.

$$n = 10 \text{ attributes} \quad \{A_1, A_5, A_7\}$$

$$m = 3$$

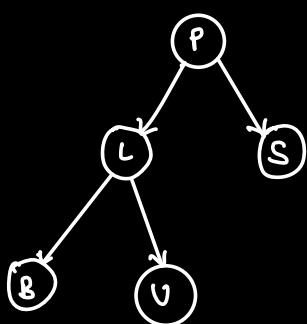


$$m \text{ can be } \sqrt{n}$$

4. Accuracy tested on "Out-of-bag" examples.

BAYESIAN NETWORK CLASSIFIER:

Threshold, T



Instance

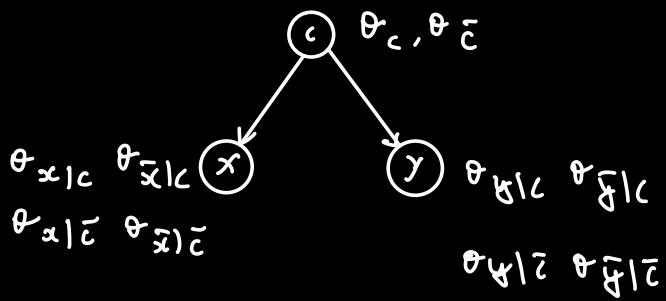
$$S = +, B = -, U = -$$

$$\Pr(P | S=+, B=-, U=-) \geq T$$

↳ Return Pregnant

else Not Pregnant.

Naïve Bayes Classifier:



$$\Pr(c|x,y) = \frac{\Pr(x,y|c) \Pr(c)}{\Pr(x,y)}$$

↳ $\Pr(x,y|c) \cdot \Pr(c) + \Pr(x,y|\bar{c}) \cdot \Pr(\bar{c})$

$$= \frac{\theta_{x,y|c} \cdot \theta_c}{\theta_{x,y|c} \cdot \theta_c + \theta_{x,y|\bar{c}} \cdot \theta_{\bar{c}}}$$

\nwarrow_{CPT}
 \nearrow_{CPT}

In Naive Bayes

x, y are independent given C

$$\left. \begin{aligned} \text{So } \theta_{x,y|C} &= \theta_{x|C} \cdot \theta_{y|C} \\ \theta_{x,y|\bar{C}} &= \theta_{x|\bar{C}} \cdot \theta_{y|\bar{C}} \end{aligned} \right\} \rightarrow \text{CPT}$$

So, output Probability \rightarrow directly from CPT!

Also, Time complexity for inference: $O(n \cdot d^\omega)$

here $\omega = 1$

So $O(n \cdot d) //$

EXPLAINABLE AI:

Bayesian Network \rightarrow Compiles into tractable circuit

\downarrow
(NNF) with same input-output
behavior.

Numerical

Symbolic

This circuit can help explain
why we got an output-

"PI-Explanation": e.g.: $s=1, b=0, f=1, m=1$ is true

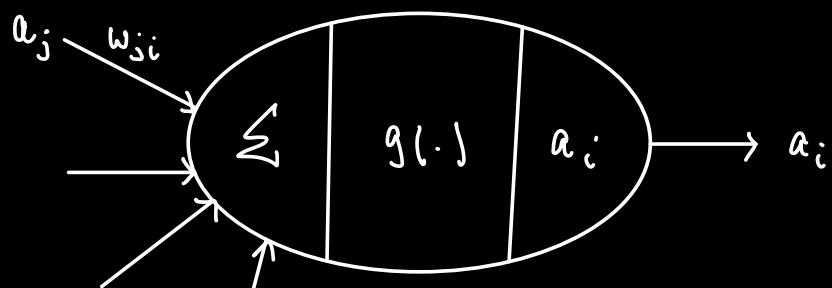
but we see all $F=M=1$ is true

So, only F, M contributed

NEURAL NETWORKS:

- Model-free approach.

NEURONS:



$$i_{in,i} = \left(\sum_j w_{ji} a_j \right) + b$$

$$a_i = g(i_{in,i})$$

a_i, a_j : activations

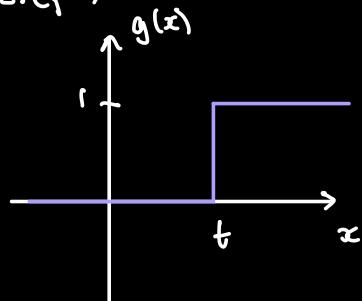
w_{ji} : weights

b : bias

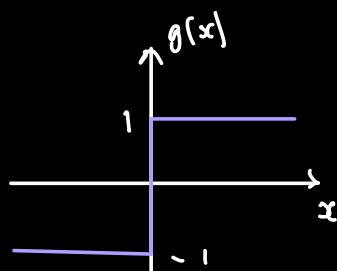
g : activation function

ACTIVATION FUNCTIONS:

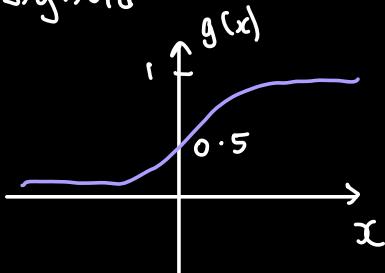
Step :



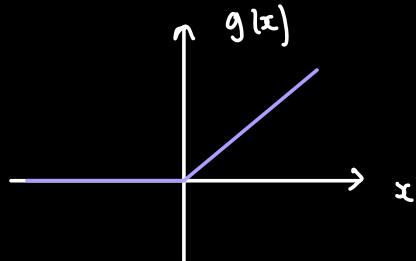
Sigmo :



Sigmoid:



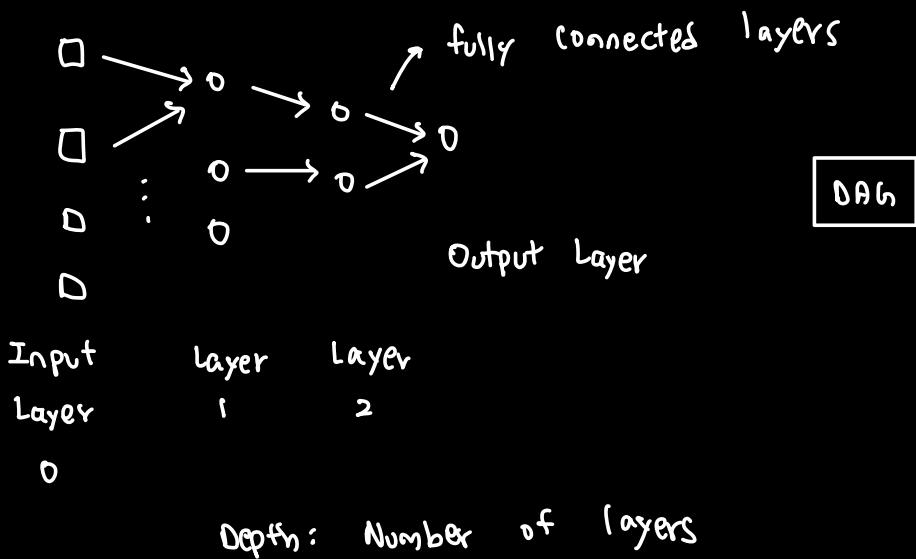
ReLU (Rectified Linear Unit):



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

(Probabilistic)

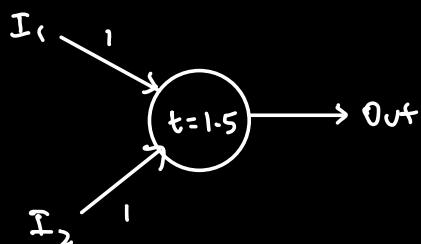
FEEDFORWARD NN:



NN are "Universal Function Approximators".

NEURONS WITH STEP ACTIVATION FUNCTIONS:

→ AND



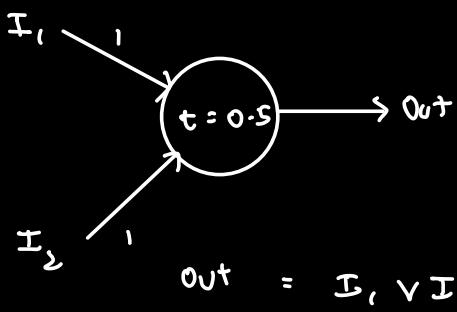
$$Out = I_1 \wedge I_2$$

So Out is true only when $I_1 = I_2 = \text{true}$

$$\text{So if } w_1 = 1, w_2 = 1, \text{ out} = 2$$

$$t = 1.5$$

→ OR

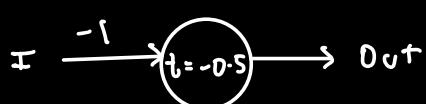


Out is false only when $I_1 = I_2 = \text{false}$

$$\text{If } w_1 = 1, w_2 = 1, \text{ out} = 0.$$

$$t = 0.5$$

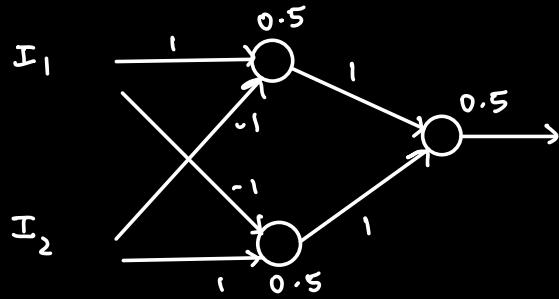
→ NOT



\rightarrow XOR :

$$I_1 \oplus I_2$$

$$(I_1 \wedge \neg I_2) \vee (I_2 \wedge \neg I_1)$$



$$x = I_1 \wedge \neg I_2$$

↓

If $w_1 = 1 \quad w_2 = -1$

$I_1 \wedge \neg I_2$ is true only when

$I_1 = \text{true}, I_2 = \text{false}$

which is $1 \times 1 + (-1 \times 0) = 1$

$$\text{So } t = 0.5$$

$$y = I_2 \wedge \neg I_1$$

↓

$w_1 = -1 \quad w_2 = 1$

$$(-1 \times 0 + 1 \times 1) = 1$$

$$t = 0.5$$

$$\text{out} = x \vee y$$

Neuron can only solve "linearly separable functions".

Not XOR.

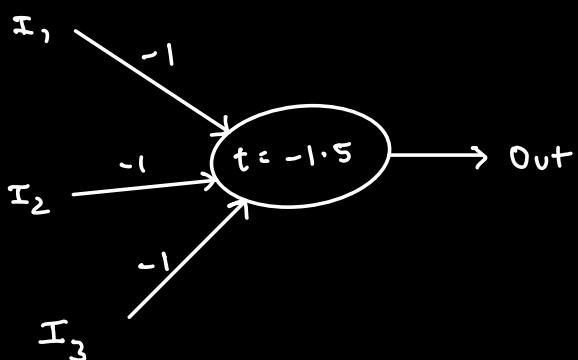
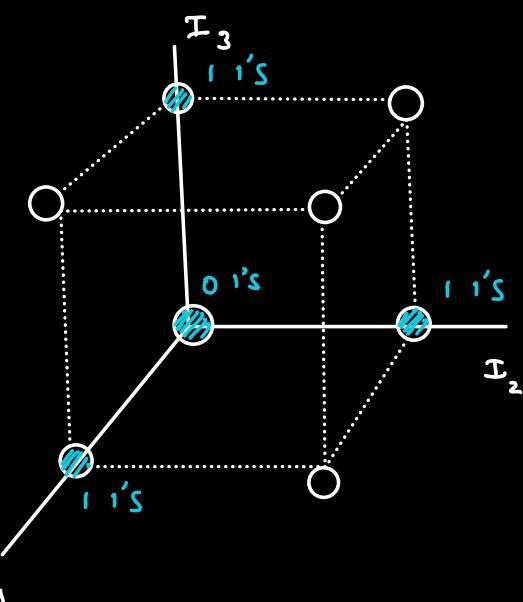
MINORITY :

$$001 \rightarrow 1$$

Output = 1 if 1 is

$$011 \rightarrow 0$$

minority.



$$0 \quad -1 \quad -1 \quad -1 \\ 000, 001, 010, 100 \rightarrow 1$$

$$011, 110, 101, 111 \rightarrow 0$$

$$-2 \quad -2 \quad -2 \quad -3 \qquad \geq -1.5$$

NEURAL NETWORK AS A FUNCTION:

$$\text{out} = f \left(\underbrace{\text{in}_1, \text{in}_2}_{\text{input}}, \underbrace{\omega_{13}, \omega_{14}, \dots, \omega_{45}}_{\text{weight}} \right)$$

[Complex non-linear function]

of inputs and weights.

Given an input $I_i : \text{in}_1, \text{in}_2$

$$\text{out}_i = f_i(\text{weights})$$

TRAINING NEURAL NETWORKS:

Loss function: [Optimize]

Cross Entropy (CE), Mean-Square Error (MSE).

	in_1	in_2	out	Labels (l_i)	$\text{NN}(I_i)$
I_1					
I_2					
:					
I_N					

$$\text{MSE} : \frac{1}{n} \sum_{i=1}^n \left(\text{NN}(I_i) - l_i \right)^2$$

$f_i(\omega_1, \dots, \omega_k)$

GRADIENT DESCENT:

$$f = f(w_1, w_2, \dots, w_k)$$

- loss function
- * Initialize random values
 - * gradient: $(\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \dots, \frac{\partial f}{\partial w_k})$
 - * Update (Step: η)

OPTIMIZATION: loss function [Training Data]

PERFORMANCE METRIC: Accuracy [Testing Data
+ Validation Data]

Epoch: One pass through data.

Batch: Divide training data into batches.
↓
Epoch 1 : Divide Training data into batches (eg: Batch size:
randomly) b examples

Gradient Descent over each batch

↓
Output: Loss 1, Metric 1

Epoch 2 : Divide ...

↓
Gradient Descent ...

↓
STOP

STOPPING CRITERION:

1. Epoch Count

2. Till Loss Decreases [Overfitting]

3. VALIDATION: Training \leq Training 80%
validation 20%

After every epoch

↓ metric on validation
data.

Stop when accuracy (metric) stops

increasing / stagnant.

PATIENCE FACTOR: No improvement over
some x epochs - STOP.

Epoch 1 loss 1 metric 1

Epoch 2 loss 2 metric 2

:

Epoch K Stop.

[SGD: Stochastic Gradient Descent].