

DATA SCIENCE & ARTIFICIAL INTELLIGENCE



Data Science

1. Sample collection
2. Data Preparation
3. Exploratory Analysis
4. Analytic Visualisation
5. Algorithmic Optimisation
6. Information Presentation
7. Ethical Consideration

- Practical Motivation
- Problem Formulation
- Statistical Description
- Pattern Recognition
- Machine Learning
- Statistical Inference
- Intelligent Decision

Common Problems

1. Numeric → How much? How many? (Regression)
2. Classes → Type A? Type B? (Classification)
3. Structure → How is this organised? (Clustering)
4. Anomaly → Is it weird behavior? (Anomaly Detection)
5. Action → what should be done next? (Adaptive Learning) → decision

→ prediction

→ detection

Common Data Types

1. Structured Data

↳ Highly Organised Data, Clearly Defined Variables, Easy to Mine & Analyse

a) Numeric Data → Numeric Continuous Variables

b) Categorical Data → Factor / Level / Class Variables (e.g. safety → high, med, low. seats → 2, 4, more)

c) Mixed Data → Numeric & Categorical

d) Time Series Data → Numeric with TimeStamps (e.g. Stock & Equity Markets, Weather Data over Time, Prices & Promotions)

e) Network Data → Nodes (Social Networks & Web, Transport Networks (MRT), Financial Transactions)

Spreadsheet (Excel, CSV)
Standard SQL Databases
Sensors & Devices



2. Unstructured Data

↳ Highly Unorganised Data, Non-Obvious Variables, Highly Context-Sensitive (Words, Phrases, Emoticons)

a) Image Data → Pixels & Objects (Social networks, mobile phone cameras, blogs, documents) & Text

b) Video Data → Images, Frames, Objects (youtube, video messages & calls)

c) Voice Data → Voice Signal & Waves (songs, microphones & cameras, recordings, announcements)

Uni-Variate Statistics

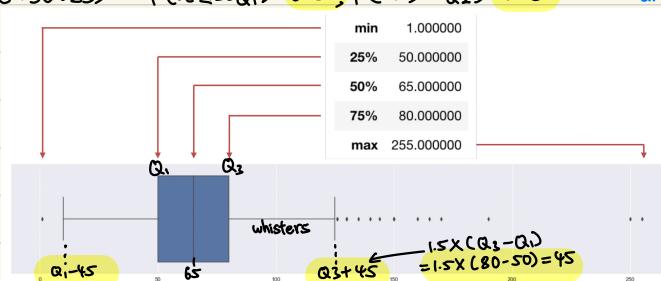
→ Mean (Average) → $\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n}$ → no. of datapoints

→ Standard Deviation (Average Deviation from mean) → $\sqrt{\frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}}$

→ Median (Mid-Value) → $P(x \leq x_m) = P(x \geq x_m) = 0.5$ $\xrightarrow{\text{min} \rightarrow x_1 \rightarrow \text{max}}$ median

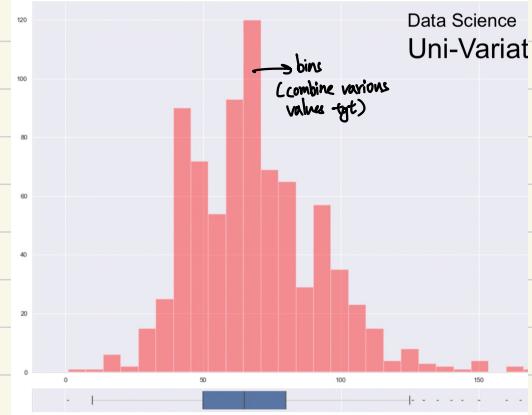
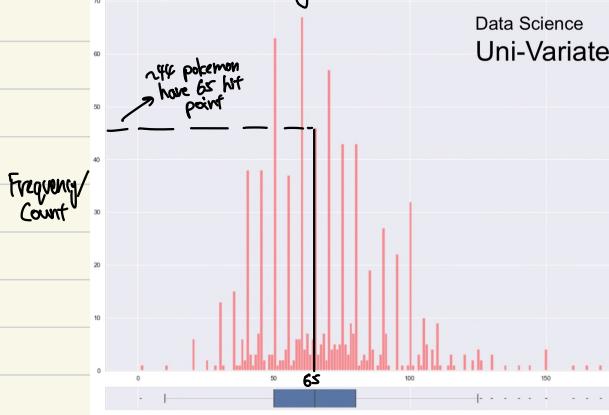
→ Quantiles (Divide Data 25:50:25) → $P(x \leq x_{Q1}) = 0.25$, $P(x \geq x_{Q2}) = 0.25$ $\xrightarrow{\text{min} \rightarrow x_{Q1} \rightarrow x_{Q2} \rightarrow \text{max}}$

count	800.000000
mean	69.258750
std	25.534669
min	1.000000
25%	50.000000
50%	65.000000
75%	80.000000
max	255.000000

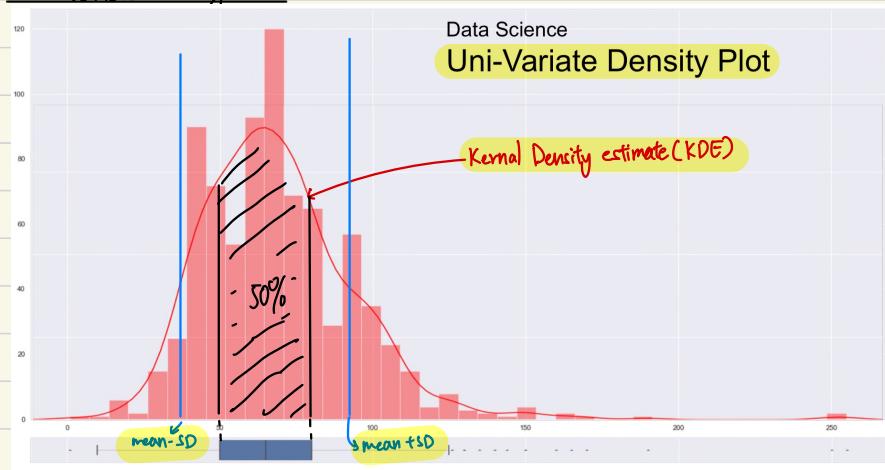


← Uni-Variate Box Plot

Uni-Variate Histogram



Uni-Variate Density Plot



~~need memorise~~ (only for Normal Distribution)

mean - SD to mean + SD $\rightarrow 68.27\%$

mean - 2*SD to mean + 2*SD $\rightarrow 95.45\%$

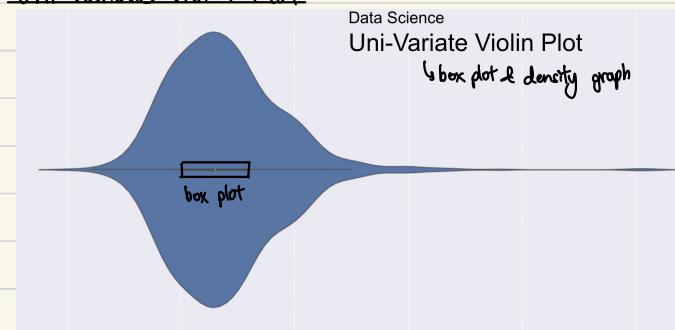
mean - 3*SD to mean + 3*SD $\rightarrow 99.73\%$

% of data that lies in b/w the points

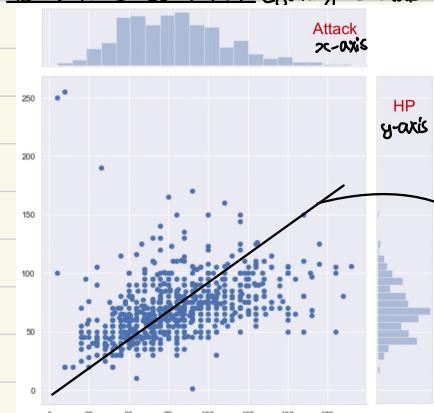
Confidence	Interval in terms of SD
90 %	Mean +/- 1.645 * SD
95 %	Mean +/- 1.96 * SD
99 %	Mean +/- 2.576 * SD

90% confidence that data lies b/w mean +/- 1.65*SD

Uni-Variate Violin Plot



Bi-Variate Joint Plot (joining 2 stats tgt)



HP \rightarrow Plotted along Y axis

Attack \rightarrow Plotted along X axis

\rightarrow HP \uparrow as Attack \uparrow

\rightarrow Dependence moderately strong [Correlation Coefficient $\rightarrow r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$]

\rightarrow The closer the dots to line & the closer the line to $y=x \rightarrow$ the stronger the correlation

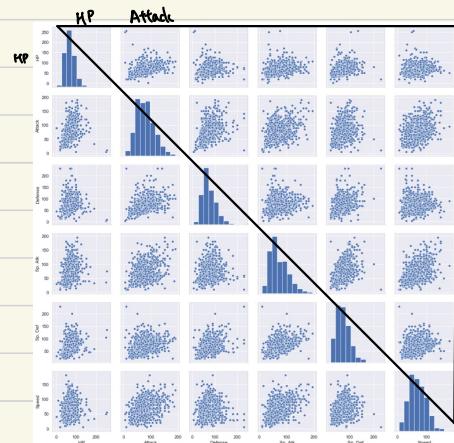
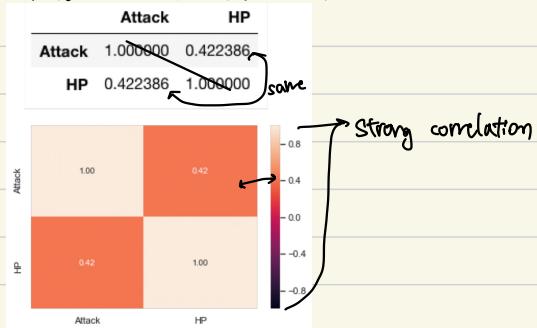
- No dependence \rightarrow corr = 0

(covered below)

- Perfect Positive \rightarrow corr = +1

- Perfect Negative \rightarrow corr = -1 (one \downarrow , one \uparrow)

Bi-Variate Relation (Correlation Matrix & Plot)



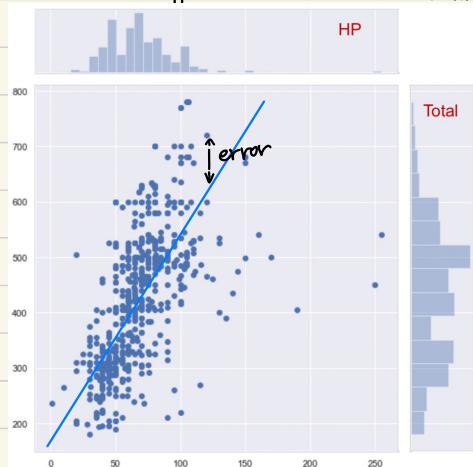
Multi-Variate Statistics

	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed
HP	1.000000	0.422386	0.239622	0.362380	0.378718	0.175952
Attack	0.422386	1.000000	0.438687	0.396362	0.263990	0.381240
Defense	0.239622	0.438687	1.000000	0.223549	0.510747	0.015227
Sp. Atk	0.362380	0.396362	0.223549	1.000000	0.506121	0.473018
Sp. Def	0.378718	0.263990	0.510747	0.506121	1.000000	0.259133
Speed	0.175952	0.381240	0.015227	0.473018	0.259133	1.000000

Bi-Variate Exploration

- Split the dataset into train (use to train the model) & test data (use to test the model)
- Learn relationship from Train & try to predict the value on test

Uni-Variate Regression (Statistical Modeling)



Hypothesize a Linear Model

$$\text{Total} = a \times \text{HP} + b + \epsilon$$

gradient y-intercept

Learn from Train Data & Predict Test Data (a & b same)

Goodness of Fit of Model

1. Mean Squared Error (MSE)

$$MSE = \frac{1}{n} \sum (Total - a \times HP - b)^2$$

model test data.

The lower the MSE the better the model

2. Explained Variance (R^2)

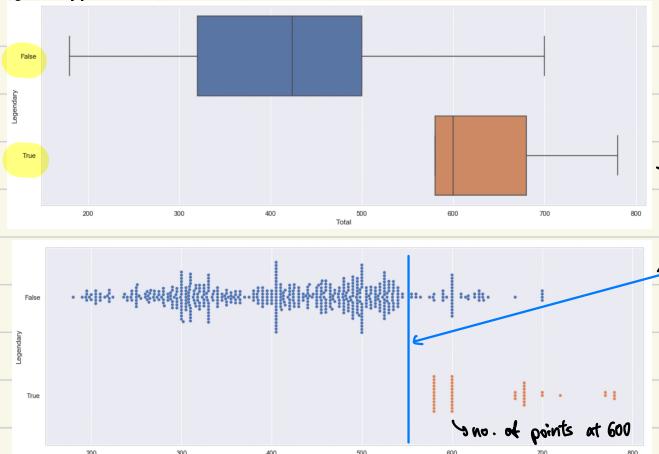
$$R^2 = 1 - \frac{\sum (Total - a \times HP - b)^2}{\sum (Total - \bar{Total})^2}$$

Residual sum of squares (RSS) → how much error we make
Total sum squared (TSS)

The higher the R^2 , the better the model, $0 \leq R^2 \leq 1$

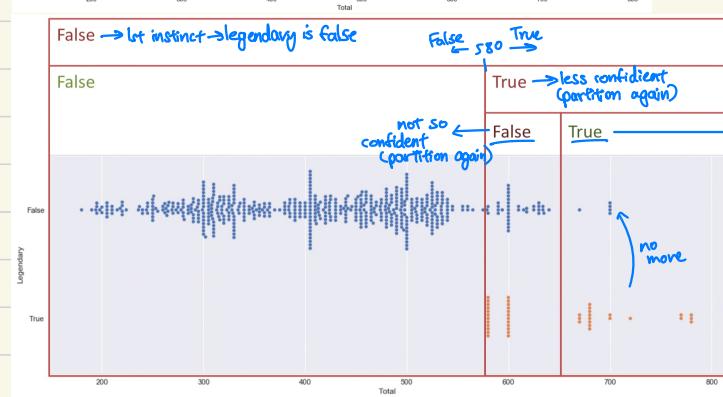
$$R^2 = 1 - \frac{MSE}{var(total)}$$

Binary Classification



Box Plot

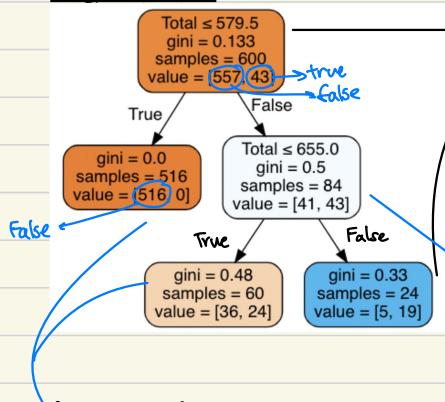
→ legendary vs Total have a strong relationship
draw a line at 550, when total < 550, most likely False



SwarmPlot

Data Space

Decision Tree

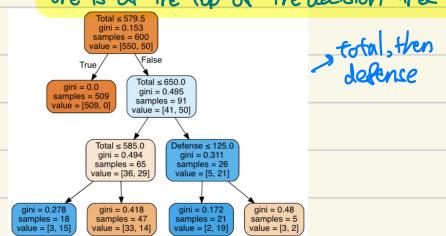


- darken red → confident of it being false
- darken blue → confident of it being true

Gini Index (metric of misclassification)

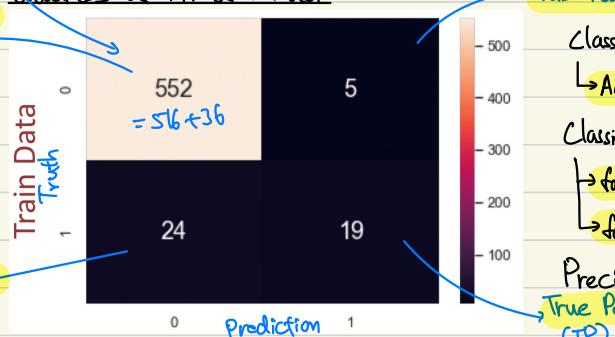
$$\text{gini} = \frac{x}{n}(1 - \frac{x}{n}) + \frac{y}{n}(1 - \frac{y}{n})$$

As Most impf factor \rightarrow based on which one is at the top of the decision tree



Goodness of Fit of Model

True Negatives
(TN) ←
correct prediction of it being false



۱۶

$$\hookrightarrow \text{Accuracy} : \frac{552+19}{552+5+24+19} = \frac{552+19}{600} = 0.952$$

Classification Errors

$$\rightarrow \text{false positive rate} = \frac{FP}{TN+FP} = \frac{5}{552+5} = \frac{5}{557}$$

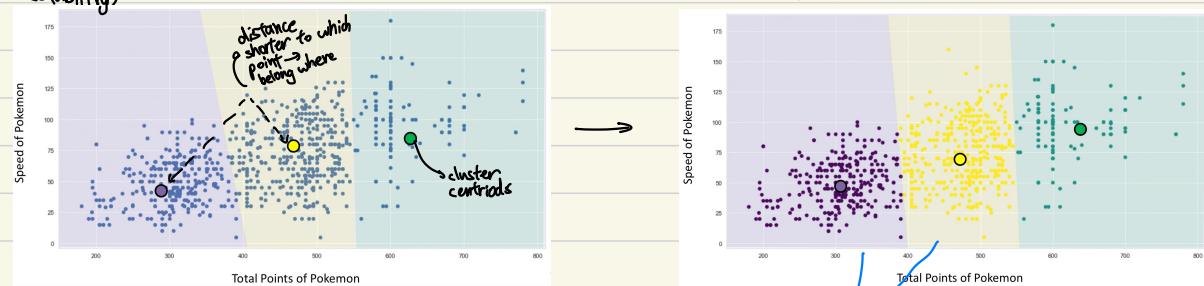
$$\rightarrow \text{false negative rate} = \frac{FN}{TP+FN} = \frac{24}{24+19} = \frac{24}{43}$$

Precision = $\frac{TP}{TP+FP}$ → predicted positive (TP + FP)

Precision = $\frac{TP}{TP+FP}$ → predicted positive ($TP + FP$)
True Positive
(TP)

K-Means Clustering

1. Choose $K \rightarrow$ potential number of clusters & choose K cluster centroids from dataset
2. for each point in dataset \rightarrow re-label according to nearest centroid \rightarrow recompute the centroid of cluster \rightarrow until convergence (stability)



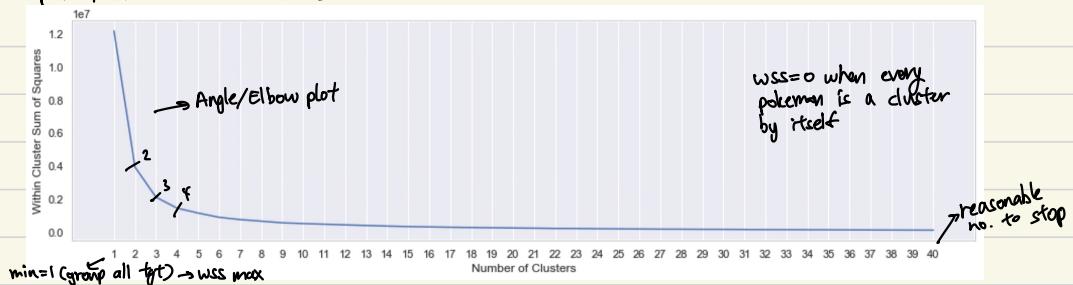
- Within Cluster Sum of Squares = 2118651

↳ sum of square of distance of each data point

in cluster from centroid ($\sum [dist(c_i - c)]^2 + \sum [dist(c_2 - c)]^2$)

Features	Total	Speed
Cluster 0:	305.65	49.36
Cluster 1:	622.57	97.08
Cluster 2:	474.27	73.55

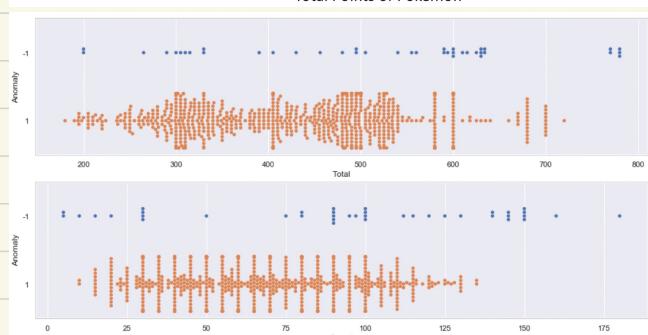
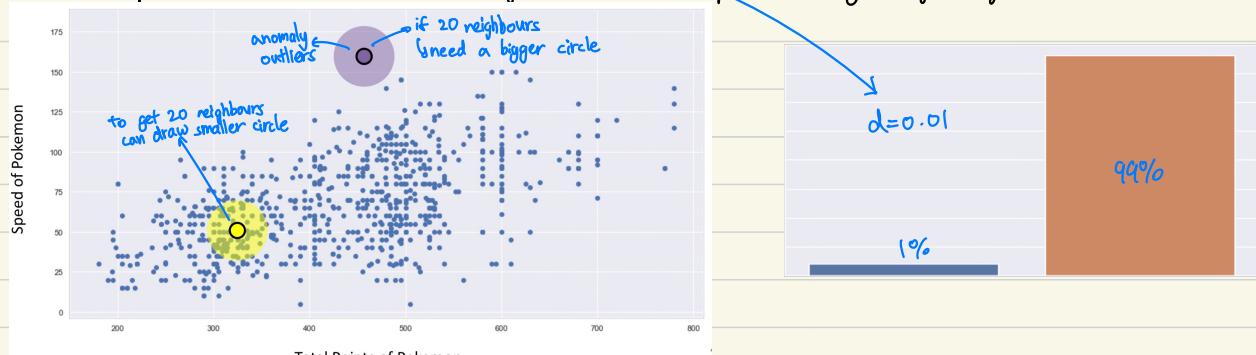
Optimal Number of Clusters



↳ choose no. of clusters by where it last bent \rightarrow 4

Anomaly Detection (Nearest Neighbor)

1. Choose $K \rightarrow$ total no. of neighbors \rightarrow if K is 1 \rightarrow emphasizes local anomalies within data, even if they are not at the boundary.
2. Choose $d \rightarrow$ fraction of anomalies in data
3. for each point in dataset \rightarrow find K nearest neighbors in data & compute if density is high enough



↳ swarm plot

Views of AI

cognitive modeling

laws of thought

rational behavior: doing the right thing

↳ falls into 4 categories: Thinking Humanly, Acting Humanly, Thinking Rationally, Acting Rationally

Turing Test → ask question to both machine & human & try to determine them

Agent

→ entity that perceives through sensors (eyes, ears, cameras, infrared range sensors)

→ acts through effectors (hands, legs, motor)

→ perceives environment → action by effectors → perceives environment

1. Rational Agents

→ do right thing (rational action → maximise expected value of an objective performance measure given percept sequence so far)

→ rationality depends on: performance measure, everything agent perceived so far, built-in knowledge about environments, actions that can be performed

2. Autonomous Agents

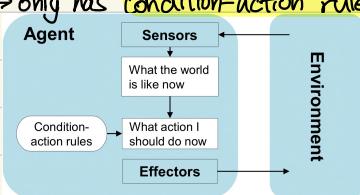
→ do not rely entirely on built-in knowledge about the environment → agent only operates successfully

when built-in knowledge are correct

→ adapt to environment through experience

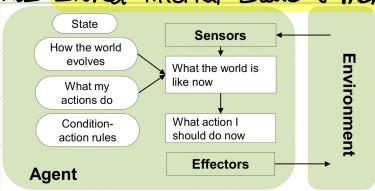
3. Simple Reflex Agents

↳ only has condition-action rules (if car in front brakes then initiate braking)



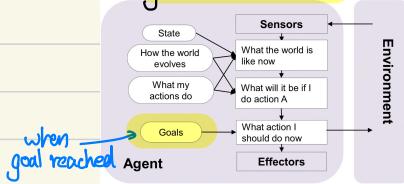
4. Reflex Agents with State

↳ has stored internal state & memory (if yesterday at NTV & no traffic jam now then go Orchard)



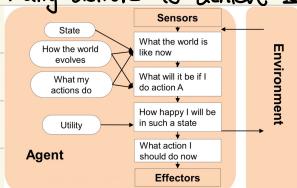
5. Goal-Based Agents

↳ has goal information



6. Utility-Based Agents

↳ many actions to achieve some goal → how happy agent will be if it attains a certain state? ⇒ utility



Types of Environments

1. Accessible (knowledge of environment) vs Inaccessible (partial/no knowledge)
2. Deterministic (outcome controlled by action) vs Non-deterministic (outcome cannot be controlled by action)
3. Episodic (previous action does not affect the next) vs Sequential (previous action affect the next)
4. Static (environment does not change when agent is deliberating) vs Dynamic (environment changes when agent deliberating)
5. Discrete (state & action integer) vs Continuous (going left → to right etc.)

E.g. Chess → accessible (all position in chessboard can be observed), deterministic (outcome of each movement can be determined), sequential (action depends on previous movements), static (no clock, considering moves, opponent can't move), discrete (all positions & movements in discrete domain)

Design of Problem-Solving Agent

1. Goal formulation
2. Problem formulation
3. Search process (no knowledge → uninformed search, knowledge → informed search)
4. Action execution

Uninformed search strategies

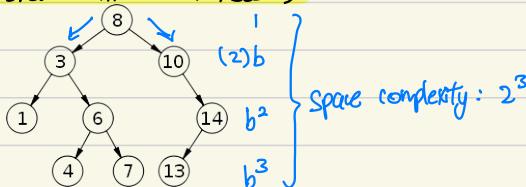
↳ use only information available in problem

Criterion	Breadth-first	Uniform-Cost	Depth-First	Depth-Limited	Iterative Deepening	Bidirectional (if applicable)
Time	b^d	b^d	b^m	b^l	b^d	$b^{d/2}$
Space	b^d	b^d	bm	bl	bd	$b^{d/2}$
Optimal	Yes	Yes	No	No	Yes	Yes
Complete	Yes	Yes	No	Yes, if $l \geq d$	Yes	Yes

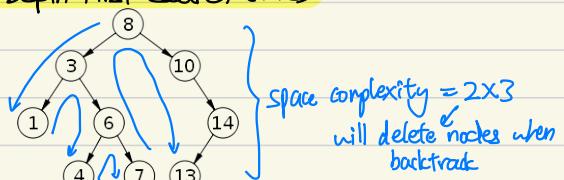
cutoff on max depth
only expand another depth when goal not found on current depth

b : branching factor, d/m : max depth (goal depth)

1. Breadth-first search (BFS)

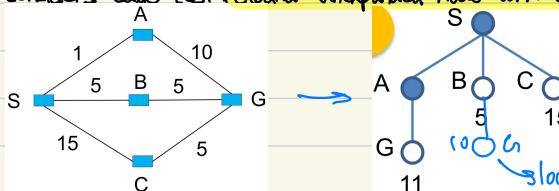


3. Depth-First Search (DFS)



2. Uniform-Cost Search

↳ considers edge cost. expand unexpanded node with least path cost, weighted graph



→ sometimes they give weighted graph, but ask for PURE BFS/DFS
slowest cost

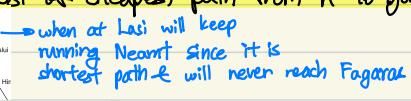
Informed search

→ use problem-specific knowledge to decide the order of node expansion

→ Uniform-cost search → path-cost function $g(n)$ → no information on cost towards goal

→ Greedy search → "heuristic" function $h(n)$ → estimate cost of cheapest path from n to goal state each round

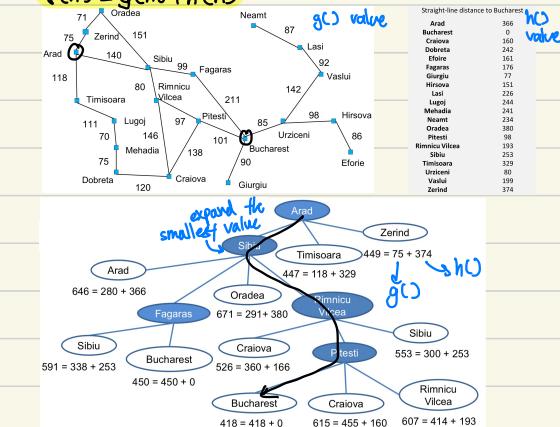
Complete	No
Time	$O(b^m)$
Space	$O(b^m)$ (keeps all nodes in memory)
Optimal	No



when at Lasi will keep running Neamt since it is shortest path & will never reach Fagaras

A* Search

$$f(n) = g(n) + h(n)$$

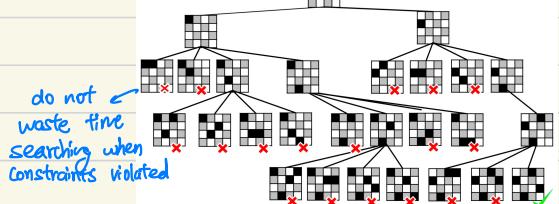


Time	Exponential in length of solution
Space	(all generated nodes are kept in memory) Exponential in length of solution

} with a good $h(n)$, significant savings are possible compared to uninformed search methods

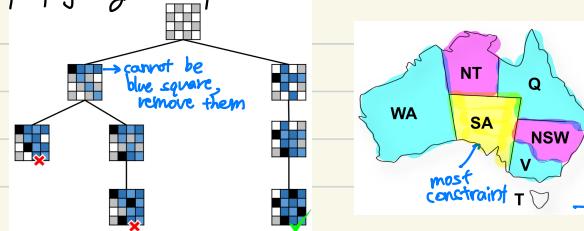
Constraint Satisfaction

1. Backtracking: Depth first search + constraint checking



2. Forward Checking & Constraint Propagation

→ propagating the implications of a constraint on one variable onto other variables

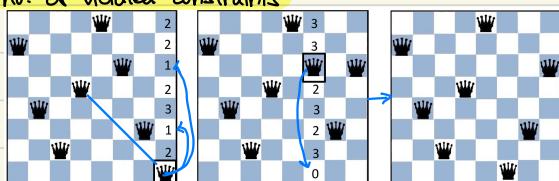


→ set the color of most constraint one first & go down list

3. Min-Conflicts Heuristic

→ given initial assignment, selects a variable in scope of violated constraint & assigns it to value that minimises

no. of violated constraints

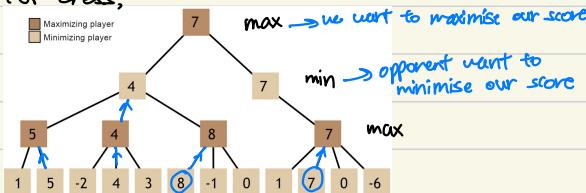


Minimax Search Strategy

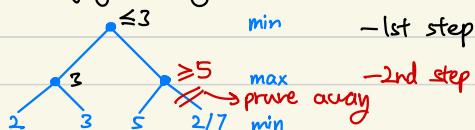
- minimax search happens everytime the agent make decision → find a sequence of moves that leads to goal
- maximise one's own utility & minimise opponent's. Assume opponent also does the same
- 3 steps:

1. Generate entire game down to terminal states
2. Calculate utility → determine best utility at terminal state up to root
3. Select best move (highest utility value)

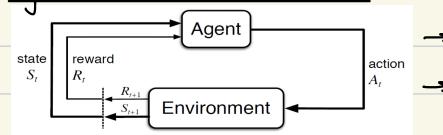
E.g. For chess,



- Game Playing = Pruning (depth-first search)



Agent-Environment Interface



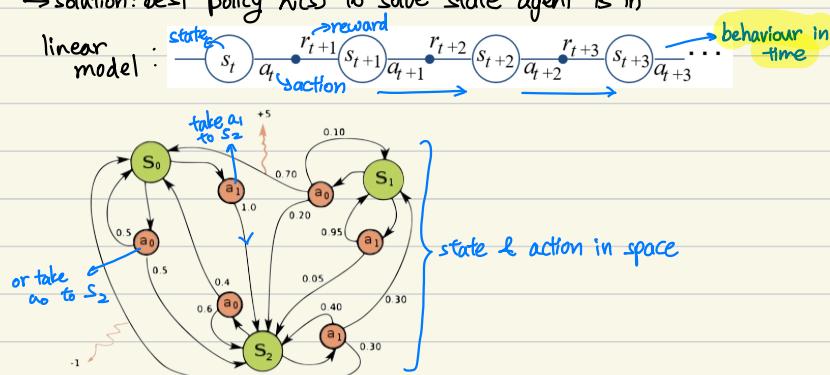
→ state, action, reward

→ maximise reward received as agent operate in environment (start to end)

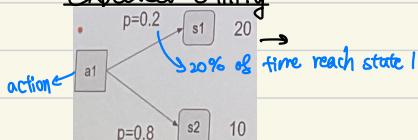
Markov Decision Process

→ components: states (s), actions (a), transition model ($P(s'|s, a)$), Reward function ($R(s)$).

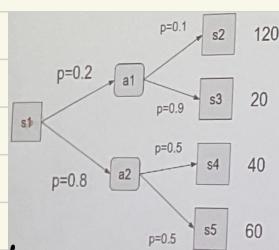
→ solution: best policy $\pi(s)$ to solve state agent is in



Expected Utility



$$\text{Expected utility of action "a1"} = 20 \times 0.2 + 0.8 \times 10 = 12$$



$$\text{Expected utility of state "s1"} = 46 \text{ (probabilistic)}$$

→ deterministic policy

$$a1 \text{ vs } a2$$

$$12 + 18 = 30$$

$$20 + 30 = 50$$

$$30 \text{ vs } 50$$

give more expected utility → take a2

State sequence Utility / Discounted Utility

→ discount the individual state rewards by a factor $0 < r < 1$

→ sooner rewards count more than later rewards → r^t will decrease to 0 as $t \uparrow$

$$U(S_0, S_1, S_2, \dots) = R(S_0) + rR(S_1) + r^2R(S_2) + \dots$$

$$\cancel{U} = \sum_{t=0}^{\infty} r^t R(S_t) \leq \frac{R_{\max}}{1-r} \quad (0 < r < 1)$$

memorise

* **Bellman Equation**

current state utility $\rightarrow U(s) = R(s) + r \max_{a \in A(s)} \sum_{s' \in S} P(s'|s, a) U(s')$

current state reward
state now
next state utility
best action of a_1, a_2
so → for a_2

$$41 = 1 + 0.8(50)$$

Value Iteration

→ start out with $U(s) = 0 \rightarrow$ iterate until convergence



Greedy policy
(pick state with highest utility)

→ apply Bellman eqn multiple times

Policy Iteration

→ Policy evaluation: calculate $V^{\pi_i}(s)$ for every state s

→ Policy improvement: calculate new policy π_{i+1} based on updated utilities

$$\pi_i^{i+1}(s) = \arg \max_{a \in A(s)} \sum_{s' \in S} P(s'|s, a) V^{\pi_i}(s')$$

} repeat

Temporal difference (TD) Prediction

→ for a given policy $\pi \rightarrow$ compute state-value function V^π , target: an estimate of return

$$TD(0): V(s_t) \leftarrow V(s_t) + \alpha [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$$

$s_t \rightarrow s_{t+1}, r_{t+1}$

→ do not require model of environment, only experience

→ can learn before knowing final outcome (less memory, peak computation)

→ can learn without final outcome (from incomplete sequences)

- In which of the following ways can you interpret LinkedIn Data as a Network (or graph).

↳ Accounts as nodes & Contacts (yes/no) as links. ✓

↳ Accounts as nodes & common Experience as links ✓

↳ Accounts as nodes & number? of individual contacts as links. X

Q need be careful to see if question got ~~say~~ say if it is a normal distribution. ↗ if not current use SD 68.27%

- Which ones of the following are decent choices for Cost Function in Linear Regression? Multiple options may be correct.

1. $\sum_i (\text{Response} - \text{Prediction})^2 \rightarrow$ sum square of errors ✓

2. $\sum_i |\text{Response} - \text{Prediction}| \rightarrow$ residual sum of squares over train set ✓

3. $\min \sum_i |\text{Response} - \text{Prediction}|^3 \rightarrow$ absolute sum of errors

4. $\max_i (\text{Response} - \text{Prediction})^2 \rightarrow$ residual absolute sum over train set

don't wanna get max/min

of train data.

- What is the relation b/w outliers of two Uni-Variate datasets & the anomalies in their joint Bi-Variate data set? That is, what is the relation b/w outliers in the boxplot of X & Y with anomalies in jointplot of $X-Y$?

Point not tagged as outliers in individual variables X & Y may be tagged anomalies in the $X-Y$ case.