

CS50: Introduction to IR with Python

Introduction

Plan

- Search
- knowledge
- Uncertainty
- Optimization
- Machine Learning
- Neural Networks
- Language.

Chapter 1: Search

Def : An **agent** is an entity that perceives its environment and acts upon that environment

Def : A **state** : a configuration of the agent and its environment

Def : An **action** : choices that can be made in a state.
functions : Actions(s) → output : set of actions that can
state input be executed

Def : Transition model : a description of what states result from performing any applicable action in any state
function : Results(s, a) : → output : state resulting from performing a in state s .

Def : State Space : the set of all states reachable from the initial state by any sequence of actions

Representation :



• nodes (state)
→ actions

Def. goal test: way to determine whether a given state is a goal state

Def. path cost: numerical cost associated with a given cost

=> goal : find goal state with low path cost.
=> optimal solution

=> **Search Problem definition**

- init state
- actions
- transit model
- goal test
- path cost function

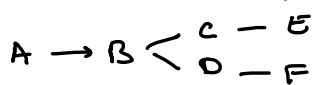
Def. a node: data structure that keep tracks of:

- a state
- a parent (previous node)
- an action (applied to parent)
- a path cost (from init state to node)

Approach:

- start with a frontier that contains the initial state.
- repeat:
 - if frontier is empty, then no solution
 - Remove a node from the frontier
 - if node contains goal state, return the solution
 - Expand node, add resulting nodes to the frontier

Ex: Find path from A to E.



Start init: Frontier = A.

- Remove A (Not goal)
- Expend

Step 1: Frontier = B

- Remove B (Not goal)
- Expend

Step 2: Frontier = C, D

- Remove D (Not goal) (arbitrarily)
- Expend

Step 3: Frontier = C, F

- Remove C
- Expend

Step 4: Frontier = E

- Remove E: find solution

A → B → C → E

▷ If we can go back to a state it can create infinite loop between 2 states.

⇒ Revised Approach

- Start frontier that contains the initial state
- Start with an empty explored set
- Repeat
 - If frontier is empty, then no solution
 - Remove a node from frontier
 - if node contains a goal state, return solution
 - Add node to explored set
 - Expend Node, add resulting nodes to the frontier if they aren't already in the frontier or the explore set

▷ How we remove a node is super important.

a. Uninformed Search

Def : **Stack** : Last-in, first-out data type

Def **Depth-first search** : search algo that always expand the deepest node in the frontier

↳ what happen if we use stack as frontier, we will go as deep as possible in a branch and change only if we reach the end.

Def **Breadth-first search** : search algo that always expand the shallowest node in the frontier

↳ use opposite of stack as frontier

Def **queue** : first in, first out data type (opposite of stack)

⇒ If final state both algo will always find a solution.

However, Breadth first search will always find the best solution (\neq Depth-first algo)

↑ Need to explore a lot of states to find the solution

Def **Uninformed search** : search strategy that uses no problem-specific knowledge.

(Ex BFS / DFS)

Def **Informed search** : search strategy that uses problem-specific knowledge to find a solution more efficiently.

b. Informed Search

Def : **Greedy best-first search**: search algo that expands the node that is closer to the goal, as estimated by a heuristic function $h(n)$

Def **A* search**: search algo that expands node with lowest value of $g(n) + h(n)$

$g(n)$: cost to reach node

$h(n)$: estimated cost to goal.

→ GBFS do not always find optimal solutio
A* search optimal if:

- $h(n)$ is **admissible** (never over estimate the true cost)
- $h(n)$ is **consistent** (for every node n and successor n' with step cost c , $h(n) \leq h(n') + c$)

c. Adversarial approach.

Adversarial search: someone try to make you fail.

ex: Numpic. (X or O)

X wins = 1, O wins = -1, no one = 0

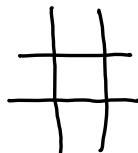
Def **Minimax**: give a value to each output possible.

- Max (X) aims to max score
- Min (O) aims to min score.

Game:

- S_0 : initial player
- Player (s): returns which player to move in state s
- Actions : returns legal moves in state s
- Results (s, a): return state after applying a in state s
- terminal (s): check if state s is a terminal state
- Utility (s): final numerical value for terminal state ..

Ex: Init



$$\text{player}(\#) = X$$

(X starts)

$$\text{player}(X) = 0$$

$$\text{action}\left(\begin{array}{|c|c|c|} \hline X & 0 & X \\ \hline 0 & 0 & \\ \hline X & 0 & 0 \\ \hline \end{array}\right) = \left\{ \begin{array}{l} \begin{array}{|c|c|c|} \hline X & 0 & X \\ \hline 0 & X & \\ \hline X & X & 0 \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline X & 0 & X \\ \hline 0 & 0 & X \\ \hline X & X & 0 \\ \hline \end{array} \end{array} \right\}$$

$$\text{Result}(X | \begin{array}{|c|c|c|} \hline X & & \\ \hline 0 & & \\ \hline & & \# \\ \hline \end{array}, \#) = \begin{array}{|c|c|c|} \hline X & 0 & \\ \hline \# & 0 & \\ \hline & & \# \\ \hline \end{array}$$

$$\text{Terminal}\left(\begin{array}{|c|c|c|} \hline 0 & 0 & X \\ \hline X & X & \\ \hline 0 & 0 & \\ \hline \end{array}\right) = \text{True}.$$

MinMax Algo:

- given a state s :
- Max picks action in $\text{Actions}(s)$ that produces highest value of Min-Value($\text{Result}(s, a)$)
- Min picks action in $\text{Actions}(s)$ that produces smallest value of Max-value($\text{Result}(s, a)$).

function Max-Value(state):

if Terminal(state):
return Utility(state).

$V = -\infty$

for action in Actions(state):

$v = \text{Max}(v, \text{Min-Value}(\text{Result}(state, action)))$

return v .

function Min-Value(state):

if Terminal(state):
return Utility(state).

$V = +\infty$

for action in Actions(state):

$v = \text{Min}(v, \text{Max-Value}(\text{Result}(state, action)))$

return v .

→ Optimization

As we maximize (or minimize) between all possibility it is possible that we don't have the best but we maximize the min value we will have.

Depth Limited Min Max:

Same but stop after n iteration even if game not stop

Need evaluation function instead of utility to find who is likely to win.
(ex: in chess: 10^{29000} possibility so it is a way to reduce calculation)

Chapter 2 : Knowledge

Def : knowledge-based agents: agent that reason by operating on internal represent' knowledge.

a. Logic

Def : sentence : assert about the world in a knowledge representation language.

⇒ Propositional Logic

→ Propositional symbol :

P Q R

→ logical connective

¬: not ∧: and ∨: or →: implicit ↔: biconditional

Ex: P = true : $\neg P$: false

P	Q	$P \wedge Q$	$P \vee Q$	$P \rightarrow Q$	$P \leftrightarrow Q$
false	false	false	false	true	true
false	true	false	true	true	false
true	false	false	true	false	false
true	true	true	true	true	true

Def: model : assignment of a truth value to every propositional symbol ("possible world")

ex : P : it is raining model : $\{P = \text{true}, Q = \text{false}\}$
Q : it is Tuesday

Def: knowledge base: a set of sentences known by a knowledge-base agent.

Def: Entailment: $\alpha \models \beta$. (read " α entails β ")

In every model in which sentence α is true, sentence β is also true

Ex. P : it is a Tuesday in February

$P \models Q$ Q : it is February

Def: inference: the process of deriving new sentences from old ones.

Ex : P : it is Tuesday

Q : it is raining

R : Harry will go for a run.

knowledge base

KB : • $(P \wedge \neg Q) \rightarrow R$: If it is Tuesday and not raining, Harry will go for run
• P
• $\neg Q$

Inference $\rightarrow R$ is true.

Algo goal: $KB \models \alpha$?

b. Model Checking algo.

Algo

- To determine if $KB \models \alpha$:
 - enumerate all possible models
 - if in every model where KB is true, α is true then KB entails α .
- otherwise, KB does not entail α .

Ex: Same P, Q and R.

$$KB: (P \wedge \neg Q) \rightarrow R, P, Q$$

Query: R (what we are looking)

P	Q	R	KB
false	false	false	false
false	false	true	false
false	true	false	false
false	true	true	false
true	false	false	false
true	false	true	true
true	true	false	false
true	true	true	false

Because P is True
Because $(P \wedge \neg Q) \rightarrow R$
Because $\neg Q$

\Rightarrow One possible world: in this world R is true
 \Rightarrow R is true. "where KB true"

c. Knowledge Engineering.

Create algo to solve sol with knowledge.

Ex: Cluedo, Logic Puzzles.

▷ Model checking is not very efficient (2^N N: variables)

d. Inference Rules

Modus Ponens:

if $\alpha \rightarrow \beta$ and α then β

And elimination

if $\alpha \wedge \beta$ then α

Double negation elimination

if $\neg(\neg \alpha)$ then α

Implication Elimination

if $\alpha \rightarrow \beta$ then $\neg \alpha \vee \beta$

Biconditional Elimination

if $\alpha \leftrightarrow \beta$ then $(\alpha \rightarrow \beta) \wedge (\beta \rightarrow \alpha)$

De Morgan's Law

if $\neg(\alpha \wedge \beta)$ then $\neg \alpha \vee \neg \beta$

if $\neg(\alpha \vee \beta)$ then $\neg \alpha \wedge \neg \beta$

Distributive property

if $(\alpha \wedge (\beta \vee \gamma))$ then $(\alpha \wedge \beta) \vee (\alpha \wedge \gamma)$
if $(\alpha \vee (\beta \wedge \gamma))$ then $(\alpha \vee \beta) \wedge (\alpha \vee \gamma)$

d. Theorem proving.

Theorem Proving:

- Initial state : starting knowledge base (KB)
- actions : inference rules
- transit. models : new KB after inference.
- goal test : check statement we are trying to prove
- path cost : number of steps in proof.

rule : if $(P \vee Q) \wedge \neg P$ then Q

generalise : if $(P \vee Q_1 \vee \dots \vee Q_n) \wedge \neg P$ then $(Q_1 \vee \dots \vee$

rule : if $(P \vee Q)$ and $(\neg P \vee R)$ then $Q \vee R$

generalise if $(P \vee Q_1 \vee \dots \vee Q_n)$ and $(\neg P \vee R_1 \vee \dots \vee R_n)$
then $Q_1 \vee \dots \vee Q_n \vee R_1 \vee \dots \vee R_n$

Def clause: a disjunction of literals (ex : $P \vee q \vee \neg r$)

Def Conjunctive normal form: logical sentence that's a conjunction of clauses. (ex: $\underbrace{(A \vee B \vee C)}_{\text{clause 1}} \wedge \underbrace{(\neg D \vee \neg E)}_{\text{clause 2}} \wedge \underbrace{(F \vee G)}_{\text{clause 3}}$)

Every logical sentence can have this form, with inference rule.

Conversion to CNF (conjunctive normal form)

- Eliminate biconditionals:
 - turn $(\alpha \leftrightarrow \beta)$ into $(\alpha \rightarrow \beta) \wedge (\beta \rightarrow \alpha)$
- Eliminate implication
 - turn $(\alpha \rightarrow \beta)$ into $\neg \alpha \vee \beta$
- Remove \neg inwards using De Morgan's Laws
 - ex turn $\neg(\alpha \wedge \beta)$ into $\neg \alpha \vee \neg \beta$
- use distributive law to distribute \vee wherever is possible.

Ex: $(P \vee Q) \rightarrow R$

$$\neg(P \vee Q) \vee R$$

eliminate implicant.

$$(\neg P \wedge \neg Q) \vee R$$

De Morgan's Law

$$(\neg P \vee R) \wedge (\neg Q \vee R)$$

Distributive Law

e. Inference by Resolution

Special Case:

$$\text{if } (P \wedge Q \vee S) \wedge (\neg P \vee R \vee S)$$

$$\text{then } (Q \vee S \vee R \vee S)$$

→ need to eliminate duplicate variable
⇒ $(Q \vee R \vee S)$

Special case: if $P \wedge \neg P$ then () (equivalent False)

Inference by resolution / Principle.

- To determine $\text{KB} \models \alpha$
- check if $(\text{KB} \wedge \neg \alpha)$ is a contradiction
 - if so, then $\text{KB} \models \alpha$
 - otherwise, no entailment.

Inference by resolution Algo: To determine if $\text{KB} \models \alpha$.

- convert $(\text{KB} \wedge \neg \alpha)$ to CNF
- keep checking to see if we can use resolution to produce a new clause.
 - if ever we produce the empty clause (equivalent to False), we have a contradiction and $\text{KB} \models \alpha$
 - otherwise, if we can't add new clauses, no entailment

Ex: Does $(A \vee B) \wedge (\neg B \vee C) \wedge (\neg C)$ entail A

$$\bullet \text{KB} \wedge \neg A = (A \vee B) \wedge \underbrace{(\neg B \vee C)}_{\hookrightarrow \neg B} \wedge \underbrace{(\neg C)}_{\hookleftarrow} \wedge (\neg A)$$

$$\rightarrow \quad \underbrace{(A \vee B)}_{\text{vs } A} \wedge (\neg B \vee C) \wedge (\neg C) \wedge (\neg A) \wedge (\neg B)$$

$$\rightarrow (A \vee B) \wedge (\neg B \vee C) \wedge (\neg C) \wedge \underline{\neg A} \wedge (\neg B) \wedge \underline{\neg A}$$

$$\rightarrow (A \vee B) \wedge (\neg B \vee C) \wedge (\neg C) \wedge (\neg A) \wedge (\neg B) \wedge (\neg A) \wedge C$$

\rightarrow Empty clause \Rightarrow contradiction

f. First-order logic.

Use Constant symbols and predicate symbols.

Rinewa	person
Pomona	House
Horace	BelongsTo
Gilderoy	
Gryffindor	
Hufflepuff	
Ravenclaw	
Slytherin	

Ex. Person (Rinewa) \Rightarrow Rinewa is a Person.

House (Gryffindor) \Rightarrow Gryffindor is a house
House (Rinewa) \Rightarrow Rinewa is not a house

BelongsTo (Rinewa, Gryffindor)

\Rightarrow minimize number of items needed.

→ Universal Quantification

$\forall x . \text{BelongsTo}(x, \text{Gryffindor}) \rightarrow \neg \text{BelongsTo}(x, \text{Hufflepuff})$
 \Rightarrow Anyone in Gryffindor is not in Hufflepuff.

→ Existential Quantification

$\exists x . \text{House}(x) \wedge \text{BelongsTo}(\text{Rinewa}, x)$
 \Rightarrow Rinewa belongs to a house.

More general:

$\forall x, \text{Person}(x) \rightarrow (\exists y, \text{House}(y) \wedge \text{BelongsTo}(x, y))$

Chapter 3: Uncertainty

a. Probability

Possible world: ω Ω : all possible worlds.

Relevant proba: $P(\omega)$

Rules: $\forall \omega, 0 \leq P(\omega) \leq 1$
↓ impossible event ↑ certain event

$$\sum_{\omega \in \Omega} P(\omega) = 1$$

Def: Unconditional probability: degree of belief in a proposition in the absence of other any evidence.

Def Conditional probability: degree of belief in a proposition given some evidence that has already been revealed.

→ Notation $P(a|b)$: "proba of a given b"

property: $P(a|b) = \frac{P(a \wedge b)}{P(b)}$ or $P(a \wedge b) = P(b) P(a|b)$

Def Random variable: a variable in proba theory with a domain of possible values it can take on

Ex: roll a dice Roll: {1, 2, 3, 4, 5, 6}
↳ random variable.

Def Probability distribution: take all values of a random variable and return the proba of each

Ex. Notation $P(\text{Roll}) = \left\langle \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6} \right\rangle$

↑ This notation implies that we know the order of values of Roll.

Def Independence: the knowledge that one even occurs does not affect the proba of the other event

Consequence: $P(a \wedge b) = P(a) P(b|a)$

$P(a \wedge b) = P(a) P(b) \rightarrow$ This a def for being independent.

$$P(a \wedge b) = P(b) P(a|b) = P(a) P(b|a)$$

Bayes' Rule: $P(b|a) = \frac{P(a|b) P(b)}{P(a)}$

Def Joint probability: probability of a combination of two values from two random variable.

Def Joint probability distribution: all proba of all possible combination of two random variable.

rule: Negation: $P(\neg a) = 1 - P(a)$

rule: Inclusion - Exclusion: $P(a \vee b) = P(a) + P(b) - P(a \wedge b)$

rule: Marginalization: $P(a) = P(a \wedge b) + P(a \wedge \neg b)$
↳ for random variable.

$$P(X=x_i) = \sum_j P(X=x_i, Y=y_j)$$

Marginalization is useful if I have joint proba.

rule: Conditioning: $P(a|b) = P(a \cap b) / P(b)$

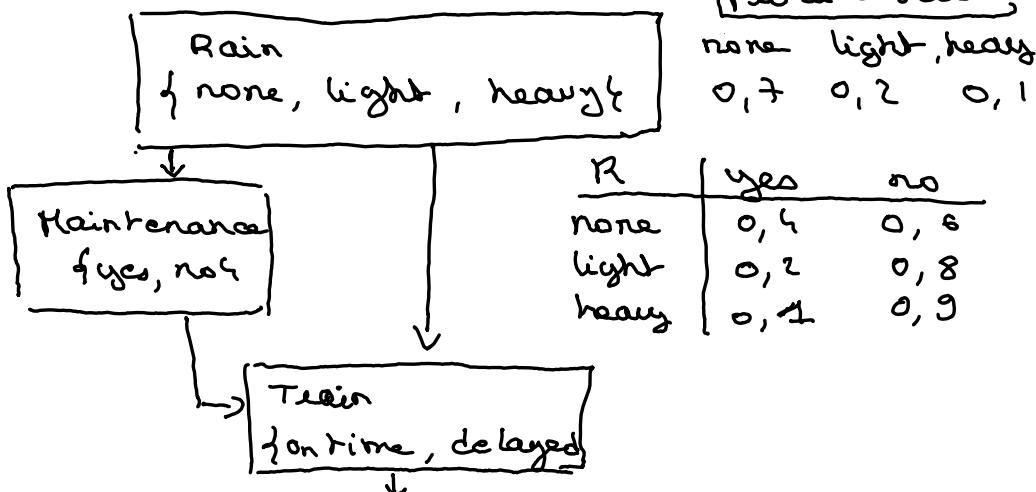
Generalize:

$$P(X = x_i) = \sum_j P(X = x_i \cap Y = y_j) P(Y = y_j)$$

Def Bayesian Network: data structure that represents the dependencies among random variables.
 → directed graph

- each node represents a random variable
- arrow from X to Y means X is a parent of Y
- each node X has probability distribution $P(X | \text{Parents}(X))$.

Ex:



R	M	On Time	delay
none	yes	0,8	0,2
none	no	0,9	0,1
light	yes	0,6	0,4
light	no	0,7	0,3
heavy	yes	0,4	0,6
heavy	no	0,5	0,5

T	attend	miss
on time	0,9	0,1
delayed	0,6	0,4

)

Prob distrob.

Rq: it is easier usually to create the dataset with dependency relationship from a node to another.

=> Computing joint probabilities with Bayes Rule.

$$\text{Ex: } P(\text{light} \wedge \text{no} \wedge \text{delayed} \wedge \text{miss})$$

$$= P(\text{light}) P(\text{no} | \text{light}) P(\text{delayed} | \text{light} \wedge \text{no}) P(\text{miss} | \text{delayed})$$

b. Inference

Can we get Prob. of new variable with what we have. . Same idea as knowledge chapter.

Inference Problem:

- Query X : variable for which to compute distribution
- Evidence: variables E : observed variables for event e
- Hidden Variables H : non-evidence, non-query variable.

=> Goal: calculate $P(X | e)$

Ex: Same as previous page.

$$P(\text{appointment} | \text{light} \wedge \text{no}).$$

Query: Appointment = A

Evidence: light, no = L, N

hidden: train. = T

$$P(\text{Appointment} | \text{light} \wedge \text{no}) = \alpha P(\text{Appointment} \wedge \text{light} \wedge \text{no})$$

iterate
over values
of hidden
variable

$$P(A \wedge L \wedge N) = \alpha P(A \wedge L \wedge N \wedge \text{onTime}) + \alpha P(A \wedge L \wedge N \wedge \text{delayed})$$

$$P(X, e) = P(X \wedge e)$$

formula : Inference by Enumeration

$$P(X|e) = \alpha P(X, e) = \alpha \sum_y P(X, e, y)$$

X : query variable
e : evidence

y : ranges over values of hidden variables
 α : normalize the result.

Rq : There is library that already implement this function.

ex : pomegranate library

- you have to describe all the nodes saying if it is conditional proba or not.
- create Bayesian network
- create links between nodes.

▷ Not particularly efficient.

c. Approximate Inference

- Sampling: Create samples from possible value of each node.
- Then when I search the proba of a value, I search the number of samples that has this value. Same if I look for a combination. If I want conditional, I only search amongst samples where there is the known value.

▷ only works if you have a great amount of sample

▷ less precise but faster.

(not 100% precise but if enough samples, good enough.)

$$\text{Prueba : } P(a) = \frac{\text{Number samples with } a}{\text{Number total samples}}$$

$$P(a \cap b) = \frac{\text{Number samples with } a \text{ and } b}{\text{Number total samples}}$$

$$P(a|b) = \frac{\text{Number samples with } a \text{ and } b}{\text{Number samples with } b}$$

d. Likelihood Weighting.

Principle:

- Start by fixing the values for evidence variables
- Sample the non-evidence variables using conditional probabilities in the Bayesian Network
- Weight each sample by its likelihood: the probability of all the evidence.

Ex: Same as previous:

We want to calculate $P(R = \text{light} | T = \text{on time})$?

→ Fix evidence variable

⇒ R = ?

M = ?

T = On Time

A = ?

→ Sample non evidence: R = light
 M = yes
 T = on time
 A = miss

→ After getting a sample, weight it
here $P(\text{On Time} | \text{light } n \text{ yes}) = 0, 6.$

⇒ Repeat procedure to have lot of sample.

⇒ Now same principle of Sampling to find proba.
↓ Approximate.

e. Prediction on Time.

Def **Markov assumption**: The assumption that the current state depends on only a finite fixed number of previous state

Def **Markov chain**: a sequence of random variable where the distribution of each variable follows the Markov assumption.

Ex: predict if sunny or raining
S R

x_t : weather at time t (today)

x_{t+1} : weather at time $t+1$

		Tomorrow (x_{t+1})	
		S R	
Today (x_t)	S	0,8	0,2
	R	0,3	0,7

⇒ Transition model.

→ create Markov Chain

$S \rightarrow S \rightarrow R \rightarrow R \rightarrow R \dots$
 $x_0 \quad x_1 \quad x_2 \quad x_3 \quad x_4$

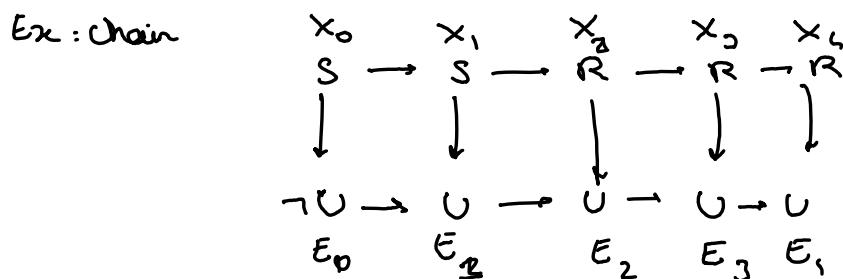
Application : Find Hidden state from observation.

ex:	Hidden state	Observation
	robot's position words spoken user engagement weather (S or R)	robot's sensor data audio waveforms website or app analytics umbrella (U)

Def Hidden Markov Model : a Naive Bayes model for a system with hidden states that generate some observed event.

Ex:	Observation (E_t)	
	U	$\neg U$
State (X_t)	S 0,2 0,8	
R 0,9 0,1		

Def Serial Naive Assumption : the assumption that the evidence variable depends only on the corresponding state.



Ex of task done:

Task	Definition
filtering	given observation from start until now calculate distribution for current state
prediction	given observation from start until now calculate distribution for future state
smoothing	given observation from start until now calculate distribution for past state
most likely explanation	given observation from start until now calculate most likely sequences of state.

Chapter 4 : optimization.

Def Optimization: choosing the best optia from a set of options

a. Local Search.

Def local search : search algo that maintain a single node and searches by moving to a neighbor node.

Def state-space landscape: tree graph where each tree represents a possible state and its value.

Then we can search :

- global max / min
- local min

⇒ for that we define an objective function (for max)
cost function (for min)

⇒ **Hill climbing Algorithm principle.**

function Hill-Climbing (problem):

 current = initial state of problem.

 repeat :

 neighbor = highest valued neighbor of current

 if neighbor not better than current

 return current

 current = neighbor

⚠ It will not always give the optimal soln. (global)
but a local min / max

Hill climbing variants:

Variant	Definition
Steepest-Ascent	choose the highest-valued neighbor
Stochastic	choose randomly from higher-valued neighbors
first-choice	choose the first higher-valued neighbor
random-start	conduct hill climbing multiple times
local beam search	choose the k highest-valued neighbors

b. Simulated annealing

Principle

- Early on, higher "temperature": more likely to accept neighbors that are worse than current state.
- Later on, lower "temperature": less likely to accept neighbors that are worse than current state.

Pseudo code

```

function simulated-annealing(problem, max):
    current = initial state of problem
    for t = 1 to max
        T = temperature (+)
        neighbor = random neighbor of current
        ΔE = how much better neighbor is than current:
        if ΔE > 0: (neighbor is better)
            current = neighbor
        with probability  $\exp(\frac{\Delta E}{T})$  set current = neighbor
        between current
    
```

Ex: Traveling sales man Pb:

visit cities and get back home by minimizing the distance.

C. Linear Programming.

Problem:

- minimize a cost function $c_1x_1 + c_2x_2 + \dots + c_nx_n$
- with constraints of form $a_1x_1 + a_2x_2 + \dots + a_nx_n \leq b$
or of form $a_1x_1 + \dots + a_nx_n = b$
- with bounds for each variables $l_i \leq x_i \leq u_i$

Ex: two machines x_1 and x_2 costs 50 \$/h and 80 \$/h to run. goal is to minimize cost

• x_1 requires 5 units of labor per hour. x_2 requires 2 units of labor per hour. Total of 20 units of labor to spend.

• x_1 produces 10 units of output/h. x_2 produces 12 units of output per hour. company needs 90 units of output.

Cost function: $50x_1 + 80x_2$ x_i : number.

constraints: $5x_1 + 2x_2 \leq 20$ \leftarrow we need \leq or \geq
 $10x_1 + 12x_2 \geq 90 \Rightarrow (-10x_1) + (-12x_2) \leq -90$

Aleg possible: simplex

d. constraint satisfaction.

⇒ Constraint satisfaction Problem

- Set of variables $\{x_1, \dots, x_n\}$
- Set of domains for each variable $\{D_1, D_2, \dots, D_n\}$
- Set of constraints C

Ex: Sudoku.

Variables: empty cells.

Domains: values between 1 and 9

Constraints: cells \neq in row, column and square.

Def hard constraints: constraints that must be satisfied in a correct solution.

Def soft constraints: constraints that express some notion of which solutions are preferred over others.

Def unary constraint: single variable not equal to a given value (ex: $A \neq \text{Monday}$)

Def binary constraint: constraint involving two variables (ex: $A \neq B$)

Def node consistency: when all the values in a variable's domain satisfy the variable's unary constraints

Ex: $A \text{ } \underline{\hspace{2cm}} \text{ } B$

$\{\text{Mon, Tue, Wed}\}$

$\{\text{Mon, Tue, Wed}\}$

constraint: $\{A \neq \text{Mon}, B \neq \text{Tue}, B \neq \text{Wed}, A \neq B\}$

$\Rightarrow A \text{ } \underline{\hspace{2cm}} \text{ } B$ node consistency. p3 row.
 $\{\text{Tue, Wed}\}$ $\{\text{Wed}\}$

Def arc consistency: when all the values in a variable's domain satisfy the variable's binary constraints.

To make X arc-consistent with respect to Y , remove elements from X 's domain until every choice for X has a possible choice for Y .

Ex: $A \longrightarrow B$
 $\{Tue, Wed\} \quad \{Wed\}$

same constraints

$\Rightarrow A \longrightarrow B$ arc consistency pb.
 $\{Tue\} \quad \{Wed\}$

Arc consistency pseudo-code:

function revise(csp, X, Y):
 revised = false (if no change in X domain, True if change)
 for x in X .domain:
 if no y in Y .domain satisfies constraint $p_{x,y}(X, Y)$:
 remove x from X .domain
 revised = True.
 return revised

↳ This function makes one variable arc-consistency not all.

To make an entire pb arc consistent, we use the following function:

function AC-3 (CSP) : CSP: constraint satisfaction pb

queue = all arcs in CSP

while queue non empty:

- (X, Y) = Dequeue (queue)
- if revise (CSP, X, Y) :
- if size . of X.domain == 0 :
- return false (no way to solve pb)
- for each Z in X.neighbors - {Y} :
- Enqueue (queue, (Z, X)).

return true.

This algo reduce the domain for each variable but do not always solve the pb. Some variable can still have multiple values.

=> CSP as search problem.

- Initial state: empty assignment (no variable)
- actions : add a {variable = value} to assignment
- transition model: shows how adding an assignm changes the assignm
- goal test: check if all variables assigned and constraint all satisfied
- path cost function: all paths have same cost

He is a example of an algo pseudo code.

Pseudo-code.

```
function Backtrack (assignment, cap) :  
    if assignment complete :  
        return assignment.  
    var = Select-Unassigned - Var (assignment, cap)  
    for value in Domain - Values (var, assignment, cap) :  
        if value consistent with assignment :  
            add {var = value} to assignment  
            result = Backtrack (assignment, cap).  
            if result ≠ failure : return result.  
            remove {var = value} from assignment  
    return failure. (no value to solve issue)
```

Rq: constraint python library has all this function already implemented.

e. Inference.

Add arc consistency to backtrace

Def maintaining arc-consistency: algorithm for enforcing arc-consistency every-time we make a new assignment

Every time we make a new assignment to x , calls AC-3, starting with a queue of all arcs (y, x) where y is a neighbor of x

New pseudo code :

Pseudo - Code.

```
function BackTrack (assignment, csp)
    if assignment is complete : return assignment
    val = Select-Unassigned-Value (assignment, csp)
    for value in Domain-Values (val, assignment, csp):
        if value consistent with assignment:
            add {val = value} to assignment.
            inferences = Inference (assignment, csp)
            if inferences ≠ failure: add inferences to assignment
            result = Backtrack (assignment, csp)
            if result ≠ failure : return result.
            remove {val = value} and inferences from assignment
    return failure.
```

f. Select - Unassigned - value.

To be more efficient, it is important to select well this variable.

Def minimum remaining values (MRV) heuristic : select the variable that has the smaller domain.

Def degree heuristic : select the variable that has the higher degree.

Def node degree: number of node that are attached to a node. number of node that are constraint by this node

g. Domain values

Choosing well the order of the value in the variable selected is also a way to go faster.

Def least-constraining values heuristic: return variables values in order by number of choices that are ruled out of neighboring values.

- Try least-constraining value first

Chapter 6: Learning

I. Supervised Learning.

Def **Supervised Learning**: given a data set of input-output pairs, learn a function to map inputs to outputs

a. classification

Def **Classification**: supervised learning task of learning a function mapping an input point to a discrete category

Ex: predict weather.

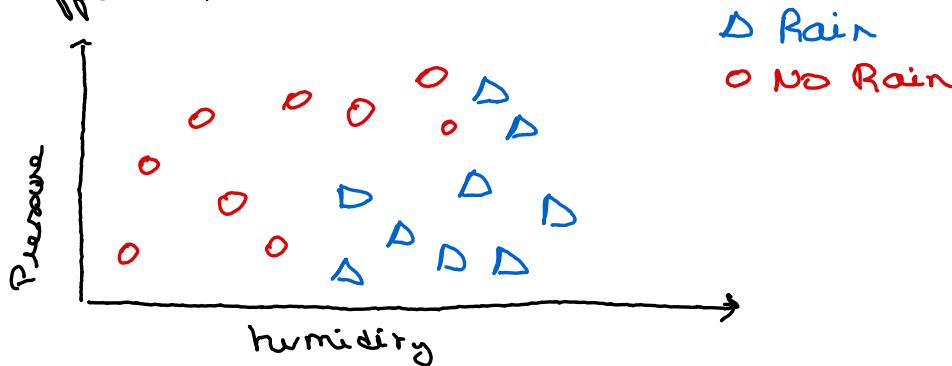
supervised.

Data:

Date	Humidity	Pressure	Rain
Jan 1	93	999,7	Yes
Jan 2	49	1015,5	No
Jan 3	79	1031,5	No
Jan 4	65	984,9	Yes
Jan 5	90	975,2	Yes

$$f(\text{humidity}, \text{pressure}) = \text{Rain}$$

↪ approximate with h.



Nearest approach

Def nearest - neighbor classification: algorithm that, given an input, chooses the class of the nearest data point to that input.

⚠ Tricky if you are close to the frontier. Where nearest point can be not good solution.

Def k-nearest - neighbor classification: algorithm that, given an input, chooses the most common class out of the k nearest data points to that input.

⚠ can be slow.

One solution is to draw a linear regression to be fast.

linear regression approach

Here: $x_1 = \text{humidity}$, $x_2 = \text{pressure}$.

$h(x_1, x_2) = \text{Rain if } w_0 + w_1 x_1 + w_2 x_2 > 0$
No Rain otherwise.

weight vector $w: (w_0, w_1, w_2)$

input vector $x: (1, x_1, x_2)$

goal: find weight vector to do prediction.

Def dot product: $x \cdot w = w_0 + w_1 x_1 + w_2 x_2$

$$\Rightarrow h_w(x) = \begin{cases} 1 & \text{if } w \cdot x > 0 \\ 0 & \text{otherwise} \end{cases}$$

Def perception learning rule: given data point (x, y) update each weight according to :

$$w_i = w_i + \alpha(y - h_w(x)) x_i \quad \alpha: \text{learning rate}$$

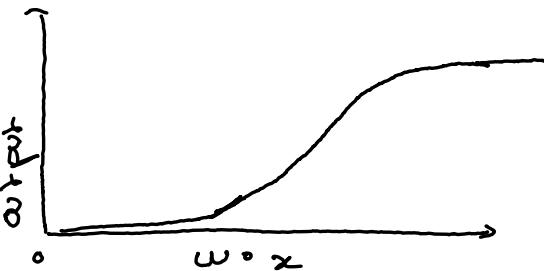
actual value estimate

Def hard threshold: does not allow uncertainty give 1 or 0 always



$$\begin{cases} 1 & \text{if } w \cdot x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Def soft threshold: can give every value between 1 and 0. give info about how confident we are



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

5. Support vector machines.

Def maximum margin separator: boundary that maximizes the distance between any of the data points.

The support vector machines (SVC) will find this one
=> good with data that are not linearly separable.

Def regression: supervised learning task of learning a function mapping an input point to a continuous value.

Ex: Impact of advertising

$$f(\text{advertising}) = \text{sales}$$

↑ money spent ↑ money make

⇒ linear regression approach possible

c. Evaluating Hypotheses.

Def loss function: function that expresses how poorly an hypothesis performs

⇒ For classification (discrete)

Def 0-1 loss function:

$$L(\text{actual}, \text{predicted}) = \begin{cases} 0 & \text{if actual} = \text{predicted} \\ 1 & \text{otherwise} \end{cases}$$

⇒ For regression (continuous)

Def L₁ loss function:

$$L(\text{actual}, \text{predicted}) = |\text{actual} - \text{predicted}|$$

Def L₂ loss function:

$$L(\text{actual}, \text{predicted}) = (\text{actual} - \text{predicted})^2$$

→ penalize more values that are too away.

Def overfitting: a model that fits too closely to a particular data set and therefore may fail to generalise to future data.

If we only want to reduce loss, we might have overfitting.

To avoid that, the best is to minimize:

$$\text{cost}(h) = \text{loss}(h) + \lambda \text{complexity}(h)$$

- Simpler model can be better.
- if λ big: penalize complexity
- if λ small: penalize loss

Def regularization: penalizing hypotheses that are more complex to favor simpler, more general hypotheses

Def holdout cross-validation: splitting data into a training set and a test set, such that learning happens on the training set and is evaluated on the test set.

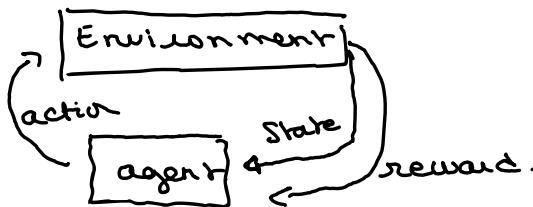
Def k-fold cross-validation: splitting data into k sets, and experimenting k times, using each set as a test set once and using remaining data as training set

d. Python libraries.

- ⇒ Scikit-Learn :
 - linear models
 - neighbor
 - SVC. ...
 - fit
 - predict
 - split data
- ⇒ CSV, pandas : read csv file.
- ⇒ random : randomize data (shuffle, select randomly...)
- ⇒ cf bank notes 1.py

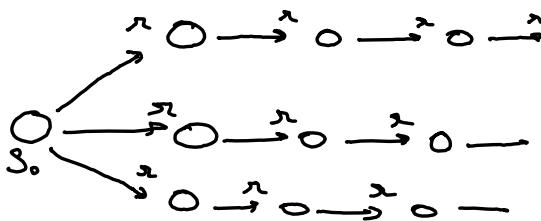
II. Reinforcement learning

Def reinforcement learning: given a set of rewards or punishments, learn what actions to take in the future.



Def Markov decision process: model for decision-making, representing states, actions and their rewards.

Rappel: Markov chain chap. 3.



- Markov model:
- Set of states S
 - Set of actions (A)
 - Transition model $P(s'|s, a)$
 - Reward function $R(s, a, s')$

Def Q-learning: method for learning a function $Q(s, a)$, estimate of the value of performing action a in state s .
(How much reward will I get?)

Principle

- Start with $Q(s, a) = 0$ for all s, a .
- When we taken an action and receive a reward:
 - estimate the value of $Q(s, a)$ based on current reward and expected future rewards
 - update $Q(s, a)$ to take into account old estimate as well as our new estimate.

Pseudo-code.

- Start with $Q(s, a) = 0$ for all s, a
- Every time we take an action a in a state s and observe a reward r , we update:

$$Q(s, a) \leftarrow Q(s, a) + \alpha((r + \gamma \max_a Q(s', a')) - Q(s, a))$$

Learning rate | new value estimate | old value estimate

\Rightarrow importance of new value compare to old (if $\alpha=1$ only new value count)

Def Greedy Decision making: when in state s , choose action a with highest $Q(s, a)$

⚠ Greedy can lead to solution that are not optimal.

Exploration vs Exploitation

↳ explore ↳ use rewards to make decision

Def ϵ -greedy:

- set ϵ equal to how often we want to move randomly
- with probability $1-\epsilon$, choose estimate best move.
- with probability ϵ , choose a random move.

You can decrease ϵ over time.

Ex: play game. aim & strategy: AI to learn how to win

Def function approximation: approximating $Q(s, a)$, often by a function combining various features, rather than storing one value for every state-action pair

III. Unsupervised learning

Def unsupervised learning: given input data without any additional feedback, learn patterns.

Def clustering: organizing a set of objects into groups in such a way that similar objects tend to be in the same group.

Ex application:

- genetic research
- image segmentation
- market research
- medical imaging
- social network analysis

Def k-means clustering: algorithm for clustering data based on repeatedly assigning points to clusters and updating those clusters' centers.

Chapter 6 : Neural Networks.

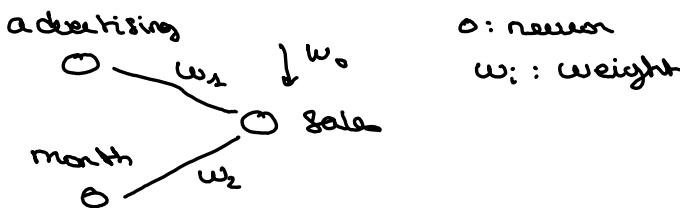
Principle :

- Neurons are connected to and receive electrical signals from other neurons
- Neurons process input signal and can be activated.

Def artificial neural network: mathematical model for learning inspired by biological neural networks.

- Model mathematical function from inputs to outputs based on the structure and parameters of the networks
- allows for learning the network's parameters based on data.

Ex:



a. One layer

Def gradient descent: algorithm for minimizing loss when training neural network.

Principle :

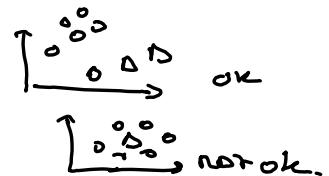
- Start with a random choice of weights
- Repeat:
 - calculate the gradient based on all data points : direction lead to decreasing loss
 - update weights according to the gradient

▷ calculate on all data points is expensive.

- alternative:
 - **Stochastic gradient descent**
 - ⇒ calculate the gradient based on one data points choose randomly
(faster but less accurate)
 - **Mini batch gradient descent**
- ⇒ calculate the gradient on a mini batch of points.

If multiple output we can see that as several neural networks and we can train them separately.

limitation of linear combination:
⇒ only for linear separable



Def Perceptron: only capable of learning linearly separable decision boundary.

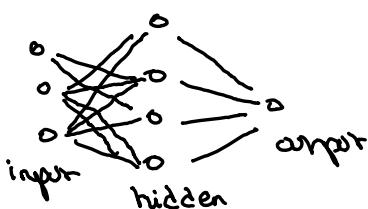
b. multilayer.

Def multilayer neural network: artificial neural network with an input layer, an output layer and at least one hidden layer.

Not like this:



but



⇒ ability to model complex situation by adding layers

Def backpropagation: algorithm for training neural networks with hidden layers

Principle:

- start with a random choice of weights:
- repeat:
 - calculate error for output layer
 - for each layer, starting with output layer and moving inward towards earliest hidden layer
 - propagate error back one layer
 - update weights.

Def deep neural networks: neural networks with multiple hidden layers

c. overfitting

Def dropout: temporarily removing units, selected at random, from a neural network to prevent over-reliance on certain units.



d. library.

⇒ **Tensorflow** (develop by google)

`playground.Tensorflow.org.` : to understand principle

klass: AP I.

e. Computer vision

Def **Computer vision**: computational methods for analysing and understanding digital images

Def **image convolution**: applying a filter that adds each pixel value of an image to its neighbors, weighted according to a kernel matrix

Ex :

kernel : $\begin{matrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{matrix}$	image : $\begin{matrix} 10 & 20 & 30 & 40 \\ 10 & 28 & 30 & 40 \\ 20 & 30 & 40 & 50 \\ 20 & 30 & 40 & 50 \end{matrix}$	Step 1 \times kernel
--	--	------------------------

First image $[3 \times 3] \times$ kernel (multiply each value by its corresponding in kernel).

$$\rightarrow 0 \times 10 + 20 + 30 \times 0 + \dots = 10$$

→ Slide kernel : $\begin{matrix} 10 & 20 \\ 40 & 50 \end{matrix}$

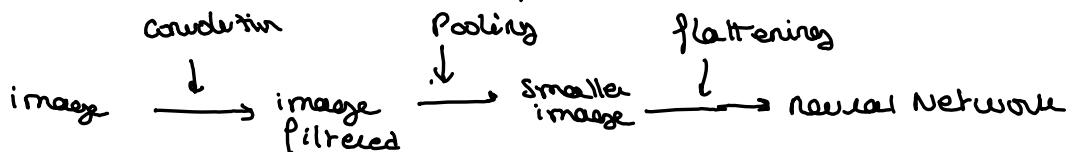
Def **pooling**: reducing the size of an input by sampling from regions in the input.

Def **max-pooling**: pooling by choosing the maximum value in each region

Ex : image : $\begin{matrix} 30 & 40 & 80 & 90 \\ 20 & 50 & 100 & 110 \\ 5 & 10 & 20 & 30 \\ 10 & 20 & 40 & 30 \end{matrix}$ $\xrightarrow{2 \times 2 \text{ Pooling}}$ $\begin{matrix} 50 & 110 \\ 20 & 40 \end{matrix}$

Def convolutional neural network (CNN): neural networks that use convolution, usually for analysing images.

⇒ CNN can learn what filter use.

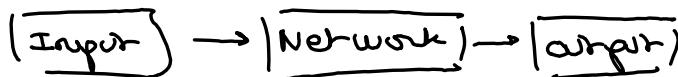


You can use it reversed time

image → convolution → pooling → convolution → pooling ... → Neural Networks
low level features (e.g., edges, curves, shape) high level features: object

+). keras.layers.Dense(activation = "softmax")
↳ poster distribution

Representation:



Def feed-forward neural network: neural network that has connections only in one direction.

Def recurrent neural network: neural network that generates output that feeds back into its own inputs



Input → network → output



network → output



network → output

:

Chapter 7: Language

I. Natural language processing

- automatic summarization
- Info extraction
- language identification
- machine translation ...
- named entity recognition
- speech recognition
- text classification
- word sense disambiguation ...

=> Need to focus on

- syntax (grammatically correct ...)
- semantic (order of words ...)

Def **Formal grammar**: a system of rules for generating sentences in a language.

=> **Context free grammar**

She saw the city
N V D N

N: noun

V: verb

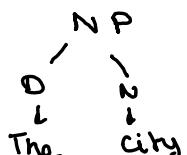
D: determinant

P: preposition

a word = terminal symbol

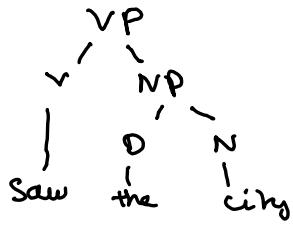
type : non terminal symbol

$NP \rightarrow N \mid DN$ (noun phrase) (noun or Determinant + noun)



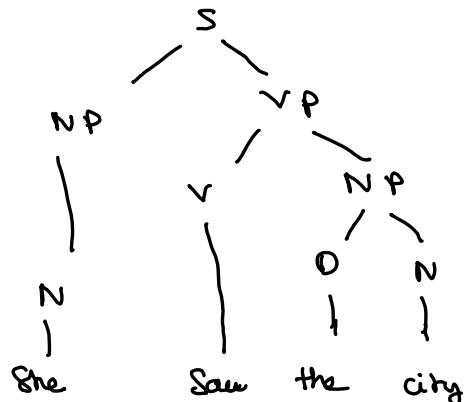
$VP \rightarrow V \mid VN$

(verb phrase)



$S \rightarrow NP VP$ (sentence)

If a sentence is not grammatically correct it will not be possible to create such a tree.



good if not that much rules.

\Rightarrow Python library: nltk.

Def **n-gram**: a contiguous sequence of n items from a sample of text.

Def **character n-gram**: a contiguous sequence of n characters from a sample of text

Def **word n-gram**: a contiguous sequence of n words from a sample of text

Def **unigram**: a contiguous sequence of 1 item from a sample of text.

Def **bigrams**: a contiguous sequence of 2 items from a sample of text.

Def **trigrams**: a contiguous sequence of 3 items from a sample of text.

can be extract n gram as all AI might have seen them already in other text.

Def tokenization: the task of splitting a sequence of characters into pieces (tokens)

Def word tokenization: the task of splitting a sequence of characters into words.

=> any time see a space create an item.

▷ pb punctuation that still there.

II N-grams model.

Once we know the most common n-gram, we can deduce some word from other.

Ex: if "it was a" is the most common trigram, if we have "it was" we can imagine having "a" after.

=> python library: markovify + nltk.

III. Text categorization

Ex: mail : spam or not
sentiment analysis.

Def bag of word model: model that represents text as an unordered collection of words.

=> **Naive Bayes approach**.

p: positive, n: negative.

Ex: Text: "my grandson loved it"

c = ["my", "grandson", "loved", "it"] = c₁, c₂, c₃, c₄

$$P(c|p) = \frac{P(c_1|p) P(c_2|p) \dots P(c_n|p)}{P(c)} \underset{\text{proportional}}{\sim} P(c_1|p) P(c_2|p) \dots P(c_n|p) = P(c \wedge p)$$

Bayes' rule.

$$P(c \wedge p) = P(c_1 \wedge c_2 \wedge c_3 \wedge c_4 \wedge p)$$

$$= P(p) P(c_1|p) P(c_2|p) P(c_3|p) P(c_4|p)$$

We assume
that words
are independent
if we know it is positive
message (Naive Bayes)

$$P(p) = \frac{\text{num positive sample}}{\text{num sample}}$$

$$P(c_i|p) = \frac{\text{num positive sample with } c_i}{\text{num positive sample}}$$

P		n
0,49		0,51

$$\Rightarrow P(c \wedge p) = 0,00014112$$

$$P(c \wedge n) = 0,00006528$$

C _i	P	n
C ₁	0,3	0,2
C ₂	0,01	0,02
C ₃	0,32	0,08
C ₄	0,3	0,4

\Rightarrow Treat value as probability distribution
and normalized them

$$\Rightarrow P(p|c) = 0.6837 \quad \Rightarrow \text{we can guess the message is positive}$$

$$P(n|c) = 0.3169$$

∴ if 0 is in the data set : ex if $p(c_2|p) = 0$. (c_2 never happen) $\Rightarrow P(c \wedge p) = 0$ (problem because c_2 never happen in positive sentence).

Def **additive smoothing**: adding a value α to each value in our distribution to smooth the data.

Def **Laplace smoothing**: adding 1 to each value in our distribution : pretending we've seen each value one more time than we actually have.

\Rightarrow avoid multiplying by 0.

=> python library : nltk.

IV . Information retrieval.

Def information retrieval : the task of finding relevant documents in response to a user query

Def topic modeling : models for discovering the topics for a set of documents.

Def term frequency : the number of times a term appears in a document.

Def function words : words that have little meaning on their own , but are used to grammatically connect other words.

Def content words : words that carry meaning independently.

=> usually when we look a document we are interested in the term frequency of content word.

=> When I do Tf, I ignore function words

Def inverse document frequency : measure of how common or rare a word is across documents

=> gives better idea of specific topic of a document compare to others - (how rare is a word in the documents)

=> inverse document frequency (word) = $\log \left(\frac{\text{Total documents}}{\text{Num documents containing(word)}} \right)$

Def Tf-idf : ranking of what words are important in a document by multiplying term frequency (Tf) by inverse document frequency (idf)

V. Semantico.

Def **Information extraction**: the task of extracting knowledge from documents -

⇒ give AI template to search .

ex : when a company was founded in a year

→ will learn name of company and year of creation
⇒ ↗ needs to have template → long to write.

→ If we give data instead of template, AI can search and learn the template by looking at sentence that relate data . Then use it for other data .

Only works if we find template .

⇒ **wordnet** : dictionnay that give all definition of a word and its relation with others.

↗ Issue with word that change

a. word representation .

Representing word with number to use them in a Neural Network for exemple .

Ex : "He wrote a book"

He : [1, 0, 0, 0] , wrote : [0, 1, 0, 0] a : [0, 0, 1, 0] book : [0, 0, 0, 1]

Def **one-hot representation** : representation of meaning as a vector with a single 1, and with other values as 0.

↗ if dictionnay have 10^5 words, each vector is 10^5 dig.

↳ each word are disconnected:

ex: write and author can be totally \neq (nothing to do)
(we want vector that are similar)

Def **distributional representation**: representation of meaning distributed across multiple values.

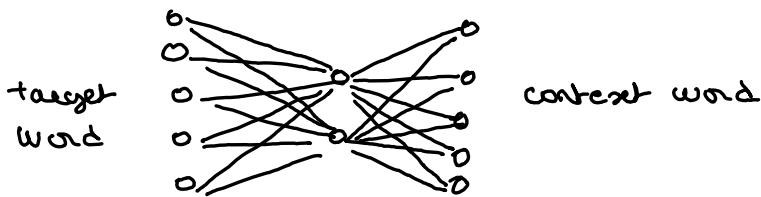
Ex: he : [-0.34, -0.08, 0.02, ...]

wrote : [-0.27, 0.40, 0.00, ...]
⋮

We can hope that similar word have similar vector.

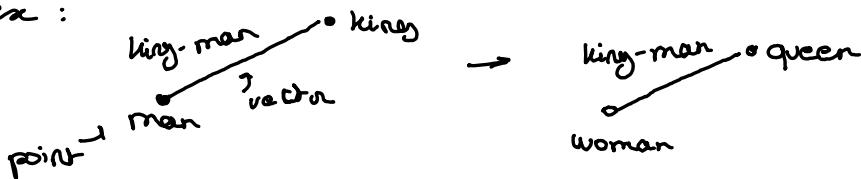
Def **Word2Vec**: model for generating word vectors

Def **Skip-gram architecture**: neural network architecture for predicting context words given a target word



if you have vector you can do some calculation.

ex :



⇒ can help to find new word with other.

$$\text{queen} = \text{woman} + \text{king-man}$$