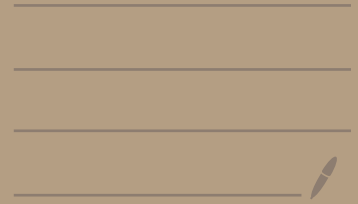


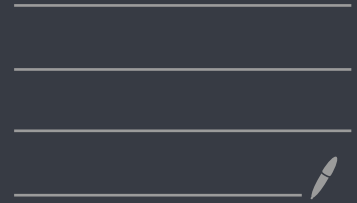
# CISC 352

---



# Unit 1

---



# Week 1: AI Origin, Agent Types and Basic Search

## Rational Decisions

- goals are expressed as utility of outcomes
  - ↳ maximizing expected utility could be a danger to other things
- So constraints are very important
- Rational agents are only good at one utility

## Heuristic Searching

- rule of thumb, helps make good decisions faster
- changes how fast you find an answer

{ which var to assign next?  
  which to try first?

- Reflex Agents react only to current state
- planning Agents have a goal, think ahead (search algorithms)

## Search Problems

- Initial state
- Actions / successor functions
- Transition model
- goal test
- path cost (optional)

## State Space

- set of all possible states
- can be huge
- can have cycles
- Search trees represent paths to a node

# CSP's

- Constraint satisfaction Problems
  - ↳ must assign values to variables
  - ↳ all constraints must be satisfied
  - ↳ There is no "best" solution. (Just valid or invalid)
- Defined by three things:
  - 1) Variables  $x, y, z$
  - 2) Domains  $x \in \{1, 2, 3\}$
  - 3) Constraints  $x \neq y$

## Examples

- Class scheduling
- Sudoku

CSP's are smart pruning huge search trees

## Search Problems

- AI explores all possibilities
- Each node = partial assignment
- Root = nothing assigned
- Leaf = all variables assigned

Backtracking  $\rightarrow$  naive DFS

$\hookrightarrow$  one variable at a time

$\hookrightarrow$  finds issue earlier by backtracking to find problem

$\hookrightarrow$  check constraints as you go (incremental goal test)

DFS with 2 improvements

$\rightarrow$  Extremely Powerful

Backtracking pseudocode \*

Improving Backtracking:

- Ordering

- Filtering

Filtering

• forward checking: cross off values of all possible affected

• variable with domain size 1, you choose that one, when you see no possible unassigned value, then stop backtrack, try again.

ARC Consistency: only cutting from tail

when you check an arc and things are good, then you change smth else, you may need to check it again. Each change is independent.

ARC - Preprocessor

ARC is more checking goal / rules

(does not detect soon to be failures)

Forward Check: + ARC consistency

Lot of relationship between  $k \rightarrow \text{node}(n)$

Must know  $k$  for ARC consistency.

# CSP Quiz Review

- goals are expressed in terms of utility (numbers of how good each outcome is)
- maximize expected utility, choose best average outcome.
- Rationale = perfectly optimized for goal, does not care about side effects.
- Rationale agent has perception and environment.
- Agents have sensors called percepts and actuators make actions
- Heuristic Search: iteratively improving solutions based on heuristic function
- Bayesian Networks: Probabilistic model representing variables

## Search

- Reflexive Agent (based on current percepts) does not care about consequences
- Planning Agents (considers future consequences, "what if" world "would look like".
- Search problem:
  - State space  $\rightarrow$  set of all possible configs
  - Successor function  $\rightarrow$  available actions and their cost
  - Start state and goal state  $\rightarrow$  Start condition and tests to try goal
  - solution is a sequence of actions that transforms start state to goal state

## State Abstraction

- World state vs State space
- includes details about environment - abstraction with only details needed for planning
- if world states are independent, can multiply the state of each object.
- State space graph: each state occurs only once  $\rightarrow$  only one
- Search tree: "what if" represents planning and their outcome  $\rightarrow$  multiple

Key Distinction: single state can appear multiple times, because they represent a distinct path, search trees can be infinite.

- Search trees can have cycles (infinite loop)
- Search tree is generated from state space graph, expands paths

## General Tree Search

- Expansion → expand potential options
- Fringe: collection of partial plans under consideration for expansion (chopping block)
- Exploration Strategy / heuristic
- Note: failure happens when no more nodes on fringe.

## Constraint Satisfaction Problems

- Find valid state **not specific path**
- Defined by: Variables, Domains, Constraints
- Standard Searches: concern about path to goal, CSPs: only care about satisfying all constraints and solutions are found at similar levels.
- State is assignment of values to variables

Partial ↙ ↘ complete

Arcs are edges = constraints between variables  
(tells you which nodes depend on each other)

## Varieties of CSP's and constraints

- Variables can be discrete/continuous
- Unary constraint: only one var involved
- Binary cons.: relates two vars
- High order: 3 or more
- Soft cons.: preferences may be better (more optimal)

## Backtracking Search:

- DFS: one var at a time
- order of assignment does not change final state
- Incremental Constraint Checking

## Techniques to help Backtracking

- Forward checking: keeps track of remaining valid domains for unassigned variables  
↳ does not provide early detection for all failures
- ARC Consistency: An arc  $x \rightarrow y$  is consistent, if for every  $x$ , there at least one  $y$ , which satisfies constraint. If value is removed, recalibrates for all neighbors.  
↳ can be run as a pre-processor.
- Backtracking is skeleton: AC-3 / FC are plug ins
- \*Try to understand code
- K-Consistency: generalized to larger groups of nodes, 1 consistency  $\rightarrow$  node level, 2 consistency  $\rightarrow$  arc level,  $k$ -consistency  $\rightarrow$   $k-1$ . Extendable  
 $k=2 \rightarrow$  Arc Consistency: • detects multivariable conflicts earlier.
- Branching factor in CSP's are size of domain of the var you are assigning
- Which variable you assign changes branching factor early.  
(matters when domain are large)

---

## GAC (generalized Arc Consistency)

- any assignment to node can be extended to all other var's sharing constraints. Involving hyper arcs.
- Hyper Arc: connects one variable to constraint involving multiple variables
- DFS: Expand deepest node first using Lifo stack for fringe  
↳ not optimal, space advantage  $O(bm)$   $b$ -branching factor,  $m$ -max depth
- BFS: Expands most shallow node first, space complexity of  $O(b^s)$  where  $s$  is depth
- Iterative Deepening: combines space efficiency of DFS with optimality of BFS. By increasing depth limit.
- Uniform Cost search: Expand cheapest node using Priority Queue.  
Both complete and optimal for positive Arc costs.
- Heuristic deciding how to move to next node.

# GAC - Pseudocode

Prunings: [ ]  
queue: [ ]

if newVar is True:

queue = CSP.get all constraints

else:

queue = CSP.get var and constraints

for Cons in queue

for Val in Cons

if Val has Support:

Pass

else: Val.prune

Prunings.append(Val)

if var.curr domain size == 0:  
return False, Prunings

for Cons in CSP.get Cons with var

if Cons not in Que  
Que.append

return True, Prunings

Value  $V$  in var  $x$  is GAC-consistent w.r.t  $C$  iff there exists at least one combination of values

- DFS > BFS, when solution is very deep (low memory)
- BFS > DFS, solution is shallow, shorter path needed
- Iterative deepening runs DFS, with depth limit

**Uniform Cost Search:** expands node with lowest path cost so far. Uses priority queue ordered by  $g(n)$  = cost from start to  $n$ .

- Effective depth of  $C^*/\epsilon$ . Takes  $O(b^{C^*/\epsilon})$ .  $C^*$  is cost,  $\epsilon$  is smallest step.  
How many cost layers before optimal solution.

• Greedy Search: expand node that seems closest to heuristic, not optimal and can be led astray.  $h(n)$  = estimated distance to goal. "what looks best"

•  $A^*$  Search: combines actual cost so far and estimated cost to goal  
 $f(n) = g(n) + h(n)$ ,  $g(n)$  is backward cost and  $h(n)$  is forward cost.

Only stop when goal is derived from fringe.

$h$  - remaining cost to goal  $A^*$  is optimal when heuristic is admissible

# Heuristic Properties

- Heuristic is **admissible** if it is optimistic. Does not over estimate.

↳  $h(n) \leq h^*(n)$ , where  $h^*(n)$  is true cheapest cost.

$h(n)$  is remaining cost from node  $n$  to goal

$h^*(n)$  is true cheapest cost to goal.

- Consistency: Heuristic is consistent if estimated cost is less than or equal to  $f$  value **never decreases along a path.**

Consistent  $\Rightarrow$  Admissible

- Relaxed Problems: easier versions of problems.

• As you move closer to goal, your heuristic estimate cannot drop by more than the cost of that step.  $f(n) = g(n) + h(n)$  never decreases along path.

- Will never find cheaper path later.

Slide 208-215

• UCS expand equally,  $A^*$  towards goal.

• all search algos are same with diff fringe strategies

---

# Graph Search

• Tree search explores paths without remembering visited states

• Graph search remembers visited states  $\rightarrow$  avoids cycles

$\rightarrow$  same tree, can appear multiple times

•  $A^*$  Tree search - Admissible heuristic is enough

•  $A^*$  Graph search - requires consistency for optimality and efficiency.

• Tree Search Pseudocode

↳ Starts from initial state

for every

Pick node from fringe

check if its goal

expand children

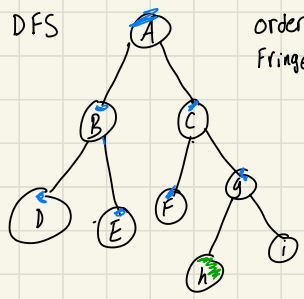
add children

• Does not remember visited states

## Graph Search

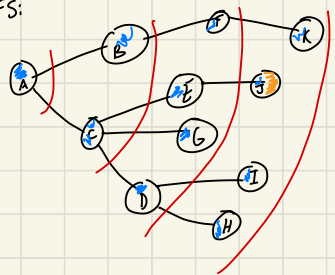
- once state expands, its never expanded again. if state expands too early with sub optimal cost

1.) DFS



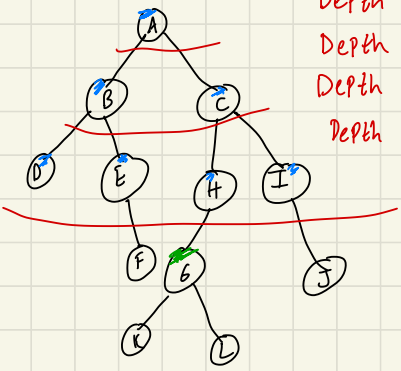
order of nodes:  $\{A, B, D, E, C, F, g, h\}$   
 Fringe:  $\{B, C\}, \{D, E, C\}, \{F, G\}, \{h, i\}$

2.) BFS:



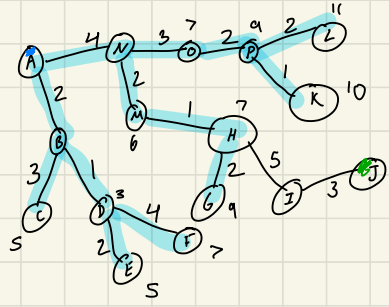
Nodes visited:  $\{A, C, B, D, G, E, F, H, I, J\}$   
 Fringe nodes:  $\{C, B\}, \{D, G, E, B\}, \{H, I, G, E, B\}, \{J, B\}$

3.) Iterative Deepening



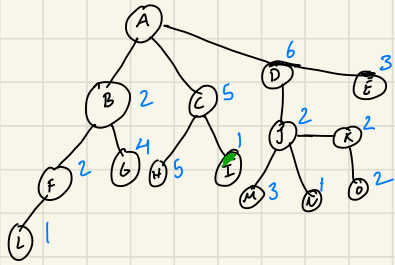
Depth Level: L1 : Nodes visited:  $\{A\}$  goal not found  
 Depth Level 2: Nodes visited:  $\{A, B\}$  goal not found  
 Depth Level 3: Nodes visited:  $\{A, B, D, E, C, H, I\}$  goal not found  
 Depth Level 4: Nodes visited:  $\{A, B, D, E, F, C, H, G\}$  goal found

4. Uniform Cost Search



Expanding Cheapest Path order:  
 of the fringe, PICK Cheapest:  
 Nodes visited:  $\{A, B, D, G, C, E, H, F, H, D, P, G, K, L, I, J\}$   
 • USES Cumulative Sum

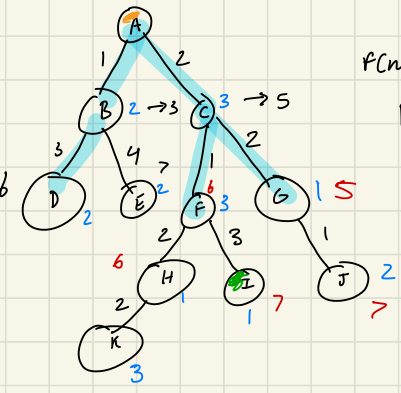
### 5. Greedy Search



$f(n) = g(n)$ , ignores path  
 always choose node that  
 looks closest to a goal.

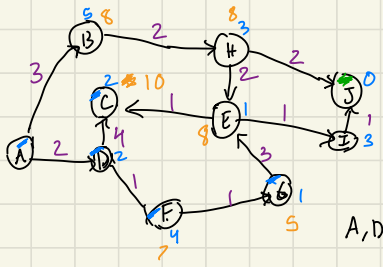
Visited nodes =  $\{B, F, L, G, C, H, I\}$   
 Fringe after G expansion  
 $\{C, D\}$

### 6. A\* Search



Nodes visited:  $\{A, B, C, D, F, H, E, I\}$   
 $f(n) = g(n) + h(n)$   
 Fringe:  $\{G, J, K\}$

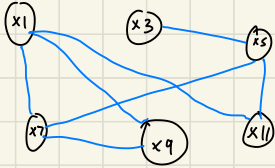
A\*  $f(n) = g(n) + h(n)$



A, D, F, G, B, E

Vars:  $\{x_1, x_3, x_5, x_7, x_9, x_{11}\}$

Cons:  $x_3 = x_5, x_5 \neq x_7$   
 $x_7 \neq x_9, x_9 \neq x_1$   
 $x_7 = 2x_1, x_{11} = x_1$   
 $x_5 \neq x_{11}$



Arc consistency:

- 1.) Empty queue, get all constraints
- 2.) take first constraint, pop it
- 3.) for every var in cons  
for every value in Var.CurDomain(c)  
if value has no support  
prune.append  
return false  
else:  
pass  
return true, prunings

PropFC / GAC - use Heuristics MRV and DH

binary\_grid / nary\_ad-grid / cagey\_csp\_model

- Transform Cagey  $\rightarrow$  CSP
- binary\_grid
  - $\hookrightarrow$  each row: all values diff
  - $\hookrightarrow$  each col: all val diff
  - double for loop add constraint
- nary\_ad\_grid:
  - One n-ary all-diff constraint
  - get scope of vars in that row check tuples adding constraint
- Cagey\_csp\_model (cagey\_grid)
  - includes col/row all-diff constraint
    - $\hookrightarrow$  cage constraints and special cages

# Unit 2

---



# Unit 2 - Planning

• adversarial search: how agents make decisions in competitive games

## Minimax Algorithm

- Zero sum games: agents have opposite utilities, one player's gain is other's loss.
- General games: independent utilities where cooperation is possible.
- **Information and Certainty**: games can be deterministic or stochastic (random)
  - ↳ may offer perfect information (can see all state details)
- could involve one player, two, two +

## To solve deterministic game with terminal utilities

- States (S) all possible configs
- Players (P): usually two taking turns
- Actions (A): may change based on state or player
- Transition Function:  $(S \times A \rightarrow S)$
- Terminal test  $(S \rightarrow \{t, f\})$ : Determines if game has ended — is game over? → once game over, numerical outcome
- Terminal utilities  $(S \times P \rightarrow R)$ : Final score or value assigned to players at end
- Solution: find policy  $(S \rightarrow A)$  mapping every state to best possible action
  - ↳ best move under worst case opponent response.

## Minimax Algorithm

- Core method for solving adversarial games: Minimax Search
  - ↳ operates on state space search tree, players alternate
- One player tries to reach state with highest possible utility, while opponent (M.M.) tries to minimize.
- Minimax value is best achievable utility **against rational adversary**?
- Recursive structure: value of state computed by looking at terminal states. Max value function chooses from its children's min results.

```
def max-value(state):  
    initialize v = -∞  
    for each successor of state:  
        v = max(v, min-value(successor))  
    return v
```

$$V(s) = \max_{s' \in \text{successors}(s)} V(s')$$



```
def min-value(state):  
    initialize v = +∞  
    for each successor of state:  
        v = min(v, max-value(successor))  
    return v
```

$$V(s') = \min_{s \in \text{successors}(s')} V(s)$$

Dispatch?

## Efficiency and Pruning

- Optimal against perfect player, but computationally expensive
- Similar complexity to DFS,  $O(b^m)$ , s.c of  $O(bm)$ , where  $b = \text{branching factor}$   
 $m = \text{max depth}$
- game tree pruning to cut off search space
- Do few minimax examples

## Alpha-Beta Pruning

- Optimization of minimax algorithm, ignores branches that can't influence final decision
- Logic: holds 2 values.  $\alpha$  (max best guaranteed value) and  $\beta$  (min best guaranteed val)  
↳ if value makes it worse than previously explored option, prune it.
- Does not change result at the end but can reduce time to  $O(b^{m/2})$
- AI uses depth limited search.?
- Evaluation Function: because search does not reach terminal states.  
↳ score goodness of current board position
- Depth tradeoff: deeper search makes quality of evaluation less critical, can't see long term benefit.

## Intro to Planning

- Planning differs from search using formal indep language (S.T.R.I.P.S)
- Key Dimension: categorized by determinism (deterministic vs probabilistic)
- Observability (full vs partial) and multi-agent nature (collab vs adversarial)
- Moving basic formalism to complex real world applications.

# Automated Planning

- STRIPS methodology: represent problem using tuple  $\langle F, A, I, G \rangle$ 
  - F: Fluents - set of variables representing, what can be true or false in world
  - A: Actions - set of maneuvers an agent can perform
  - I: Initial State - fluents that are true at start
  - G: Goal State: - specific fluents agent intends to make true
- State represents set of fluents currently true
  - Complete state: all other fluents false
  - Partial state: does not matter if other fluents T/F

## Action Definition

- a is defined by three components
  - Preconditions:  $(PRE(a))$ : fluents that must hold for action to be executable
  - Delete List:  $(DEL(a))$ : fluents removed from state after action
  - Add List:  $(ADD(a))$ : fluents being added to state

Progression: moving from one state to next

Resulting state is calculated by taking current state, removing delete list and joining it with add list:  $Progress(s, a) = (s \setminus DEL(a)) \cup ADD(a)$ .

## PDDL (Planning Domain Definition Language)

- Standardized language used to specify random planning problems. Splitting tasks in two file
- Domain File: general world physics: requirements, types, predicates and actions
- Problem File: Defines specific instance, including objects involved, initial state and specific goal.
- Solving via Search: Planning usually uses  $A^*$  which combines  $g(n)$  - cost to get to path and heuristic:  $h(n)$

↳ Because of massive state spaces, modern solvers rely on domain-indep heuristics rather than problem specific ones

# Advanced PDDL: Conditional Effects

- Conditional effects make action more dynamic
- **Concept:** Single action can have diff outcomes depending on current state of the world. Conditional on current state of the world.
- **Syntax:** implemented using  $(\text{when } \langle \text{condition} \rangle \langle \text{effect} \rangle)$
- Instead of switch on switch off. Have 1 switch, which changes based on if light is currently on/off.

## Search Strategies and Domain indep heuristics

- uniform cost search and greedy search  $F(n) = g(n) + h(n)$
- Key Heuristics:

$h^{\text{add}}$ : pessimistic heuristic that sums cost of achieving all individual goals. often overestimates true cost and is inadmissible

$h^{\text{max}}$ : optimistic heuristic, takes max cost required to reach any goal, underestimates the true cost, making it admissible

## Real world Applications

- greenhouse optimization:
- penetration testing
- Mega Printers

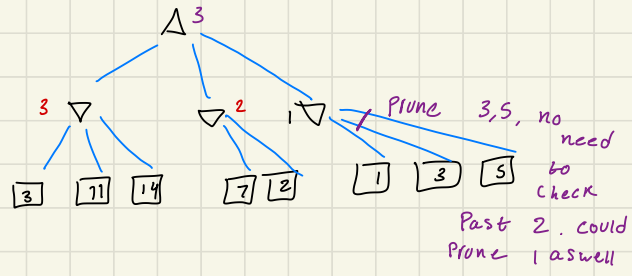
# Quiz Prep

## Minimax

- Successor of nodes Children
- if state = terminal return utility
- if next agent is max/min decide

```
def max_value(state):
    v = -∞
    for each possible option
        v = max(v, value(child))
    return v
```

## Minimax



Alpha-Beta - skips branches that can't affect decision.

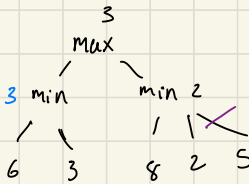
$\alpha$  - best value max can guarantee so far

$\beta$  - best value min can guarantee so far

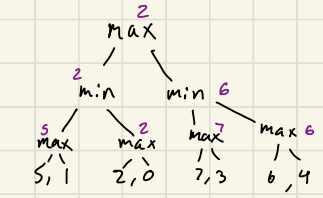
if  $\alpha > \beta$ , Prune

$\alpha = 2$   
 $\beta = -\infty$

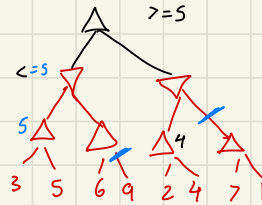
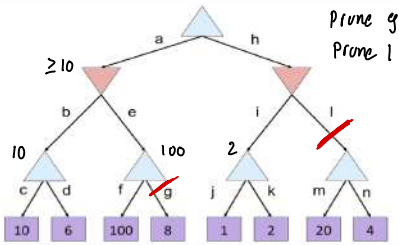
$2 \geq 7 \times$   
 $2 \geq 6 \times$



$(+\infty, 6)$   
 $\beta = 8$   
 $3 \geq \times$   
 $\dots$   
 $(8, 2)$   
 $\beta = 2$   
 $\alpha = 3$   
 $3 \geq 2 \checkmark$



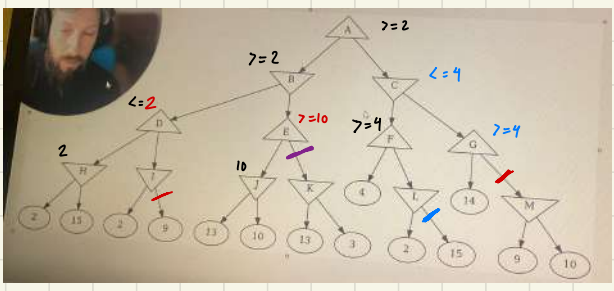
## Alpha-Beta Quiz 2



\* Practise more later



Q.8) given english convert to STRIPS/PDDL  
 Q.9) Do Reverse, convert PDDL → English



a1: T  
 a2: F  
 a3: T  
 a1, a2

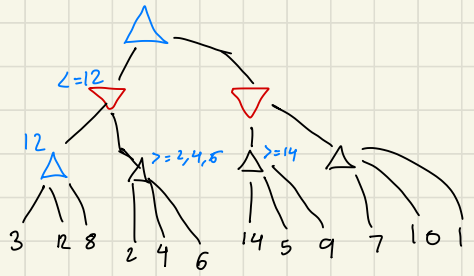
x1 is applicable  
 b, c  
 c, d yes valid

(: action register  
 : Parameters (s, c)  
 : pre (and (open c) (not (registered s, c))  
 )  
 : effect  
 (and (not (open c)) (registered s, c))  
 )  
 )

- 1.) v1, v2 are valid
- 2.) p, r, v2 are valid
- 3.) no valid plan

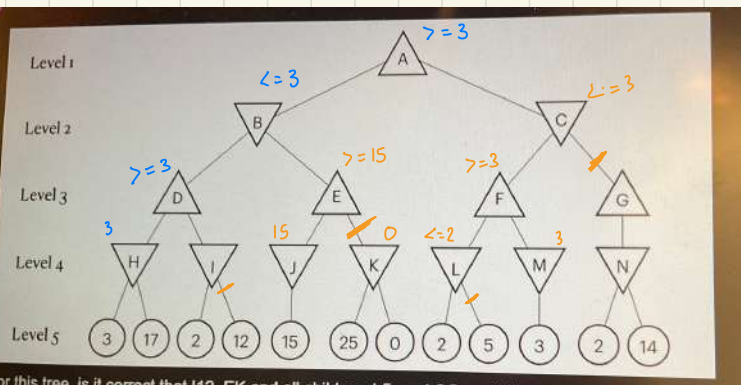
In state S, where  $S = \{a, b, d\}$   
 we execute a4, so  
 $S = \{a, c, d\}$   
 $S_2 = \{a, c, e\}$   
 $S_3 = \{a\}$

w1 is applicable  
 $\{k, \text{door open}\}$   
 $\{\text{door open}, k, \text{inside}\}$   
 $\{\text{inside}, k\}$



Practise Heuristic Tree + Alpha beta  
 + examples w words

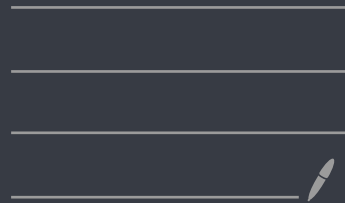
Review hmax/hadt stuff  
 Alpha beta pruning



For this tree, is it correct that H, E, K and all children, L5, and CG and all children would be pruned?

# Unit 3

---



# Unit 3: Probabilistic Inference

- Managing uncertainty: agents have observed variables (things it knows) and unobserved variables (unknown environmental factors: location of ghost)
- Probabilistic reasoning is framework to guess what we don't know, with what we do know.
- Random Variables and Distributions:
  - ↳ represent aspects of world
  - ↳ table that lists every value a variable can take and assigns a probability.
- Joint Distribution: larger table, covers multiple variables, gives every possible combination of outcomes. All probabilities must sum to 1 (Probability of  $x$  and  $y$  at same time)
- Marginalization: Big table, many variables, but only care about 1, collapse by adding probabilities for variables you want. (add rows that match what you want)

## Conditional Probabilities

- Simple relation between joint and conditional probabilities.

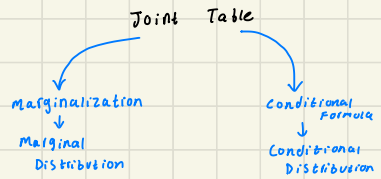
$$P(x|y) = \frac{P(x,y)}{P(y)}$$

Quiz practise:

$$P(x|y) = \frac{0.2}{0.6} = 0.33$$

$$P(-x|y) = \frac{0.4}{0.6} = 0.66$$

$$P(y|x) = \frac{0.2}{0.5} = 0.4$$



## Normalization Trick

- Normalization = rescale probabilities so they sum to 1.

Instead of computing full formula:  $P(A|B) = \frac{P(A,B)}{P(B)}$

- select rows with matching evidence
- ignore denominator
- normalize at the end

Reads like: what's the distribution of all possible values  $x$  can take, given evidence  $y$ .

Find  $P(W|Cold) =$  rows with Cold, 0.2, 0.3

↳ normalize  $\frac{0.2}{0.5}, \frac{0.3}{0.5}$   
0.4, 0.6

$$P(W|Cold) = \frac{P(W,Cold)}{P(Cold)}$$

Distribution, Sums to 1

Avoids calculating

$$P(-x|y) = \frac{P(-x,y)}{P(y)}$$

# Probabilistic Inferencing

- inference is computing a conditional probability (what you want to know) based on evidence (what you already have)
- Evidence: variables you observed already
- Query: variable you want to find
- goal: update beliefs based on evidence

## Inference by Enumeration

- answering probability questions based on listing, summing rows from joint distribution table
- limitation:  $O(d^n)$  for big tables
- sum all possibilities that produce that evidence
- $P(Q|E) = \frac{P(Q, E)}{P(E)}$  → can use normalized version or not

•  $P(W=rain | winter) = 0.25$   
 ↳ normalize all  $P(W | winter)$   
 $= \frac{0.25}{0.50} = 0.50$

## Product Chain, Bayes Rule

- Product rule:  $P(y) \cdot P(x|y) = P(x, y)$   
 ↳ useful when given conditional distribution but want joint
- Chain rule: break down complex situation into simpler conditionals

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

$$P(x_1, x_2, x_3) = P(x_1) \cdot P(x_2|x_1) \cdot P(x_3|x_1, x_2)$$

- Bayes Rule: Two ways to factor joint distribution over 2 variables:  $P(x, y) = P(x|y) \cdot P(y) = P(y|x) \cdot P(x)$

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)}$$

### Example

$$\begin{aligned} P(w|dry) &= \\ P(sun|dry) &= P(dry|sun) \cdot P(sun) \\ &= 0.9 \cdot 0.8 \\ &= 0.72 \\ P(rain|dry) &= P(dry|rain) \cdot P(rain) \\ &= 0.3 \cdot 0.2 \\ &= 0.06 \\ P(sun|dry) &= 12/13 \rightarrow \text{normalized} \\ P(rain|dry) &= 1/13 \rightarrow \text{normalized} \end{aligned}$$

## Inference with Bayes' Rule

- Example: Diagnostic probability from causal probability:

$$P(\text{cause}|\text{effect}) = \frac{P(\text{effect}|\text{cause})P(\text{cause})}{P(\text{effect})}$$

- Example:

• M: meningitis, S: stiff neck

$$\begin{aligned} P(+m) &= 0.0001 \\ P(+s | +m) &= 0.8 \\ P(+s | -m) &= 0.01 \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Example gives}$$

$$P(+m | +s) = \frac{P(+s | +m)P(+m)}{P(+s)} = \frac{P(+s | +m)P(+m)}{P(+s | +m)P(+m) + P(+s | -m)P(-m)} = \frac{0.8 \times 0.0001}{0.8 \times 0.0001 + 0.01 \times 0.9999}$$

all possible ways to get P(+s)

$$= 0.0079$$

• Note: posterior probability of meningitis still very small

• Note: you should still get stiff necks checked out! Why?

Can calculate updated probabilities using causal probability

# Independence

- Two variables don't affect each other

$$\hookrightarrow P(x, y) = P(x) \cdot P(y)$$

- Conditional Independence: Two events are only related because of a third event.  
 $\hookrightarrow$  once you know third event, first two don't know anything new about each other.

$X \perp\!\!\!\perp Y \mid Z$  X is conditionally independent of Y given Z.

# Bayes Net

- describes complex relationships using simple, local interactions.
- nodes are variables
- arcs indicate a direct influence between variables
- $R \rightarrow T$ , rain affects traffic

Represent uncertainty  
Instead of guessing  
Compute probabilities

# Conditional Probability Tables

- $\rightarrow$  if no parent  $\rightarrow$  simple probability
- $\rightarrow$  if has parents  $\rightarrow$  conditional prob

## To Build Bayes Net

- 1) Identify variables
- 2) find causes, connect
- 3) Graph, Tables

# For independence

$$R \rightarrow T \rightarrow L$$

Late does not depend on rain, simplifying our calculation  $P(L \mid T, R) = P(L \mid T)$

$P(\text{all variables}) = \text{Product of each node given its parents}$

# Factor Zoo ?

x and y are indep if  $P(x) \cdot P(y) = P(x, y)$

x and y are conditionally indep if given z

is  $x, y \perp\!\!\!\perp z$  if

$$P(a, b) = P(b \mid a) \cdot P(a)$$

$$P(a \mid b) = \frac{P(a, b)}{P(b)}$$

$$P(a \mid b) = \frac{P(b \mid a) \cdot P(a)}{P(b)}$$

$$P(+t) =$$

$$P(+w) =$$



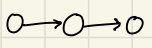
# D-Separation

- in bayes net find indep var based on graph structure
- Causal chains ( $X \rightarrow Y \rightarrow Z$ ): if you observe Y, path is blocked and x and z become indep
- Common Cause ( $X \leftarrow Y \rightarrow Z$ ): one cause has 2 effects. if you have Y, the effects become indep
- Common Effect/V-Structure ( $X \rightarrow Z \leftarrow Y$ ): Two indep causes share one effect. Causes are indep unless observe effect. Observing effect 'activates' path

## D-Separation - Algorithm

- Are x and y independent given evidence  $\{z\}$
- 1) Find all paths
- 2) check Active path (every triple)
- 3) if even one active path, indep is not guaranteed. if all paths are blocked, x and y are guaranteed to be indep
- Complexity:  $2^N \rightarrow O(N \times 2^{k+1})$

## Active Triples

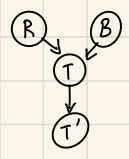


Path is active where  $A \rightarrow B \rightarrow C$ , where B is unobserved  
 common cause  $A \leftarrow B \rightarrow C$ , where B is unobserved  
 Common effect:  $A \rightarrow B \leftarrow C$ , where B is observed.

Note: if atleast one path open, not indep

EX

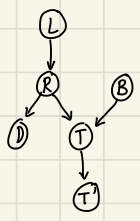
$R \perp\!\!\!\perp B$   
 $R \perp\!\!\!\perp B \mid T$   
 $R \perp\!\!\!\perp B \mid T'$



- $R \perp\!\!\!\perp B$  has  $R \rightarrow T \leftarrow B$ , collider no T, so not observed, independent
- $R \perp\!\!\!\perp B \mid T$ , same thing given T, so not indep.
- $R \perp\!\!\!\perp B \mid T'$   
 $R \rightarrow T \rightarrow T'$   
 $B \rightarrow T \rightarrow T'$   
 causal and we have  $T'$   $\therefore$  blocked so not indep

EX

- $L \perp\!\!\!\perp T' \mid T$   
 $L \rightarrow R \rightarrow T \rightarrow T'$   
 T is observed, causal blocked  
 - indep

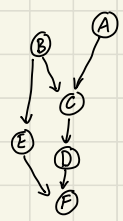


- $L \perp\!\!\!\perp B$   
 $L \rightarrow R \rightarrow T \leftarrow B$   
 - V affect, not observed, blocked  
 - indep

- $L \perp\!\!\!\perp B \mid T'$   
 $L \rightarrow R \rightarrow T \leftarrow B$   
 - V affect opens path  
 - not indep

- $L \perp\!\!\!\perp B \mid T$   
 $L \rightarrow R \rightarrow T \leftarrow B$   
 - V affect, observed, valid  
 $\rightarrow$  not indep

- $L \perp\!\!\!\perp B \mid T, R$   
 given both, blocked, indep



$A \rightarrow C \leftarrow B$   
 $A \rightarrow C \leftarrow B \rightarrow E$

# Naive Bayes'

- Predict a class  $Y$ , from features
- All features are independent given class

$$P(w_1, w_2, w_3 | \text{spam}) = P(w_1 | \text{spam}) \cdot P(w_2 | \text{spam}) \cdot P(w_3 | \text{spam})$$

easy to compute and uses  $\text{Score}(\text{class}) = P(\text{class}) \cdot \text{Product of word probabilities.}$

Ex: give 2 words in spam

$$P(\text{spam}) = 0.33$$

$$P(w_1 | \text{spam}) = 0.0156$$

$$P(w_2 | \text{spam}) = 0.0153$$

$$\left. \begin{array}{l} P(\text{spam}) = 0.33 \\ P(w_1 | \text{spam}) = 0.0156 \\ P(w_2 | \text{spam}) = 0.0153 \end{array} \right\} \text{Spam Score} = 0.33 \times 0.0156 \times 0.0153$$

$$\text{Ham Score} = 0.66 \times P(w_1 | \text{ham}) \cdot P(w_2 | \text{ham})$$

Compare Probabilities, bigger Score = predicted label

$$P(\text{positive}) = 0.55 \cdot 0.08 \cdot 0.05 \cdot 0.02 = \frac{11}{250000}$$

$$P(\text{neg}) = 0.45 \cdot 0.01 \cdot 0.03 \cdot 0.07 = 9.45 \times 10^{-6}$$

## Laplace Smoothing

- When estimating Probabilities if word is not in relevant group. Probability is 0, but instead of  $\times 0$ , add  $K$  fake observations to all category.

-  $\text{Count}(x) / N$

↳  $(\text{Count}(x) + K) / (N + K \times \text{number\_of\_categories})$

→ add  $K$  to each count, increase total accordingly.

Ex Data:  $\{r, r, b, b, g, b, b, g, r, b, g, r\}$

count:  $r=4, b=4, g=3, N=11$

So  $|X|=3$ , from 3 possible values

$$\left. \begin{array}{l} P(r) = \frac{4}{11} \\ P(b) = \frac{4}{11} \\ P(g) = \frac{3}{11} \end{array} \right\} \begin{array}{l} \text{use } K=1, \\ \text{now } \rightarrow \end{array} \begin{array}{l} P(r) = \frac{5}{14} \\ P(b) = \frac{5}{14} \\ P(g) = \frac{4}{14} \end{array}$$

Now  $Y=0$ , gives  $P(Y) = \frac{1}{(N + |X|)}$

$$\begin{array}{l} \text{red} = \frac{4}{11} \rightarrow \frac{6}{14} \\ \text{blue} = \frac{3}{11} \rightarrow \frac{5}{14} \\ \text{green} = \frac{1}{11} \rightarrow \frac{2}{14} \\ \text{yellow} = \frac{1}{11} \rightarrow \frac{2}{14} \\ \text{purple} = 0 \rightarrow \frac{2}{14} \end{array}$$

$n=9, K=2, |X|=3$

$(N + K|X|) \rightarrow \text{denominator}$

$\frac{K}{(N + K|X|)}$

# Variable Elimination with Bayes Net

Trying to find  $P(L)$

Must remove hidden variables

1) Join R Tables, multiply matching rows

+r	+t	0.08
+r	-t	0.02
-r	+t	0.09
-r	-t	0.81

2) Sum out R

+t	0.17
-t	0.83

3) Join with L

+t	+L	0.051
+t	-L	0.119
-t	+L	0.083
-t	-L	0.747

4) Eliminate T

+L	0.134
-L	0.866

Join A and B

+A	+B	0.025
+A	-B	0.012
-A	+B	0.12
-A	-B	0.056

Eliminate A

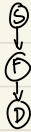
+B	0.3
-B	0.7

Join C

+B	+C	0.18
+B	-C	0.12
-B	+C	0.07
-B	-C	0.72

Sum C

+C	0.25
-C	0.9



1) Join S, F  
 $P(S, F)$

S	f	0.06
S	-f	0.24
-S	f	0.42
-S	-f	0.28

2) Sum out S

$P(f)$

+f	0.48
-f	0.52

3) combine with F, D

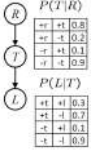
$P(f, d)$

+f	+d	0.336
+f	-d	0.144
-f	+d	0.052
-f	-d	0.468

4) Sum out f

$P(d)$

+d	0.388
-d	0.612



Join R

+r	0.1
-r	0.9

$P(T|R)$

+r	+t	0.2
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

$P(L|T)$

+t	+L	0.3
+t	-L	0.7
-t	+L	0.1
-t	-L	0.9

# Unit 4

---



# Unit 4: Deep Learning

## Decision Tree

- model that uses series of binary or multiway questions to classify data. 3 main node types
- Root node
- Internal nodes - test a specific attribute "question nodes" 2 or more branches
- Leaf Nodes: - final answers, that give class label

## How to Construct

- induction uses greedy strategy: making best decision at current split
- **Hunts Algorithm**: recursive procedure to split data into smaller subsets, until all records in subset belong to same class
- **ID3**: Creating shallowest tree by using information gain to select best feature to split on each step

## Choosing Split

- to choose which algorithm makes best split, measure **impurity** - how mixed classes are in a group
- Entropy: the uncertainty or pureness of set. 1 class  $\rightarrow$  0 entropy
- Information gain: reduction in entropy is achieved by splitting **specific feature**
- Gini Index and Misclassification Error

## Pros and Cons

- Easy interpretation
- Efficiency
- Robust
- Lower Accuracy

## The Neuron

- Perceptron (neuron) Flow:
- Inputs and Weights: receives feature values as input and assigns weight to determine importance
- Summation and Activation: sum inputs to create activation value.
- Output: passed through activation function to produce classification label

# Deep Neural Networks

- deep learning involves stacking many neurons
- Learning features: learns themselves from raw data
- Power of depth: more layers  $\rightarrow$  better generalization and higher accuracy?
- Universal Approximators: network with even 2 layers can approximate any continuous function

## Training Network

- finding best weights to maximize the probability of a correct predictions

### Training Loop:

- Forward pass  $\rightarrow$  make prediction
- Compare to true answer  $\rightarrow$  compute loss
- Backpropagation  $\rightarrow$  compute gradients
- Update weights  $\rightarrow$  gradient descent

While training may feel automated, developer chooses:

- number of layers (more layers, more complex)
- width (neurons per layer)
- Learning rate (most important)
- Epoch, batch size, regularization

## Training Techniques

- gradient ascent
- Backprop: use chain rule to calculate derivatives by passing information backward
- Early stopping to prevent overfitting

## Real-world challenges

- Bias: toxic, racist, discriminatory content
- correlation vs causation
- Robustness and Manipulation (apple labelled "iPod")

## Activation Functions

- nonlinear functions
- Sigmoid and tanh
- ReLU ( $\max(0, x)$ )
- Leaky ReLU

# Loss Functions

- Score that explains how good/bad prediction was (Try to minimize)
- Classification:
  - ↳ Cross entropy: reward confident correct answers, punish confident wrong answers hard
  - ↳ Negative log-likelihood: take probability of correct answer, punish low prob, scales nicely
- Regression:
  - ↳ L1: Mean Absolute error  $(\text{prediction} - \text{actual})$
  - ↳ L2: Mean Squared error  $(\text{prediction} - \text{actual})^2$

## Regularization - prevent overfitting

- reducing generalization by adding L1/L2: penalties to avoid complexity
- Data Augmentation: Create synthetic training data by flipping/turning/color shift in images
- Dropout: randomly turn off neurons

## Backpropagation

- figure out where error comes from and fix it
  - 1.) make prediction
  - 2.) compute loss
  - 3.) send error backward through network
  - 4.) Adjust weights reduce error

## CNN'S

- Process grid like data, same weights across different parts of input, much more efficient
- Replaces matrix multiplication with convolution, slides over input to create feature map
- Filter is group of weights, model learns these, tells us what patterns to look for
- Stride: how many pixels you jump each time, (low stride = walking, big stride = running)
- Padding: adding border to keep output consistent, and not to lose edges
- Pooling: Summarizing small regions

## RNN'S

- Used for sequential data, order matters. Reads one piece at time, combines memory, update and pass forward
- loops: allow output to get fed back in as input
- Unfolding: Multiple copies of same network, each passing message to successor
- Configuration: many to one, one to many, many to many

## Long-Short Term Memory

- RNN's struggle to remember information because of vanishing gradient. (longer it gets, harder to remember)
- Gating is a smart technique with a forget gate, input gate and output gate.
- Real world Applications: handwriting and action recognition, medical and robotics

## Autoencoders - type of neural network

- models complex data by learning how to compress and then reconstruct it.
- The encoder takes input, shrinks it into (latent space representation).
- The decoder takes that code and tries to reconstruct the original as accurately as possible
- Model is trained by minimizing difference between original and reconstructed (MSE)

## Variational Autoencoders (VAE's)

- instead of single fixed points, the encoder predicts probability distribution  
↳ gives a whole range of possible codes (more info)
- Gumbel-softmax: turns discrete categorical data into soft smooth approximation

## LatPlan

- learn how system works by watching it, then solve
- 1) Autoencodes turns raw images into simple binary code
- 2) Learns rules using PDDL
- 3) Planner finds solution in latent space, decodes back into sequence images

## Transformers

- look at everything at once, uses Attention
- Self-Attention: model looks at every word simultaneously to see most relevant
- Bert: looks in both direction
- Gpt: only reads from left → right. (can use zero-shot, few shot)

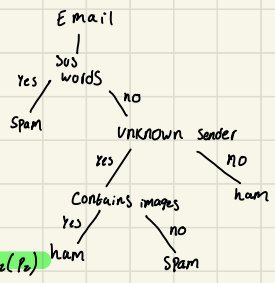
# Possible types of Questions

- 1) given Decision tree, classify test data (trivial likely will not show up)
- 2) given some data, figure out best splits for tree

Ex Each row is email: Predict ham/spam

Features: suspicious words, unknown sender, Contains images

Hunts Algorithm: If "Pure" stop, else pick feature, split, repeat.



ID3: Split at feature with most information gained.

Entropy: Calculates how mixed a dataset is.  $Entropy = -p_1 \log_2(p_1) - p_2 \log_2(p_2)$

Compute entropy for a dataset (cleaner split = higher gain)

→ if  $\frac{4}{10}$  spam,  $\frac{6}{10}$  ham  $Entropy = -0.4 \log_2(0.4) - 0.6 \log_2(0.6)$

## ID3 Method

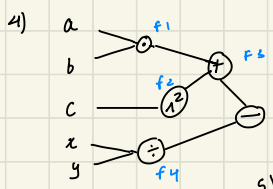
- 1) find entropy of full dataset
- 2) try each feature calculate entropy if 2 spam, 1 ham  $H(\text{Feature A}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$   
↳ calculates for both true/false.  $H_{\text{After}} = \frac{\# \text{ of T}}{\text{dataset size}} (\text{Entropy}) + \frac{\# \text{ of F}}{\text{dataset size}} (\text{Entropy})$
- 3) Information gain for that feature is data entropy -  $H_{\text{After}}$

Full entropy =  $-0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$

Feature A =  $-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \rightarrow \frac{3}{8} (0.75) +$

Feature B =

Then choose feature with most information gained



$f_1 = a \cdot b$   
 $f_2 = c^2$   
 $f_3 = f_1 + f_2$   
 $f_4 = x/y$   
 $z = f_3 - f_4$

given variables, just plug and play to compute forward pass

- 6) Calculate gradients (back pass for same model) How does changing each input affect z

Start at end

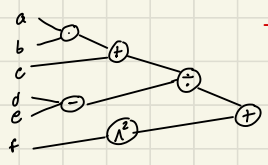
- 1)  $\frac{\partial z}{\partial z} = 1$
- 2) from  $z = f_3 - f_4$   
 $\frac{\partial z}{\partial f_3} = 1$      $\frac{\partial z}{\partial f_4} = -1$
- 3) from  $f_3 = f_1 + f_2$   
 $\frac{\partial z}{\partial f_1} = 1 \cdot 1 = 1$      $\frac{\partial z}{\partial f_2} = 1 \cdot 1 = 1$
- 4) from  $f_1 = a \cdot b$   
 $\frac{\partial z}{\partial a} = 1 \cdot b = b$      $\frac{\partial z}{\partial b} = 1 \cdot a = a$
- 5) from  $f_2 = c^2$   
 $\frac{\partial z}{\partial c} = 1 \cdot 2c = 2c$
- 6) from  $f_4 = x/y$   
 $\frac{\partial z}{\partial x} = -1 \cdot (1/y) = -1$      $\frac{\partial z}{\partial y} = -1 \cdot (-x/y^2) = +2$

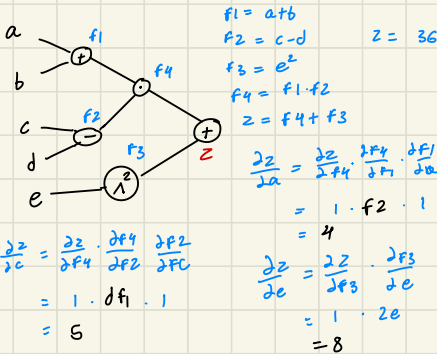
Ex  $z = \frac{(a \cdot b + c)}{(d \cdot e)} + p^2$

- 5) Compute the value z, step by step compute all intermediate terms

given  $a=5, b=4, c=3, x=2, y=1$  }  $z=27$

Chain rule:  $\frac{\partial z}{\partial a} = \frac{\partial z}{\partial f_3} \cdot \frac{\partial f_3}{\partial f_1} \cdot \frac{\partial f_1}{\partial a}$   
 $\frac{\partial z}{\partial b} = \frac{\partial z}{\partial f_3} \cdot \frac{\partial f_3}{\partial f_1} \cdot \frac{\partial f_1}{\partial b}$   
 $= 1 \cdot 1 \cdot a = a$





7.) Slide small matrix filter, over big matrix input, multiply and sum

Input: 4 by 5 matrix  
 Filter: 2 by 2 filter  
 Stride = 2 → skip 1 cell each move  
 1 → move 1 step  
 Padding = 0, no border  
 Padding = 1, zeroes around

output height =  $\left\lfloor \frac{H - F_h + 2P}{S} \right\rfloor + 1$  (round down)  
 output width =  $\left\lfloor \frac{W - F_w + 2P}{S} \right\rfloor + 1$   
 output height × output width

Ex  
 input =  $\begin{bmatrix} 1 & 2 & 3 & 0 \\ 4 & 5 & 6 & 1 \\ 7 & 8 & 9 & 2 \\ 1 & 3 & 2 & 0 \end{bmatrix}$   
 filter =  $\begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix}$   
 output =  $\begin{bmatrix} 14 & 16 \\ 12 & 13 \end{bmatrix}$

output =  $\begin{bmatrix} 6 & 8 \\ 12 & 14 \end{bmatrix}$   
 output =  $\begin{bmatrix} 10 & 16 \\ 7 & 10 \end{bmatrix}$   
 output =  $\begin{bmatrix} 6 & 12 & 8 \\ 14 & 25 & 76 \\ 12 & 21 & 14 \end{bmatrix}$

Input =  $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 2 & 3 & 4 & 5 \\ 0 & 5 & 4 & 3 & 2 & 1 \\ 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$   
 5 by 7

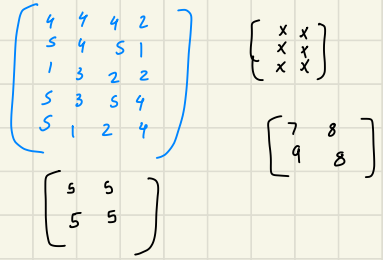
Filter  $\begin{bmatrix} 1 & 0 & 1 \\ 2 & 1 & 0 \end{bmatrix}$

Output =  $\begin{bmatrix} 1 & 4 & 7 & 10 & 13 \\ 1 & 2 & 2 & 2 & 1 \end{bmatrix}$

• CNN share weights, focus on small regions, allows pattern learning, sliding across image, receptive field: how much of image neuron sees.

8. Apply pooling function (max/sum) and dimensions, stride padding (Does not take filter matrix)

• Applies simple functions, outputs one number  
 • Given: pool size, function, stride, padding



3) run simple train process  
 Objective: correct class score > all other class score  
 Score =  $w \cdot x$   
 goal: update weight vectors, so each sentence gets classified properly

- 1) compute scores
- 2) predict class
- 3) compare with true label
- 4) if correct, nothing
- 5) if wrong update

"win the vote"  
 sports score = 1  $[1 \ 0 \ 0 \ 0] - [1 \ 1 \ 0 \ 1]$   
 for 1 =  $[0 \ -1 \ 0 \ -1]$   
 for politics =  $[1 \ 1 \ 0 \ 1]$  "win election"  
 for 1 =  $[0 \ -1 \ 0 \ -1] \cdot [1 \ 1 \ 0 \ 1]$   
 = -2  
 for politics =  $[1 \ 1 \ 0 \ 1] \cdot [1 \ 1 \ 0 \ 1]$   
 = 3

comeback to

Entropy for whole set = 1

Feature A: T: 4w, 1L

$$-(0.8 \log_2 0.8) - (0.2 \log_2 0.2)$$

Entropy = 0.72

$$f: 4L \quad 1w$$

$$= 0.72$$

$$= 0.5(0.72) + 0.5(0.72)$$

$$= 1 - 0.72 = 0.28$$

Feature B: 6 True: 4w 2L

$$T: -(-0.390) - (-0.528) = 0.918$$

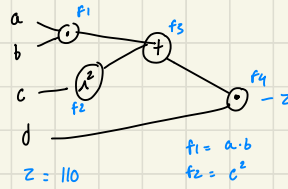
$$F: 4: 1w \quad 3L$$

$$-(-0.481) - (-0.3) = 0.781$$

$$= (0.918)(0.6) + (0.4)(0.781)$$

$$= 0.8632$$

$$= 1 - 0.8632 = 0.16$$



$$Z = 110$$

$$f_1 = a \cdot b$$

$$f_2 = c^2$$

$$f_3 = f_1 + f_2$$

$$f_4 = d$$

$$Z = f_3 \cdot f_4$$

$$\frac{\partial Z}{\partial a} = \frac{\partial Z}{\partial f_1} \cdot \frac{\partial f_1}{\partial a} = \frac{\partial Z}{\partial f_4} \cdot \frac{\partial f_3}{\partial f_1} \cdot \frac{\partial f_1}{\partial a}$$

$$= 1$$

$$\frac{\partial Z}{\partial f_3} = f_4 = 5, \quad \frac{\partial Z}{\partial f_4} = f_3 = 22$$

$$\frac{\partial Z}{\partial f_1} = \frac{\partial Z}{\partial f_3} \cdot \frac{\partial f_3}{\partial f_1} = 5 \cdot 1 = 5$$

$$\frac{\partial Z}{\partial c} = \frac{\partial Z}{\partial f_2} \cdot \frac{\partial f_2}{\partial c} = \frac{\partial Z}{\partial f_3} \cdot \frac{\partial f_3}{\partial f_2} \cdot \frac{\partial f_2}{\partial c}$$

$$= 5 \cdot 1 \cdot 2 \cdot 4$$

$$= 40$$

## Neural Nets Conceptually

Batching - sending in multiple inputs at once {5, 10, 15, 20}, update weights after every batch

Feature expansion 1 sample  $\rightarrow$  10 features  $\rightarrow$  50 features (every layer gets more features)

network makes features by doing new feature = **activation (old feature  $\times$  weight + bias)**

CNN's number of features = number of filters

After every batch, compute loss, update weights, features get smarter

$\rightarrow$  Faster computation, smooth learning  
 $\rightarrow$  Large batch = overfitting, memory limits