

Course Organization

Teacher

Thomas Dyhre Nielsen, tdn@cs.aau.dk, office 1.2.34

Machine Intelligence

Lecture 1: Introduction to MI and Agents

Thomas Dyhre Nielsen

Aalborg University



D. Poole and A. Mackworth: *Artificial Intelligence. Foundations of Computational Agents* (2nd edition)

<http://artint.info/>

Literature

- Mondays, 8:15-12:00 (typically Frøb 7H, Aud): Exercises 8:15-10:00, Lecture 10:15-12:00
- Wednesdays 8:15-12:00 (group rooms): Extended exercise sessions

Exam

Exam in January



1

MI Autumn 2019

2

Why this book?



Topics:

- Introduction
- Problem solving as search
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Machine learning
- Planning
- Reinforcement learning
- Multiagent systems

Several reasons, but:

We decided that it is better to clearly explain the foundations upon which more sophisticated techniques can be built, rather than present these more sophisticated techniques. This means that a larger gap may exist between what is covered in this book and the frontier of science. But it also means that the student will have a better foundation to understand current and future research.

Tentative course overview



MI Autumn 2019

3

4

Learning goals



After having followed the course you should

- have knowledge about basic techniques within the field of machine intelligence and computational agents.
- be able to apply key machine intelligence techniques to a specific problem domain.
- be able to reason about computational agents that can operate in domains of varying complexity.

After having followed the course you should

- have knowledge about basic techniques within the field of machine intelligence and computational agents.
- be able to apply key machine intelligence techniques to a specific problem domain.
- be able to reason about computational agents that can operate in domains of varying complexity.

... and be well-prepared for the aMI course and the machine intelligence specialization!



5

MI Autumn 2019

5

Jeopardy!: Watson

In 2011, Watson beat Brad Rutter, the biggest all-time money winner on Jeopardy!, and Ken Jennings, the record holder for the longest championship streak (75 days)



- Application of natural language processing, information retrieval, knowledge representation and reasoning, and machine learning
- Made up of a cluster of 90 IBM Power 750 servers with a total of 2880 POWER7 processor cores and 16 Terabytes of RAM. Each Power 750 server uses a 3.5 GHz POWER7 eight core processor, with four threads per core.

Taken from Wikipedia, August 25th, 2011

https://en.wikipedia.org/w/index.php?title=IBM_Watson&oldid=8750000

AI Autumn 2019 11

Computer games

AI in computer games (e.g. Starcraft)



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 12



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 13



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 14



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 15



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 16



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 17



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 18



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 19



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 20



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 21



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 22



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 23



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 24



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 25



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 26



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 27



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 28



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 29



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 30



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 31



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 32



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 33



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 34



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 35



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 36



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 37



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 38



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 39



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 40



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 41



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 42



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 43



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 44



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 45



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 46



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 47



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 48



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 49



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 50



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 51



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 52



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 53



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 54



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 55



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 56



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 57



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 58



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 59



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 60



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 61



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 62



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 63



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 64



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 65



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 66



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 67



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 68



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 69



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 70



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 71



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 72



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 73



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 74



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 75



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 76



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 77



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 78



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019 79



<http://ais.sjtu.edu/StarcraftAICompetition>

AI Autumn 2019

Automatic Translation

Autonomous driving

Achievements: Automatic Translation

Google translate:



MII Autumn 2019 | AI | 15

Autonomous driving

Achievement: DARPA grand challenge 2005



June 2016: approx. 2.8 mill. km. autonomous driving.

MII Autumn 2019 | AI | 16

Image analysis/processing

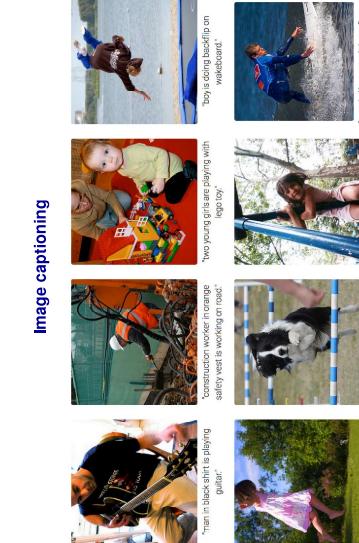
Colorization of images



Zhang, Isla. Eros: Colorful Image Colorization. In: ECOV 2016.

MII Autumn 2019 | AI | 17

Image analysis/processing



Andrej Karpathy, Li FeiFei. Deep Visual-Semantic Alignments for Generating Image Descriptions. CVPR 2015

MII Autumn 2019 | AI | 17

Other application areas

- Medical diagnosis and advisory systems
- Information processing and filtering,
- Display of information for time-critical decisions
- Spam filtering
- Optical character recognition
- Profiling/Credit scoring : profiling customers
- Bioinformatics
- Real estate: Prediction of house prices
- Computer networks: intrusion detection
- Alert and monitoring systems
- Speech recognition
- Face recognition, image annotation
- Action recognition (in video sequences)
- ...

18

Some queries to try:

Turing Test (Loebner Competition)

- Loebner Competition: (Non-scientific) competition for computer systems performing under Turing Test conditions.
- <http://www.pandorabots.com/pandora/talk?bot.id=f5d922d97e34aa1>
- **Online help**
- Combine natural language interface with expert knowledge.
- Restricted Domain (Geography): CHAT-80 (1982)
- Broad Domain ('all factual knowledge'): Wolfram Alpha
<http://www.wolframalpha.com/>

MI Autumn 2019

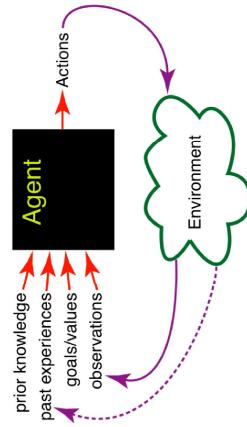
19

AI

AI and Agents

AI is the field that studies the synthesis and analysis of computational agents that act intelligently. [PM, p.3]

A coupling of perception, reasoning, and acting comprises an agent. [PM, p.10]



AI

20

AI and Agents

AI is the field that studies the synthesis and analysis of computational agents that act intelligently. [PM, p.3]

A coupling of perception, reasoning, and acting comprises an agent. [PM, p.10]



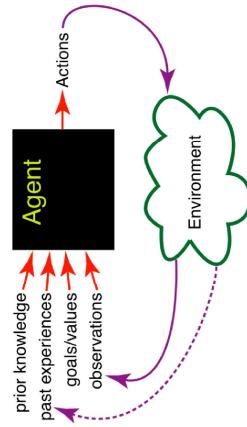
AI

21

AI and Agents

AI is the field that studies the synthesis and analysis of computational agents that act intelligently. [PM, p.3]

A coupling of perception, reasoning, and acting comprises an agent. [PM, p.10]



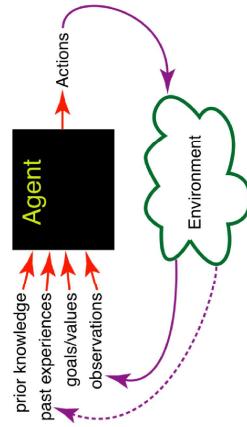
The Agent view of AI

21

AI and Agents

AI is the field that studies the synthesis and analysis of computational agents that act intelligently. [PM, p.3]

A coupling of perception, reasoning, and acting comprises an agent. [PM, p.10]



The Agent view of AI

21

- Some special flavors:
- Autonomous agents
 - Intelligent agents
 - Software agents
 - Multi-agent systems

- Some special flavors:
- perception ~ input
 - reasoning ~ computation
 - acting ~ output
 - "agent" a design metaphor, not a strict technical concept

- What is *not* an agent?
- perception ~ input
 - reasoning ~ computation
 - acting ~ output
 - "agent" a design metaphor, not a strict technical concept

MI Autumn 2019

The Agent view of AI

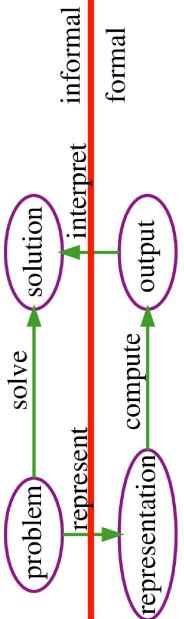
21

The Agent view of AI

21

Represent and Compute

Represent and Compute



- A representation should be:
- Sufficiently rich to encode the required knowledge.
 - Be "close" to the problem.
 - Amenable to efficient computation.
 - Able to be acquired from people, data, or experience.

MI Autumn 2019

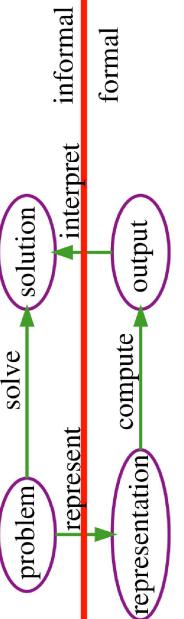
22

The Agent's View of AI

22

Represent and Compute

Represent and Compute



- A representation should be:
- Sufficiently rich to encode the required knowledge.
 - Be "close" to the problem.
 - Amenable to efficient computation.
 - Able to be acquired from people, data, or experience.
- Questions to be considered:
- What is a solution and how good should it be (optimal, satisfying, approximately, probable)?
 - How can the problem be represented?
 - How can an output be computed (what properties should a solution have)?

MI Autumn 2019

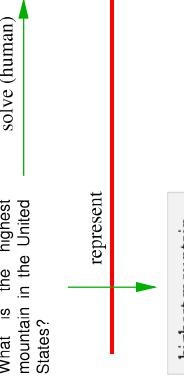
22

The Agent's View of AI

22

Represent and Compute

Represent and Compute

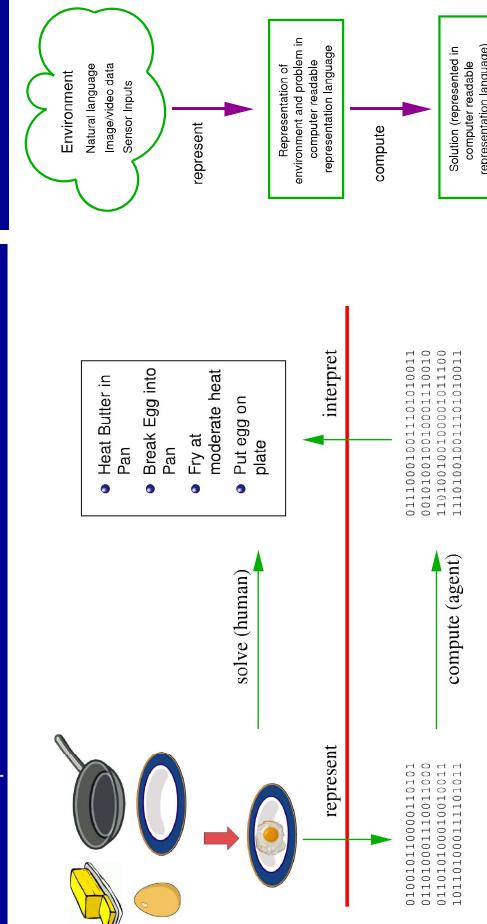


	highest mountain in	United States	solve (agent)	solve (human)
1	Mount McKinley		20 322 ft	
2	Mount Rainier		18 212 ft	
3	Mount Saint Elias		18 009 ft	
4	Mount Foraker		17 402 ft	
5	Bona		16 421 ft	



Represent and Compute

Problem Formalization



MI Autumn 2019

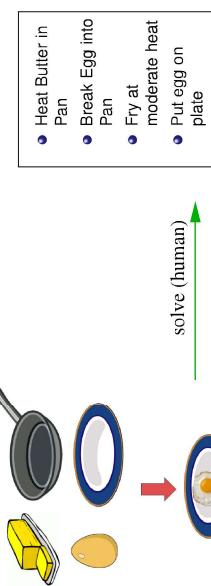
24

The Agent's View of AI

24

Problem Formalization

Problem Formalization



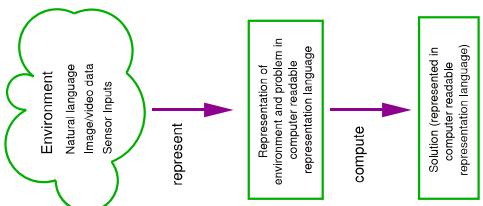
010010110000110101
011010001011001100
001010001000001100
01101001001000010110
110100101100101011

25

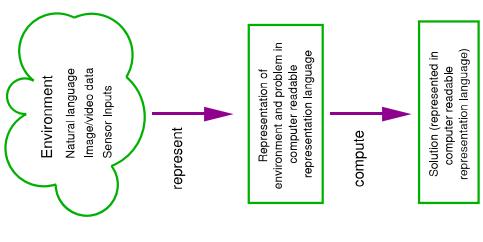
The Agent's View of AI

25

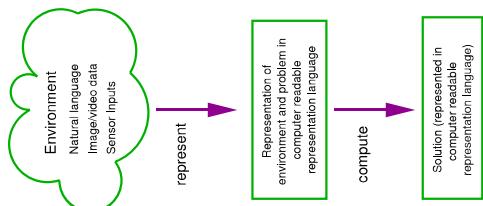
Problem Formalization



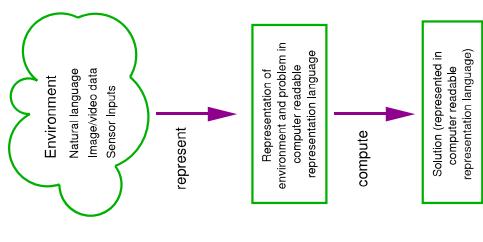
Problem Formalization



Problem Formalization



Problem Formalization

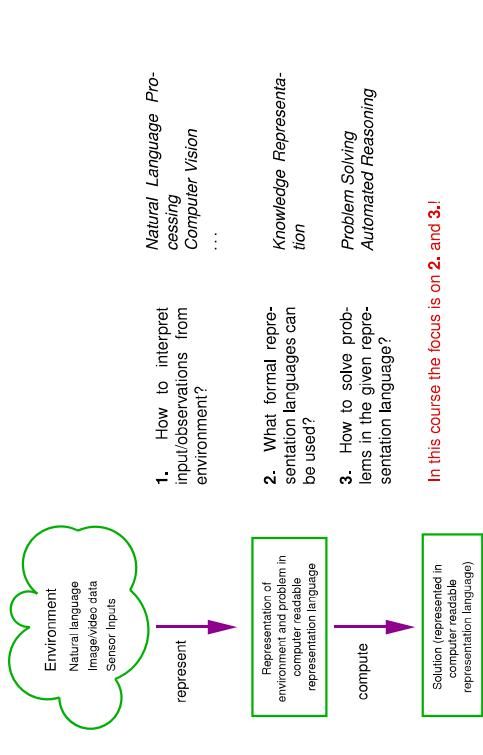


Levels of Representation

10

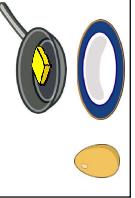
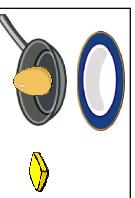
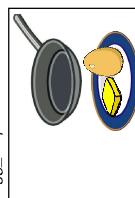
We consider 3 representation schemes

- State based
 - Feature based



State based

Feature based

					
State 01	State 04	State 08	State 12	State 14	State 18
<small>MU Autumn 2019</small>	<small>MU Autumn 2019</small>	<small>MU Autumn 2019</small>	<small>MU Autumn 2019</small>	<small>MU Autumn 2019</small>	<small>MU Autumn 2019</small>
<small>Representing the problem</small>	<small>Representing the problem</small>	<small>Representing the problem</small>	<small>Representing the problem</small>	<small>Representing the problem</small>	<small>Representing the problem</small>

30 binary features represent $2^{30} = 1.073.741.824$ states.

MU Autumn 2019 Representing the problem 28

Relational

Levels of Detail

	
<small>state(egg,whole), in(butter,pan), in(egg,table)</small>	<small>state(egg,broken), in(butter,pan), in(egg,plate)</small>
<small>state(egg,whole), in(butter,pan), in(egg,table)</small>	<small>state(egg,broken), in(butter,pan), in(egg,plate)</small>

1 binary relation and 100 individuals give $100^2 = 10.000$ boolean features, or $2^{10.000}$ states.

MU Autumn 2019 Representing the problem 29

Other dimensions of complexity

Modularity

- Flat, modular, hierarchical

Modularity

- Flat, modular, hierarchical

Planning horizon

- Non-planning, finite, indefinite, infinite

Planning horizon

- Non-planning, finite, indefinite, infinite

MU Autumn 2019 Representing the problem 31

MU Autumn 2019 Representing the problem 31

31

Other dimensions of complexity

Modularity

- Flat, modular, hierarchical

Planning horizon

- Non-planning, finite, indefinite, infinite

Domain uncertainty

- Fully observable vs. partially observable world.

- Deterministic vs. stochastic.

Preferences

- Achievement goals, maintenance goals.

- Deterministic vs. stochastic.

- Complex ordinal or cardinal preferences.

MI Autumn 2019

Representing the problem

31

Other dimensions of complexity

Modularity

- Flat, modular, hierarchical

Planning horizon

- Non-planning, finite, indefinite, infinite

Domain uncertainty

- Fully observable vs. partially observable world.

- Deterministic vs. stochastic.

Preferences

- Achievement goals, maintenance goals.

- Deterministic vs. stochastic.

Number of agents

- Single agent or multiple agents

Learning

- Knowledge is given at design time or learned from experience.

MI Autumn 2019

Representing the problem

31

Other dimensions of complexity

Modularity

- Flat, modular, hierarchical

Planning horizon

- Non-planning, finite, indefinite, infinite

Domain uncertainty

- Fully observable vs. partially observable world.

- Deterministic vs. stochastic.

Preferences

- Achievement goals, maintenance goals.

- Complex ordinal or cardinal preferences.

Number of agents

- Single agent or multiple agents

Learning

- Knowledge is given at design time or learned from experience.

Computational limits

- Perfect or bounded rationality

MI Autumn 2019

Representing the problem

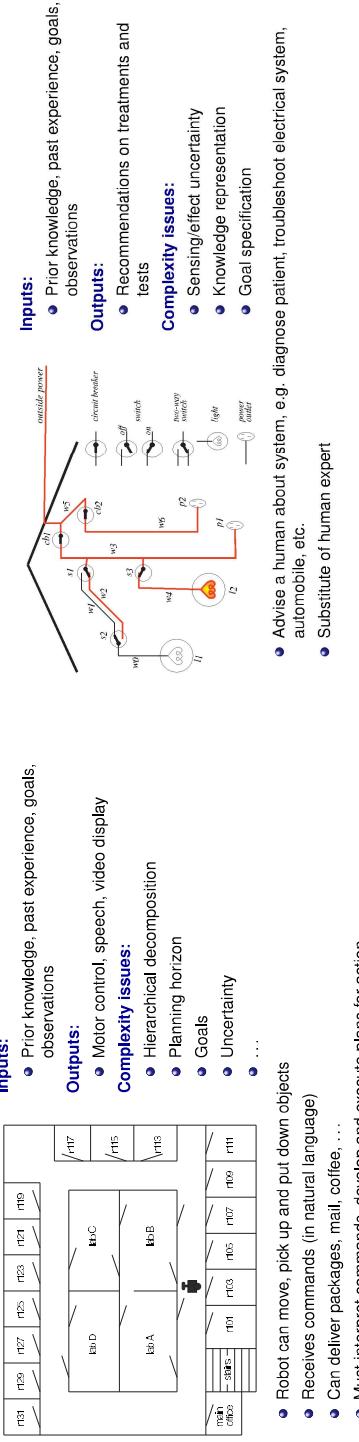
31

Prototypical Applications

MI Autumn 2019

Representing the problem

32



MI Autumn 2019 | Prototypical Applications | 33

Trading Agent

Characteristics

- Automatically buy/sell goods for user/company (possibly at auction)
- Determine good strategy to procure necessary goods in time at best price

Inputs

- Prior knowledge about goods, past experience, preferences, observations

Output

- Proposals to the user

Tentative course overview

Topics:

- Introduction
- Problem solving as search
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Machine learning
- Planning
- Reinforcement learning
- Multagent systems

MI Autumn 2019 | Prototypical Applications | 34

Tentative course overview

Topics:

- Introduction
- Search-based methods**
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Machine learning
- Planning
- Multi-agent systems

MI E19 | 1

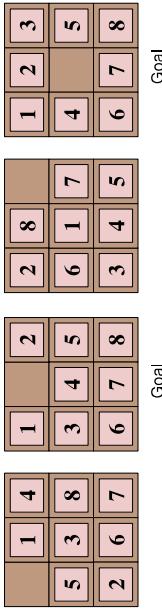
Problem Description

Problem Solving as Search

- has a state-based representation of its environment
 - can observe with certainty which state it is in
 - has a certain goal it wants to achieve
 - can execute actions that have definite effects (no uncertainty)
- The agent needs to find a sequence of actions that lead it to a **goal state**: a state in which its goal is achieved.

Example: 8 Puzzle

Problem: re-arrange tiles into goal configuration:

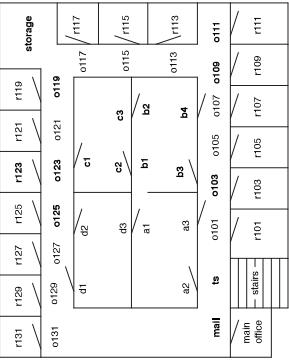


Actions: move *up*, move *down*, move *left*, move *right*

There are 362880 states

- States: locations, e.g. r131, storage, o117, c3, ...
- Actions: move to neighboring locations, e.g. move_r131_o131, move_o119_storage, move_b2_c3, ...

Example: Office Robot



Actions: move *up*, move *down*, move *left*, move *right*

There are 362880 states

- States: locations, e.g. r131, storage, o117, c3, ...
- Actions: move to neighboring locations, e.g. move_r131_o131, move_o119_storage, move_b2_c3, ...

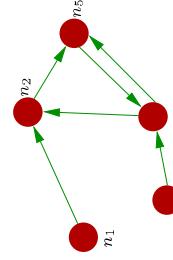
State-Space Problem

A **State-Space Problem** consists of

- A set of states
- A subset of **start states**
- A set of actions (not all actions available at all states)
- An **action function** that for a given state s and action a returns the state reached when executing a in s
- A **goal test** that for any state s returns the boolean value $goal(s)$ (true if s is a goal state)
- (optional) a **cost function** on actions
- (optional) a **value function** on goal states

A **Solution** consists of

- For any given start state, a sequence of actions that lead to a goal state
- (optional) a sequence of actions with minimal cost
- (optional) a sequence of actions leading to a goal state with maximal value

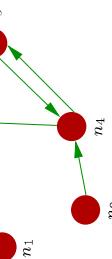


- A **directed graph** consists of
 - a set of **nodes**
 - a set of **arcs** (ordered pairs of nodes)

Graphs

State-Space Problems as Graphs

- Nodes: states
- Arcs: possible state-transitions from actions (arcs can be labeled with actions)



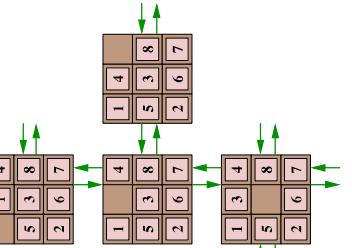
A **directed graph** consists of

- a set of **nodes**
- a set of **arcs** (ordered pairs of nodes)

Further terminology:

- n_2 is a **neighbor** of n_4 (not the other way round!).
- n_3, n_4, n_2, n_5 is a **path** from n_3 to n_5 .
- n_2, n_5, n_4, n_2 is a path that is a **cycle**.
- a graph is **acyclic** if it has no cycles.

8 puzzle graph (part)

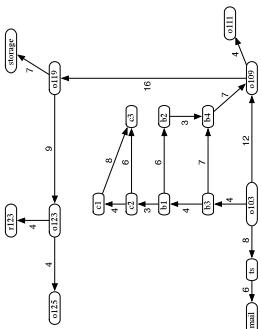


MI E 19

Problem Solving as Search 8

Delivery Robot Graph

Graph Search



Arcs labeled with costs (time to travel, fuel costs, ...)

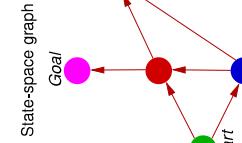
- Forward branching factor of a node = number of arcs leaving that node.
- Backward branching factor of a node = number of arcs entering that node.

MI E 19

Problem Solving as Search 9

Graph Search

From Graph to Search Tree



Search tree

State-space graph

Goal

Start

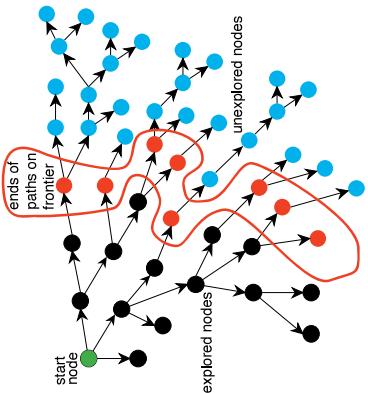
- A state-space problem can be solved by searching in the state-space graph for paths from start states to goal states.
- This does not require the whole graph at once: search may only locally generate neighbors of currently visited node.

- Tree: special graph with
 - exactly one node that has no incoming arc (the **root**)
 - all other nodes have exactly one incoming arc (may have 0, 1, 2, 3, ... outgoing arcs)
- Nodes in the search tree correspond to paths in the graph beginning in the start state
- Nodes in the search tree also are labeled with states: the last state of the path

MI E 19

Problem Solving as Search 10

Problem Solving as Search 11



Input: a graph,
a set of start nodes,
Boolean procedure $goal(n)$ that tests if n is a goal node.
 $frontier := \{(s) : s \text{ is a start node}\}$;
while $frontier$ is not empty,
 select and remove path $\langle n_0, \dots, n_k \rangle$ from $frontier$,
 if $goal(n_k)$
 return $\langle n_0, \dots, n_k \rangle$;
 for every neighbor n of n_k
 add $\langle n_0, \dots, n_k, n \rangle$ to $frontier$,
end while

MIE19 Problem Solving as Search 12

Generic Search Algorithm

Depth-first search

- Select from the frontier that path that was most recently added to the frontier (frontier implemented as **stack**).

Example

• Explored nodes with order of exploration
• Frontier (colored)
• Unexplored nodes



MIE19 Problem Solving as Search 13

Input: a graph,
a set of start nodes,
Boolean procedure $goal(n)$ that tests if n is a goal node.
 $frontier := \{(s) : s \text{ is a start node}\}$;
while $frontier$ is not empty,
 select and remove path $\langle n_0, \dots, n_k \rangle$ from $frontier$,
 if $goal(n_k)$
 return $\langle n_0, \dots, n_k \rangle$;
 for every neighbor n of n_k
 add $\langle n_0, \dots, n_k, n \rangle$ to $frontier$,
end while

The algorithm does not require the complete graph as input: only needed are:

- List of start nodes (often only one)
- Boolean function $goal(\text{Node } n)$
- Function $\text{get_neighbors}(\text{Node } n)$ returning list of neighbors of n .

MIE19 Problem Solving as Search 14

Depth-first search

- Select from the frontier that path that was earliest added to the frontier (frontier implemented as **queue**).

Example

• Explored nodes with order of exploration
• Frontier (colored)
• Unexplored nodes



MIE19 Problem Solving as Search 15

Properties

- Space used is linear in the length of the current path.
- May not terminate if state-space graph has cycles
- With a forward branching factor bounded by b and depth n , the worst-case time complexity of a finite tree is b^n .

MIE19 Problem Solving as Search 16

- Will always find a solution if one exists
- Size of frontier always increases during search up to order of magnitude of total size of search tree.

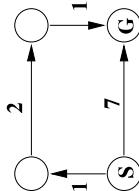
- Will always find a solution if one exists
- Size of frontier always increases during search up to order of magnitude of total size of search tree.
- Can be adapted to find a minimum cost path.

Problems with cost function

Problem

- Assume that for each action at each state we have an associated **cost**
- The cost of a solution is the sum of the costs of all actions on the path from start to goal state.
- A **minimum cost solution** is a solution with minimal cost.

Example



Breadth first search will find the *shortest*, but not the *cheapest* solution.

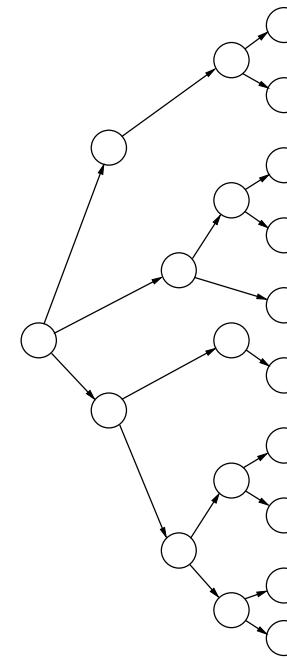
Depth first search may find either solution, depending on order of neighbor enumeration

Lowest-Cost-First Search

Properties

- If all actions have non-zero cost, and solution exists, then a minimal cost solution will be found.
- with each path in *frontier* store the cost of the path
- Space requirement depends on cost structure, but usually similar to breadth-first search.

Iterative Deepening: an example



Goal

- Termination guarantee of breadth-first search
- Space efficiency of depth-first search

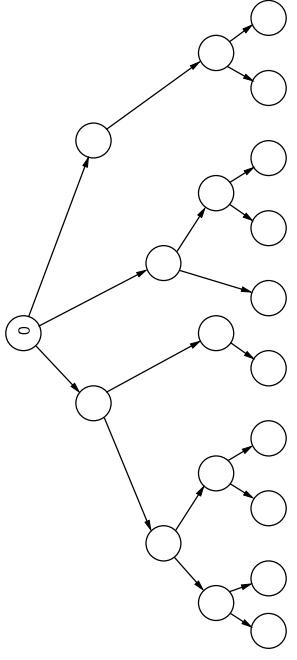
Algorithm

Depth-bounded search k

- As depth-first search, but
 - do not add neighbors of selected node to frontier if selected node has depth k .
- For $k = 1, 2, 3, \dots$ perform depth-bounded search k .

Iterative deepening: an example

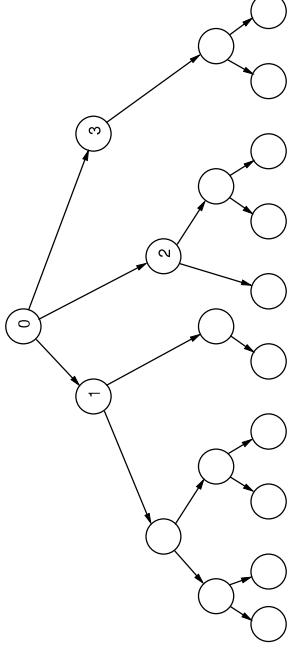
Iterative deepening: an example



Perform depth-bounded search to level $k = 1$.

NB: The node numbering corresponds to when the states are first visited.

MI E19 Problem Solving as Search 21



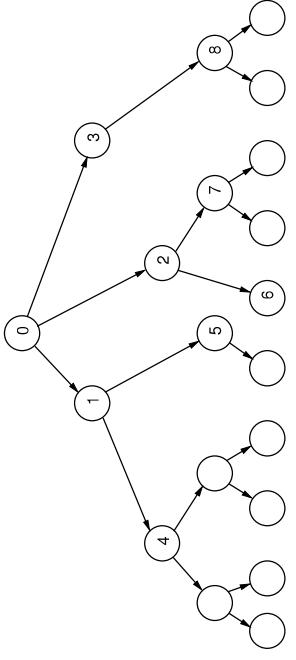
Perform depth-bounded search to level $k = 2$.

NB: The node numbering corresponds to when the states are first visited.

MI E19 Problem Solving as Search 21

Iterative deepening: an example

Iterative deepening: an example



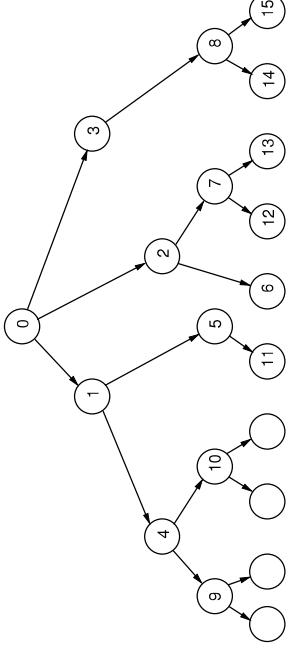
Perform depth-bounded search to level $k = 3$.

NB: The node numbering corresponds to when the states are first visited.

MI E19 Problem Solving as Search 21

Iterative deepening: an example

Iterative deepening: an example



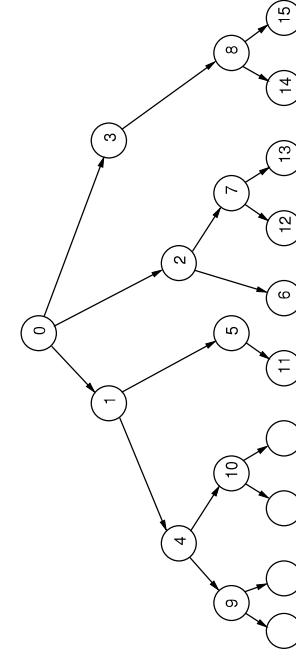
Perform depth-bounded search to level $k = 2$.

NB: The node numbering corresponds to when the states are first visited.

MI E19 Problem Solving as Search 21

Iterative deepening: an example

Uninformed Search



Depth-first, Breadth-first and iterative deepening are **uninformed** search strategies: they do not assume/use any knowledge of the search space except the pure graph structure.

Properties

- Has desired termination and space efficiency properties
- Duplicates computations (depth-bounded search k repeats computations of depth-bounded search $k - 1$). Not as problematic as it looks: constant overhead of $(b/(b - 1))$.

MI E19 Problem Solving as Search 21

MI E19 Problem Solving as Search 22

Heuristic Search

Idea

- Lowest-Cost First Search only considers cost of already constructed partial solution.
- Idea: Try to estimate the cost of optimal path from current state to goal.

Informed Search

Informed Search

MI E 19 23 Informed Search

Heuristic Search

Idea

- Lowest-Cost-First Search only considers cost of already constructed partial solution.
- Idea: Try to estimate the cost of optimal path from current state to goal.

Actual Cost

Given a cost function on actions, can define for any node n in the search tree:
 $\text{opt}(n) := \text{cost of optimal (minimal cost) path from } n \text{ to a goal state (infinite if no path to goal exists).}$

- The opt function can usually not be computed
- $\text{opt}(n)$ only depends on the state at node n .

Heuristic Function

$h(n)$ (non-negative number): estimate of $\text{opt}(n)$, $h(n)$ is an **underestimate** if for all nodes n :
$$h(n) \leq \text{opt}(n)$$

Informed Search

MI E 19 23 Informed Search

Heuristic Search

Idea

- Lowest-Cost-First Search only considers cost of already constructed partial solution.
- Idea: Try to estimate the cost of optimal path from current state to goal.

Actual Cost

Given a cost function on actions, can define for any node n in the search tree:
 $\text{opt}(n) := \text{cost of optimal (minimal cost) path from } n \text{ to a goal state (infinite if no path to goal exists).}$

- The opt function can usually not be computed
- $\text{opt}(n)$ only depends on the state at node n .

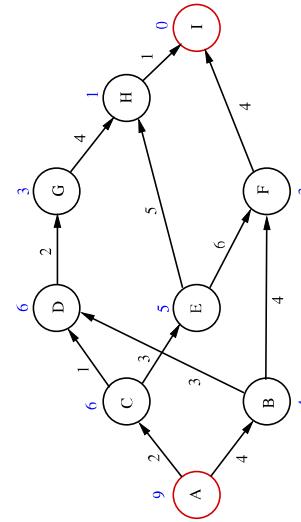
Heuristic Function

$h(n)$ (non-negative number): estimate of $\text{opt}(n)$, $h(n)$ is an **underestimate** if for all nodes n :
$$h(n) \leq \text{opt}(n)$$

Informed Search

MI E 19 23 Informed Search

A^* search: an example



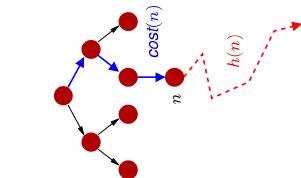
A^* search: always expand fringe node with minimal f -value.

For node D:

- $\text{cost}(D) = 3$
- $h(D) = 6$
- $f(D) := \text{cost}(D) + h(D) = 3 + 6 = 9$

Informed Search

MI E 19 25 Informed Search

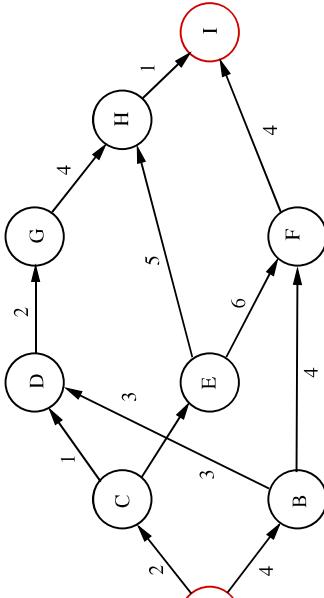


For any node n in the search tree we have:

- $\text{cost}(n)$: (true) cost of reaching n from the root node.
- $h(n)$: a heuristic function
- $f(n) := \text{cost}(n) + h(n)$

Dynamic programming: Example

Dynamic programming: Example



Idea: for statically stored graphs, build a table of $dist(n)$ the actual distance of the shortest path from node n to a goal.
This can be built backwards from the goal.

$$dist(n) = \begin{cases} 0 & \text{if } is_goal(m) \\ \min_{(n,m)} \{ |(n,m)| + dist(m) \} & \text{otherwise.} \end{cases}$$

This can be used locally to determine what to do.

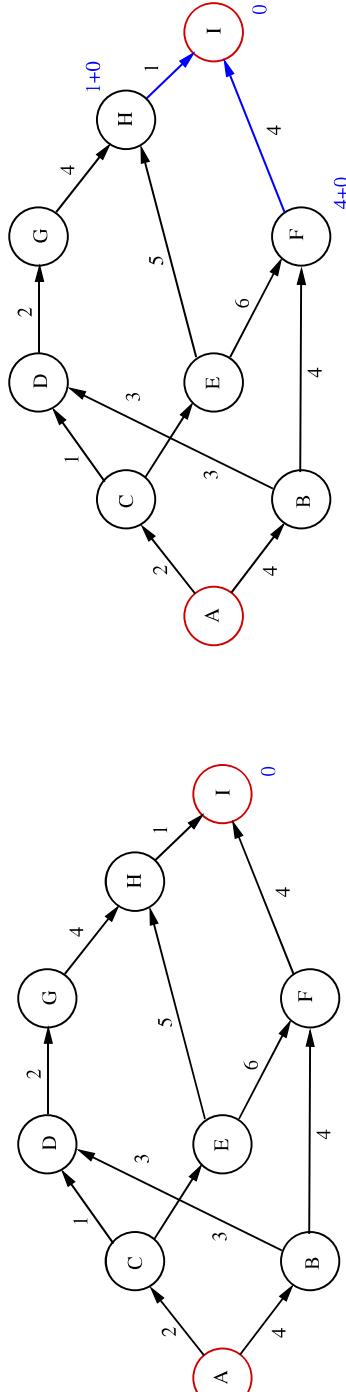
MI E 19

Informed Search

MI E 19 Informed Search 30

Dynamic programming: Example

Dynamic programming: Example



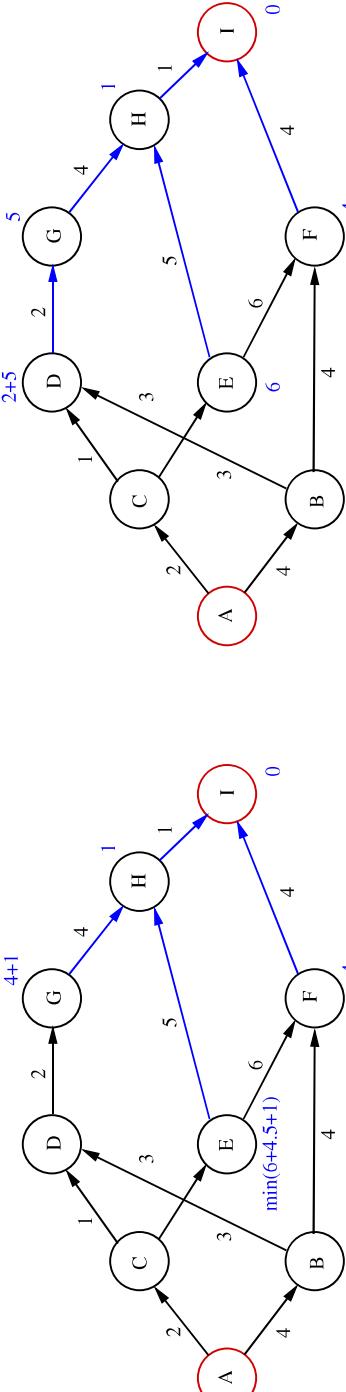
MI E 19

Informed Search

MI E 19 Informed Search 30

Dynamic programming: Example

Dynamic programming: Example



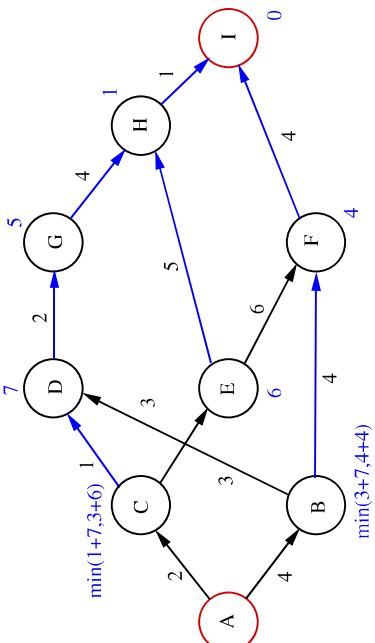
MI E 19

Informed Search

MI E 19 Informed Search 30

Dynamic programming: Example

Dynamic programming: Example



MI E19

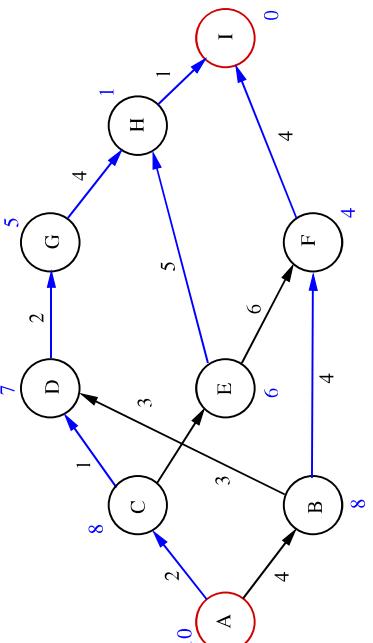
Informed Search

30

MI E19 Informed Search 30

Dynamic programming: Example

Direction of search



There are two main problems:

- You need enough space to store the graph.
- The *distf* function needs to be recomputed for each goal.

MI E19

Informed Search

30

MI E19 Informed Search 30

Direction of search

Cycle checking

- The definition of searching is symmetric: find path from start nodes to goal node or from goal node to start nodes.
- Search complexity is b^n . Should use forward search if forward branching factor is less than backward branching factor, and vice versa.
- Note: sometimes when graph is dynamically constructed, you may not be able to construct the backwards graph.

Bi-directional search

- You can search backward from the goal and forward from the start simultaneously.
- This wins as $2b^{k/2} \ll b^k$. This can result in an exponential saving in time and space.
- The main problem is making sure the frontiers meet.
- This is often used with one breadth-first method that builds a set of locations that can lead to the goal. In the other direction another method can be used to find a path to these interesting locations.
- A searcher can prune a path that ends in a node already on the path, without removing an optimal solution.
- Using depth-first methods, with the graph explicitly stored, this can be done in constant time.
- For other methods, the cost is linear in path length.

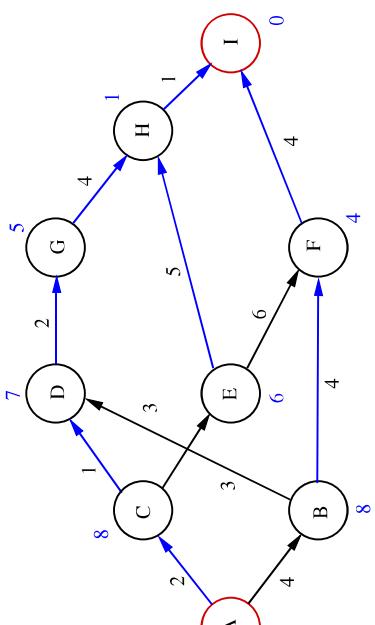
MI E19

Informed Search

31

MI E19 Informed Search 31

Dynamic programming: Example



MI E19

Informed Search

30

MI E19 Informed Search 30

Direction of search

- The definition of searching is symmetric: find path from start nodes to goal node or from goal node to start nodes.
- Search complexity is b^n . Should use forward search if forward branching factor is less than backward branching factor, and vice versa.
- Note: sometimes when graph is dynamically constructed, you may not be able to construct the backwards graph.

MI E19 Informed Search 30

Machine Intelligence

Lecture 3: Constraint satisfaction problems

Thomas Dyrre Nielsen

Aalborg University

Topics:

- Introduction
- Search-based methods
- **Constraint satisfaction problems**
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Machine learning
- Planning
- Multi-agent systems

Autumn 2019

1

Features and Variables

Describing the world (environment) by features:

Name (Algebraic Variable)	Symbol	Domain
Symbol_on_square_1	$\{1, 2, 3, 4, 5, 6, 7, 8, \text{empty}\}$	
Robot_battery	$\{\text{full}, \text{half}, \text{empty}\}$	
Robot_position	$\{r1, r2, \dots, r11\}$	
Coffee_ready	$\{\text{true}, \text{false}\}$	
No_of_undelivered_packages	$\{1, 2, 3, \dots\}$	
Temperature	$[-25, 40]$	

- We will be mostly concerned with (algebraic) variables that have a *finite* domain.
- Special interest: **boolean** variable with domain $\{\text{true}, \text{false}\}$
- Numerical variables can be approximated:

$$\begin{array}{ccc} \{1, 2, 3, \dots\} & \mapsto & \{1, 2, 3, 4, 5, > 5\} \\ [-25, 40] & \mapsto & \{-25, -24, \dots, -1, 0, 1, \dots, 40\} \end{array}$$

Autumn 2019
2

2

Example: Cooking

Variables:

egg	$\{\text{whole}, \text{broken}\}$
butter_in	$\{\text{pan}, \text{plate}, \text{table}\}$
egg_in	$\{\text{pan}, \text{plate}, \text{table}\}$

From variables to possible worlds

A **possible world** for a set of variables is an assignment of a value to each variable.

Connection with state spaces

The set of all possible worlds for a given set of variables defines a state space (we can also call a possible world simply a state).



egg=whole
butter_in=pan
egg_in=table

Autumn 2019
3

3

=Features and Possible Worlds

Autumn 2019
4

4

=Features and Possible Worlds

Autumn 2019
5

5

=Features and Possible Worlds

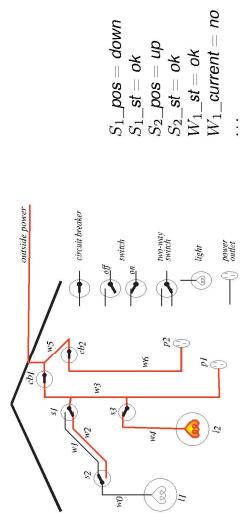
Example: Electrical

Example: Schedule

Variables:

$S_1.\text{pos}$: up, down}	$S_1.\text{st}$: ok, broken, short)
$S_2.\text{pos}$: up, down}	$S_2.\text{st}$	{ok, broken, short}
$W_1.\text{broken}$	{yes, no}		
...			

One out of many possible worlds:



Autumn 2019

Features and Possible Worlds

6

Variables:

$T_{\text{teacher}}.M_l$	$\{P.D, M.I, T.D.N\}$
$T_{\text{time}}.M_l$	$\{M.o.m, M.o.a, \dots, F.r.m, F.r.a\}$
$T_{\text{room}}.M_l$	$\{0.2, 12, 0.2, 13, 0.2, 90\}$
$T_{\text{teacher}}.AD$	$\{P.D, M.I, T.D.N\}$
$T_{\text{time}}.AD$	$\{M.o.m, M.o.a, \dots, F.r.m, F.r.a\}$
$T_{\text{room}}.AD$	$\{0.2, 12, 0.2, 13, 0.2, 90\}$

One out of many possible worlds:

Mo	Tue	Wed	Thu	Fri
	M.I, TD.N, 0.2,13			
		AD, PD, 0.2,90		

One possible world:

Mo	Tue	Wed	Thu	Fri
	M.I, TD.N, 0.2,13			
		AD, PD, 0.2,90		

Features and Possible Worlds

7

Constraints

A constraint is a condition on the values of variables in a possible world.

Extensional Constraint Specification

Constraint Satisfaction Problems

Explicitly list all allowed (or disallowed) combination of values:

$T_{\text{teacher}}.M_l$	$T_{\text{time}}.M_l$	$T_{\text{room}}.M_l$	$T_{\text{teacher}}.AD$	$T_{\text{time}}.AD$	$T_{\text{room}}.AD$
PD	$M.o.m$	0.2,12	PD		0.2,12
PD	$M.o.m$	0.2,12	PD		0.2,13
...

Not on the list of allowed possible worlds:

$T_{\text{teacher}}.M_l$	$T_{\text{time}}.M_l$	$T_{\text{room}}.M_l$	$T_{\text{teacher}}.AD$	$T_{\text{time}}.AD$	$T_{\text{room}}.AD$
PD	$M.o.m$	0.2,12	PD		0.2,12
PD	$M.o.m$	0.2,12	MJ		0.2,12
...

Features and Possible Worlds

8

Constraint Satisfaction Problems

Autumn 2019

Constraints

Example: Sudoku

1	2	3	4	5	6	7	8	9
4	5	6	7	8	9	1	2	3
7	8	9	1	2	3	4	5	6
2	3	4	5	6	7	8	9	1
5	6	7	8	9	1	2	3	4
8	9	1	2	3	4	5	6	7
4	5	6	7	8	9	1	2	3
1	2	3	4	5	6	7	8	9
6	7	8	9	1	2	3	4	5

Constraints:

$A_1 = 2 \vee A_2 = 2 \vee A_4 = 2 \vee A_5 = 2 \vee A_6 = 2 \vee A_7 = 2 \vee A_8 = 2 \vee A_9 = 2$
$A_1 = 3 \vee A_2 = 3 \vee A_4 = 3 \vee A_5 = 3 \vee A_6 = 3 \vee A_7 = 3 \vee A_8 = 3 \vee A_9 = 3$
$A_1 = 3 \vee A_2 = 3 \vee B_1 = 3 \vee C_1 = 3 \vee C_2 = 3 \vee C_3 = 3$
...

Autumn 2019

9

Constraint Satisfaction Problems

Autumn 2019

1	2	3	4	5	6	7	8	9
4	5	6	7	8	9	1	2	3
7	8	9	1	2	3	4	5	6
2	3	4	5	6	7	8	9	1
5	6	7	8	9	1	2	3	4
8	9	1	2	3	4	5	6	7
4	5	6	7	8	9	1	2	3
1	2	3	4	5	6	7	8	9
6	7	8	9	1	2	3	4	5

Constraint Satisfaction Problems

Autumn 2019

10

Constraint Satisfaction Problems

Autumn 2019

- A CSP can be represented as a state space problem:
- States are all partial assignments of values to variables that are consistent with the constraints
 - For a state s : select some variable V not assigned a value in s , and let the neighbors of s be all states that assign a value to V (if any exist).
 - The start state is the state that does not assign any values
 - A goal state is a state that assigns values to all variables

Other tasks:

- Determine the number of models of the constraints
- Find an *optimal* model (given also a value function on possible worlds).
- ...

Autumn 2019

Constraint Satisfaction Problems

11

Constraint Satisfaction Problems

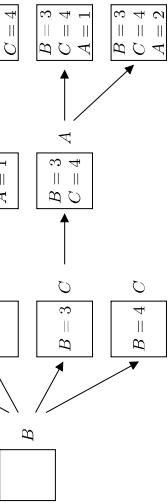
12

CSP as State Space Problem

Example [PM 4.13]

A CSP can be represented as a state space problem:

- States are all partial assignments of values to variables that are consistent with the constraints
 - For a state s : select some variable V not assigned a value in s , and let the neighbors of s be all states that assign a value to V (if any exist).
 - The start state is the state that does not assign any values
 - A goal state is a state that assigns values to all variables
- Solving the CSP**
- A solution to the state space problem is a path with a goal state at the end; a solution to the CSP problem
 - To solve the state space problem need only be able to:
 - enumerate all partial assignments that assign a value to one more variable than s
 - check whether a partial assignment is consistent with the constraints (that is sufficient to implement the `get_neighbors` and `goal` functions needed in the generic search algorithm)



- The state space graph is a tree (= search tree)

Autumn 2019

Constraint Satisfaction Problems

12

Autumn 2019

13

Consistency Algorithms

Example: Variables A, B, C ; all with domain $\{1, 2, 3, 4\}$.

Constraints: $A < B, B < C$.

Observation: There is no solution with $A = 4$.

Approach to solving CSPs: iteratively eliminate value assignments that cannot be part of a solution.

Idea

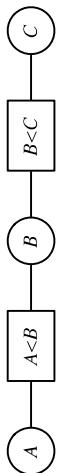
Autumn 2019

Consistency Algorithms

14

Autumn 2019

14



The **constraint network** for a CSP consists of

- One (oval) node for each variable X
- One (rectangular) node for each constraint c
- An (undirected) arc $\langle X, c \rangle$ between every constraint and every variable involved in the constraint

With each variable node X is associated a (reduced) **domain** D_X :

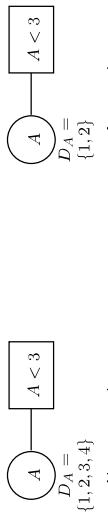
- Initially the domain of the variable
- Reduced by successively deleting values that cannot be part of a solution

An arc $\langle X, c \rangle$ is **arc consistent**, if

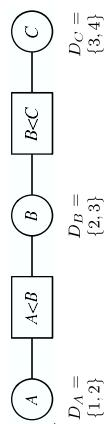
- for all $x \in D_X$ there exists values y_1, \dots, y_k for the other variables involved in c , such that x, y_1, \dots, y_k is consistent with c .

A constraint network is **arc consistent**, if all its arcs are arc consistent.

Examples

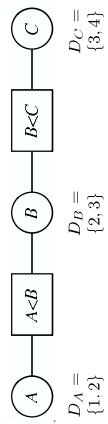


Arc Consistency Examples



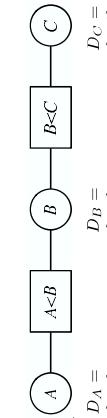
Arc consistent. Not every combination of values from D_A, D_B, D_C is a solution!

Arc Consistency Examples



Arc consistent. Not every combination of values from D_A, D_B, D_C is a solution!

Arc Consistency Examples



Arc consistent. Not every combination of values from D_A, D_B, D_C is a solution!



Arc consistent. There exists no solution!

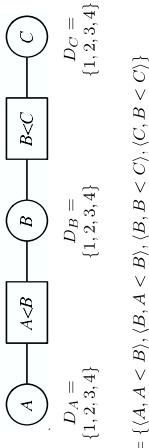
Generalized Arc Consistency Algorithm

Generalized Arc Consistency Algorithm

Algorithm Outline

1. **To-do-arcs**= all arcs in constraint network // Potentially inconsistent arcs
2. **while** **To-do-arcs** $\neq \emptyset$
3. select and delete one arc $\langle X, c \rangle$ from **To-do-arcs**
4. make arc consistent by deleting values from D_X , if necessary
5. if values were deleted: add all other arcs $\langle Z, c' \rangle$ ($c \neq c'$, $X \in dom(c')$) to **To-do-arcs**

Example



- **To-do-arcs** = $\{\langle A, A < B \rangle, \langle B, A < B \rangle, \langle B, B < C \rangle, \langle C, B < C \rangle\}$
- Selecting $\langle A, A < B \rangle$: For $A = 4$, no value of B satisfies $4 < B$.
- Remove $\langle A, A < B \rangle$ from **To-do-arcs**.
- Update D_A .

Autumn 2019 18 Consistency Algorithms

Algorithm Outline

1. **To-do-arcs**= all arcs in constraint network // Potentially inconsistent arcs
2. **while** **To-do-arcs** $\neq \emptyset$
3. select and delete one arc $\langle X, c \rangle$ from **To-do-arcs**
4. make arc consistent by deleting values from D_X , if necessary
5. if values were deleted: add all other arcs $\langle Z, c' \rangle$ ($c \neq c'$, $X \in dom(c')$) to **To-do-arcs**

Example



- **To-do-arcs** = $\{\langle B, A < B \rangle, \langle B, B < C \rangle, \langle C, B < C \rangle\}$
- Selecting $\langle B, A < B \rangle$: $B = 1$ can be pruned.
- Remove $\langle B, A < B \rangle$ from **To-do-arcs**.
- Update D_B .

Autumn 2019 18 Consistency Algorithms

Algorithm Outline

1. **To-do-arcs**= all arcs in constraint network // Potentially inconsistent arcs
2. **while** **To-do-arcs** $\neq \emptyset$
3. select and delete one arc $\langle X, c \rangle$ from **To-do-arcs**
4. make arc consistent by deleting values from D_X , if necessary
5. if values were deleted: add all other arcs $\langle Z, c' \rangle$ ($c \neq c'$, $X \in dom(c')$) to **To-do-arcs**

Example



- **To-do-arcs** = $\{\langle B, B < C \rangle, \langle C, B < C \rangle\}$
- Selecting $\langle B, B < C \rangle$: $B = 4$ can be pruned.
- Add $\langle A, A < B \rangle$ to **To-do-arcs** and remove $\langle B, B < C \rangle$.
- Update D_B .

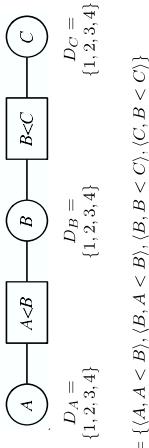
Autumn 2019 18 Consistency Algorithms

Generalized Arc Consistency Algorithm

Algorithm Outline

1. **To-do-arcs**= all arcs in constraint network // Potentially inconsistent arcs
2. **while** **To-do-arcs** $\neq \emptyset$
3. select and delete one arc $\langle X, c \rangle$ from **To-do-arcs**
4. make arc consistent by deleting values from D_X , if necessary
5. if values were deleted: add all other arcs $\langle Z, c' \rangle$ ($c \neq c'$, $X \in dom(c')$) to **To-do-arcs**

Example



- **To-do-arcs** = $\{\langle A, A < B \rangle, \langle B, A < B \rangle, \langle B, B < C \rangle, \langle C, B < C \rangle\}$
- Selecting $\langle A, A < B \rangle$: For $A = 4$, no value of B satisfies $4 < B$.
- Remove $\langle A, A < B \rangle$ from **To-do-arcs**.
- Update D_A .

Autumn 2019 18 Consistency Algorithms

Algorithm Outline

1. **To-do-arcs**= all arcs in constraint network // Potentially inconsistent arcs
2. **while** **To-do-arcs** $\neq \emptyset$
3. select and delete one arc $\langle X, c \rangle$ from **To-do-arcs**
4. make arc consistent by deleting values from D_X , if necessary
5. if values were deleted: add all other arcs $\langle Z, c' \rangle$ ($c \neq c'$, $X \in dom(c')$) to **To-do-arcs**

Example



- **To-do-arcs** = $\{\langle B, A < B \rangle, \langle B, B < C \rangle, \langle C, B < C \rangle\}$
- Selecting $\langle B, A < B \rangle$: $B = 1$ can be pruned.
- Remove $\langle B, A < B \rangle$ from **To-do-arcs**.
- Update D_B .

Autumn 2019 18 Consistency Algorithms

Algorithm Outline

1. **To-do-arcs**= all arcs in constraint network // Potentially inconsistent arcs
2. **while** **To-do-arcs** $\neq \emptyset$
3. select and delete one arc $\langle X, c \rangle$ from **To-do-arcs**
4. make arc consistent by deleting values from D_X , if necessary
5. if values were deleted: add all other arcs $\langle Z, c' \rangle$ ($c \neq c'$, $X \in dom(c')$) to **To-do-arcs**

Example



- **To-do-arcs** = $\{\langle B, B < C \rangle, \langle C, B < C \rangle\}$
- Selecting $\langle B, B < C \rangle$: $B = 4$ can be pruned.
- Add $\langle A, A < B \rangle$ to **To-do-arcs** and remove $\langle B, B < C \rangle$.
- Update D_B .

Autumn 2019 18 Consistency Algorithms

Variable Elimination

- Arc Consistency: simplify problem by eliminating values
- Variable Elimination: simplify problem by eliminating variables

Autumn 2019

Variable Elimination

19

19

Relational Operations

Projection

Variable Elimination operates on extensional (table) representations of constraints:

$B < C$:

A		B		C	
A	B	1	2	1	3
1	2	1	3	1	4
1	3	1	4	2	3
1	4	1	4	2	4
2	3	2	3	3	4
2	4	2	4	3	4
3	4	3	4		

Projection of a table:

Course	Year	Student	Grade
cs322	2008	fran	77
cs111	2009	billie	88
cs111	2009	jess	78
cs444	2008	billie	88
cs322	2009	fran	92
		jordan	92

Algorithm requires **projection** and **join** operations on tables.

Variable Elimination

19

19

$A < B$:

A		B		C	
A	B	1	2	1	3
1	2	1	3	1	4
1	3	1	4	2	3
1	4	1	4	2	4
2	3	2	3	3	4
2	4	2	4	3	4
3	4	3	4		

Projection of a table:

Course	Year	Student	Grade
cs322	2008	fran	77
cs111	2009	billie	88
cs111	2009	jess	78
cs444	2008	billie	88
cs322	2009	fran	92
		jordan	92

Algorithm requires **projection** and **join** operations on tables.

Join

Variable Elimination Algorithm

Given two tables r_1, r_2 for variables $\text{vars}_1, \text{vars}_2$. The **join** is the table $r_3 = r_1 \times r_2$ for variables $\text{vars}_1 \sqcup \text{vars}_2$ that

- contains all tuples, which restricted to vars_1 are in r_1 , and restricted to vars_2 are in r_2 .

Algorithm Outline

Course	Year	Student	Grade	TA
cs322	2008	fran	77	yuki
cs111	2009	billie	88	sam
cs111	2009	jess	78	chris
cs444	2008	billie	88	chris
cs322	2009	fran	92	chris
		jordan	92	yuki

Course	Year	Student	Grade	TA
cs322	2008	fran	77	yuki
cs111	2009	billie	88	sam
cs111	2009	jess	78	chris
cs444	2008	billie	88	chris
cs322	2009	fran	92	chris
		jordan	92	yuki

Intuition: the constraint constructed in line 5, summarizes the effect that all the constraints involving X have on variables other than X .

19

19

Variable Elimination

19

Autumn 2019

Variable Elimination

22

22

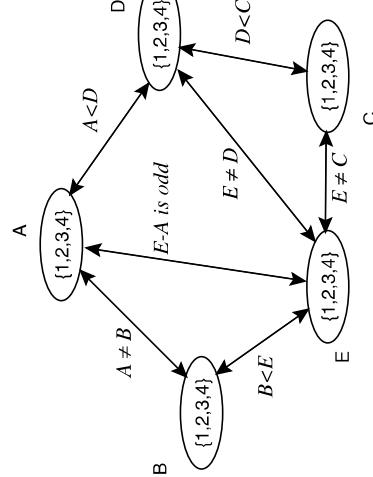
22

22

23

Example network

Example network



Autumn 2019

Variable Elimination

24

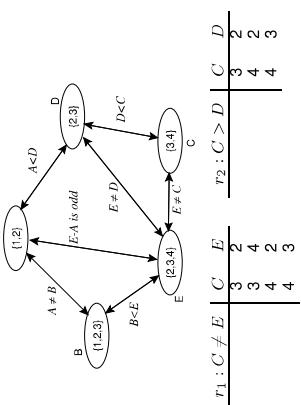
25

Example network

Example network

Example: eliminating C

Example: eliminating C



Autumn 2019

Variable Elimination

26

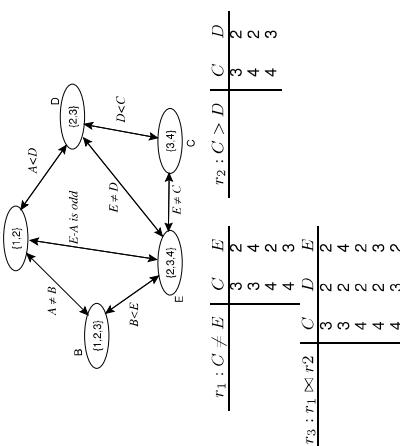
27

Example: eliminating C

Example network

Example: eliminating C

Example: eliminating C



Autumn 2019

Variable Elimination

28

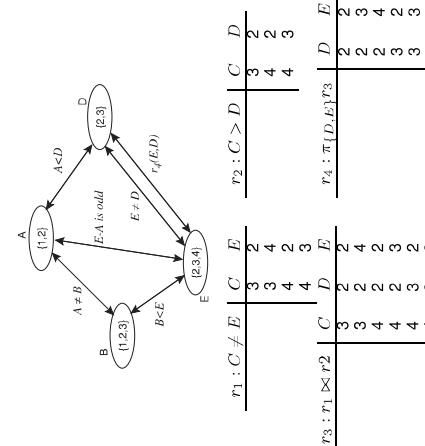
29

Example network

Example: eliminating C

Example network

Example: eliminating C



Autumn 2019

Variable Elimination

29

30

- The algorithm terminates
- The CSP has a solution if and only if the final constraint is non-empty
- The set of all solutions can be generated by joining the final constraint with the intermediate ‘summarizing’ constraints generated in line 5.
- Algorithm operates on extensional constraint representations, therefore
 - constraints must not contain too many tuples (initial and constructed constraints)
 - Worst case: VE is not more efficient than enumerating all possible worlds and checking whether they are solutions.

Properties

Properties

VE Properties

VE Properties

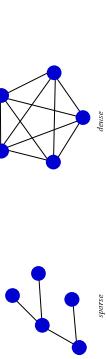
Properties

- The algorithm terminates
- The CSP has a solution if and only if the final constraint is non-empty
- The set of all solutions can be generated by joining the final constraint with the intermediate "summarizing" constraints generated in line 5.
- Algorithm operates on extensional constraint representations, therefore
 - constraints must not contain too many tuples (initial and constructed constraints)
 - Worst case: VE is not more efficient than enumerating all possible worlds and checking whether they are solutions.

Constraint Graph

Consider the graph where

- there is one node for each variable
- two variables are connected when they appear together in one constraint



Autumn 2019

27

Variable Elimination

Then: VE will work better if the constraint graph is sparsely connected

Variable Elimination

Autumn 2019

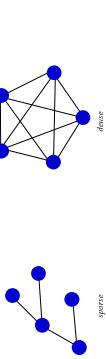
27

Variable Elimination

Constraint Graph

Consider the graph where

- there is one node for each variable
- two variables are connected when they appear together in one constraint



Autumn 2019

27

Variable Elimination

Then: VE will work better if the constraint graph is sparsely connected

Variable Elimination

Autumn 2019

27

Variable Elimination

VE Properties

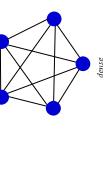
Properties

- The algorithm terminates
- The CSP has a solution if and only if the final constraint is non-empty
- The set of all solutions can be generated by joining the final constraint with the intermediate "summarizing" constraints generated in line 5.
- Algorithm operates on extensional constraint representations, therefore
 - constraints must not contain too many tuples (initial and constructed constraints)
 - Worst case: VE is not more efficient than enumerating all possible worlds and checking whether they are solutions.

Constraint Graph

Consider the graph where

- there is one node for each variable
- two variables are connected when they appear together in one constraint



Autumn 2019

27

Variable Elimination

Then: VE will work better if the constraint graph is sparsely connected

Variable Elimination

Autumn 2019

27

Variable Elimination

So far: all methods systematically explored the state space (possible worlds).

Local Search

Problem: Time and space when search space is large.

Local Search approach:

- explore state space without 'bookkeeping' (where have we been? what still needs to be explored?).
- no success/termination guarantees
- in practice, often the only thing that works

Autumn 2019

Local Search

28

Autumn 2019

Local Search

28

Local Search

Autumn 2019

Local Search

28

Local Search Approach

Local Search Approach

Algorithm Outline

1. Select some node in state space graph as *current state*
2. **while** *current state* is not a solution
 3. *current state* = some neighbor of *current state*



Autumn 2019

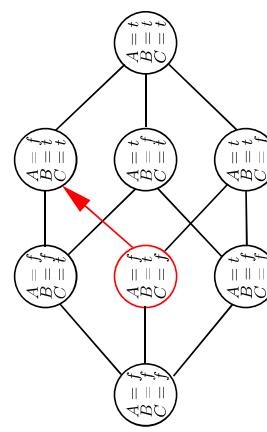
30

Local Search

Local Search Approach

Algorithm Outline

1. Select some node in state space graph as *current state*
2. **while** *current state* is not a solution
 3. *current state* = some neighbor of *current state*



Autumn 2019

Local Search

30

Local Search

Algorithm Outline

1. Select some node in state space graph as *current state*
2. **while** *current state* is not a solution
 3. *current state* = some neighbor of *current state*



Autumn 2019

Local Search

30

Local Search

Local Search Approach

Local Search Approach

Algorithm Outline

1. Select some node in state space graph as *current state*
2. **while** *current state* is not a solution
 current state = some neighbor of *current state*



Autumn 2019

30

Local Search

Local Search Approach

Algorithm Outline

1. Select some node in state space graph as *current state*
2. **while** *current state* is not a solution
 current state = some neighbor of *current state*



Autumn 2019

Local Search

30

Algorithm Outline

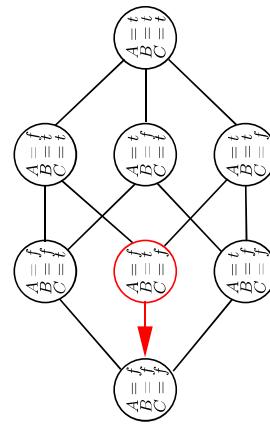
30

Local Search Approach

Algorithm Outline

1. Select some node in state space graph as *current state*
2. **while** *current state* is not a solution
 current state = some neighbor of *current state*

- Make choices in line 1. and 3. completely random
- "Random walk"
- Unlikely to find a solution if state space large with only few solutions



Autumn 2019

30

Local Search

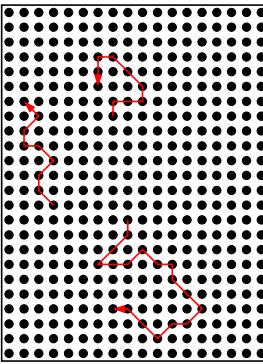
31

Greedy Search

Greedy Search or Hill Climbing:

- Use an *evaluation function* on states
- Examples for evaluation function: number of constraints not satisfied by state
- Always choose neighbor with minimal evaluation function value
- Terminates when all neighbors have higher value than current state; current state is a **local minimum**.

Possible greedy search paths starting from different states:



Autumn 2019

Local Search

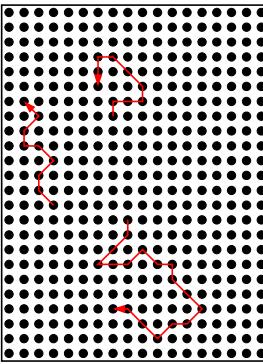
33

Escaping Local Minima

Problem

- Search terminates with local minimum of evaluation function. This may not be a solution to the CSP.

Possible greedy search paths starting from different states:



Autumn 2019

32

Escaping Local Minima

Problem

- Search terminates with local minimum of evaluation function. This may not be a solution to the CSP.

Solution Approaches

- Random restarts: repeat greedy search with several randomly chosen initial states
- Random moves: combine greedy moves with random steps

• Example (a): Small number of random restarts will find global minimum

• Example (b): Make random move when local minimum reached



Autumn 2019

Local Search

33

Solution Approaches

- Random restarts: repeat greedy search with several randomly chosen initial states
- Random moves: combine greedy moves with random steps

• Example (a): Small number of random restarts will find global minimum

• Example (b): Make random move when local minimum reached



Autumn 2019

Local Search

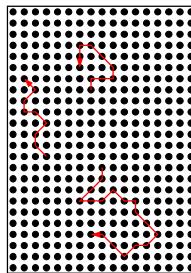
33

Local search

- Maintain an assignment of a value to each variable.

- At each step, select a "neighbor" of the current assignment (e.g., one that improves some heuristic value).

- Stop when a satisfying assignment is found, or return the best assignment found.



Principle

Select the variable-value pair that gives the highest improvement.

Naive approach

- Linearly scan all variables and for each value of each variable determine the improvement (how many fewer constraints are violated).

Requires:

- What is a neighbor?
- Which neighbor should be selected?

Autumn 2019

Local Search

34

Autumn 2019

Local Search

35

Principle

Select the variable-value pair that gives the highest improvement.

Naive approach

- Linearly scan all variables and for each value of each variable determine the improvement (how many fewer constraints are violated).

Principle

- First: choose variable
- Second: choose state

Data structure

- Maintain priority queue of variables; weight is the number of participating conflicts.
- After selecting a variable, pick the value minimizes the number of conflicts.
- Update weights of variables that participate in a conflict that is changed.

Autumn 2019

35

Local Search

36

Simulated annealing**Algorithm**

- Pick a variable at random and a new value at random.
- If it is an improvement, adopt it.
- If it isn't an improvement, adopt it probabilistically depending on a temperature parameter, T .
 - With current assignment n and proposed assignment n' , we move to n' with probability $e^{(h(n') - h(n))/T}$
- Reduce the temperature.

Simulated annealing**Algorithm**

- Pick a variable at random and a new value at random.
- If it is an improvement, adopt it.
- If it isn't an improvement, adopt it probabilistically depending on a temperature parameter, T .
 - With current assignment n and proposed assignment n' , we move to n' with probability $e^{(h(n') - h(n))/T}$

Temperature	1-worse	2-worse	3-worse
10	0.91	0.81	0.74
1	0.37	0.14	0.05
0.25	0.02	0.0003	0.00005
0.1	0.00005	0	0

Autumn 2019

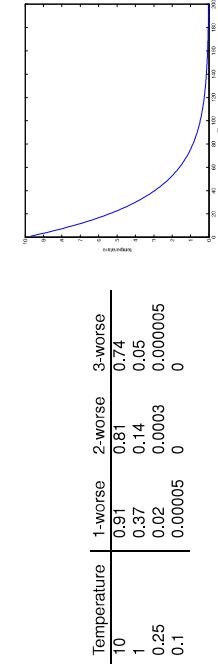
37

Local Search

37

Simulated annealing**Algorithm**

- Pick a variable at random and a new value at random.
- If it is an improvement, adopt it.
- If it isn't an improvement, adopt it probabilistically depending on a temperature parameter, T .
 - With current assignment n and proposed assignment n' , we move to n' with probability $e^{(h(n') - h(n))/T}$
- Reduce the temperature.

Propositional Logic Basics

Autumn 2019

37

Propositional Logic Basics

38

[Previously ...](#)

Intensional representation of constraints:

$$\begin{array}{l} A < B \\ \text{Teacher_AD} = \text{Teacher_MI} \rightarrow \text{Time_AD} \neq \text{Time_MI} \\ \dots \end{array}$$

CSP algorithms (arc-consistency algorithm) need to perform certain operations:

- test whether a certain value for one variable is **consistent** with a given constraint (and certain values for other variables)

To implement this:

- need a **formal language for representing constraints**

Propositional Logic

- provides a formal language for representing constraints on **binary variables**.

Boolean variables are now seen as **atomic propositions**. Convention: start with lowercase letter.

Constraints	Logic
$A = \text{true}$	a
$A = \text{false}$	$\neg a$

Propositions

Using **logical connectives** more complex propositions are constructed:

$\neg p$	$(p \wedge q)$	$(p \vee q)$	$(p \rightarrow q)$
not p	p and q	p or q	p implies q

Example

"If it rains I'll take my umbrella, or I'll stay home"

Autumn 2019

Propositional Logic: Basics

36

Propositional Logic: Syntax

Atomic Propositions

Boolean variables are now seen as **atomic propositions**. Convention: start with lowercase letter.

Constraints	Logic
$A = \text{true}$	a
$A = \text{false}$	$\neg a$

Propositions

Using **logical connectives** more complex propositions are constructed:

$\neg p$	$(p \wedge q)$	$(p \vee q)$	$(p \rightarrow q)$
not p	p and q	p or q	p implies q

A set of propositions is also called a **Knowledge Base**

Example

"If it rains I'll take my umbrella, or I'll stay home"

Autumn 2019

Propositional Logic: Basics

39

Propositional Logic: Semantics |

$\pi(a_i) \in \{\text{true}, \text{false}\}$

An interpretation defines a truth value for all propositions:

$\pi(p)$	$\pi(q)$	$\pi(p \wedge q)$	$\pi(p \vee q)$
true	true	true	true
true	false	false	true
false	true	false	true
false	false	false	false

An interpretation defines a truth value for all propositions:

$\pi(p)$	$\pi(q)$	$\pi(p \rightarrow q)$
true	true	true
true	false	false
false	true	true
false	false	true

Simple Example

Models

A **model** of a proposition (a knowledge base) is an interpretation in which the proposition (all the propositions in the knowledge base) is true.

Propositions as constraints: a model is a possible world that satisfies the constraint.

$$KB = \left\{ \begin{array}{l} p \leftarrow q, \\ q, \\ r \leftarrow s. \end{array} \right.$$

Logical consequence

Model?

$$KB \models g$$

(whenever KB is true, then g also is true).

Example

$KB = \{ \text{man} \rightarrow \text{mortal}, \text{man} \}$. Then $KB \models \text{mortal}$

Autumn 2019

Propositional Logic: Basics

41

42

Simple Example

Simple Example

$$KB = \begin{cases} p \leftarrow q, \\ q, \\ r \leftarrow s. \end{cases}$$

	p	q	r	s
I_1	true	true	false	false
I_2	false	false	false	false
I_3	true	false	false	false
I_4	true	true	false	false
I_5	true	true	false	true

Which of p, q, r, q logically follow from KB ?

Autumn 2019

Propositional Logic Basics

42

Simple Example

$$KB = \begin{cases} p \leftarrow q, \\ q, \\ r \leftarrow s. \end{cases}$$

	p	q	r	s
I_1	true	true	false	false
I_2	false	false	false	false
I_3	true	false	false	false
I_4	true	true	false	false
I_5	true	true	false	true

Which of p, q, r, q logically follow from KB ?

$KB \models p, KB \models q, KB \not\models r, KB \not\models s$

Simple Example

$$KB = \begin{cases} p \leftarrow q, \\ q, \\ r \leftarrow s. \end{cases}$$

	p	q	r	s
I_1	true	true	false	false
I_2	false	false	false	false
I_3	true	false	false	false
I_4	true	true	false	false
I_5	true	true	false	true

Which of p, q, r, q logically follow from KB ?

Propositional Logic Basics

42

Machine Intelligence

Lecture 4: Reasoning under uncertainty - probabilities

Thomas Dyrre Nielsen

Aalborg University

Autumn 2019

Propositional Logic Basics

42

1

Topics:

- Introduction
 - Search-based methods
 - Constrained satisfaction problems
 - Representing domains endowed with uncertainty
 - Bayesian networks
 - Machine learning
 - Planning
 - Multi-agent systems
- Possible worlds are possible or impossible (according to given constraints/propositions)
- Current state is fully known
- Propositions are fully known/believed, or unknown
- Actions have deterministic effects
- Generalization: soft constraints – possible worlds are more or less desirable

Certainty in Search

Assumptions for using search for solving planning problems:

Autumn 2019

2

3

Certainty in CSP and Logic

Autumn 2019

1

More realistic scenarios

- Agents do not observe the world perfectly
- Actions have uncertain effects
- Propositions are believed only with a certain confidence

Degrees of Belief for Propositions

In reality, states of knowledge may better be represented by degrees of belief:

```
Bel(light_on ← switch_on ∧ breaker_up) = 0.7
Bel(~umbrella → rain) = 1.5
Bel(umbrella → rain) = 0.2
Bel(global_warming) = 0.8
```

Question: what rules must (rational) degrees of belief obey?

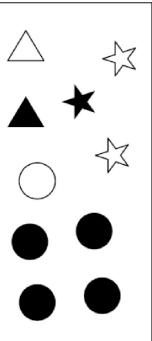
MI Autumn 2019

4

Basic Probability Calculus

Probability Measures

Probability theory is built on the foundation of variables and worlds.



Worlds described by the variables:

- Filled: {true, false}
 - Shape: {circle, triangle, star}
- as well as position.

Probability measures

Ω : set of all possible worlds (for a given, fixed set of variables). A probability measure over Ω , is a function P , that assigns **probability values**

$$P(\Omega') \in [0, 1]$$

to subsets $\Omega' \subseteq \Omega$, such that

Axiom 1: $P(\Omega) = 1$.

Axiom 2: if $\Omega_1 \cap \Omega_2 = \emptyset$, then $P(\Omega_1 \cup \Omega_2) = P(\Omega_1) + P(\Omega_2)$.

MI Autumn 2019

5

Basic Probability Calculus

Simplification for finite Ω

If all variables have a finite domain, then

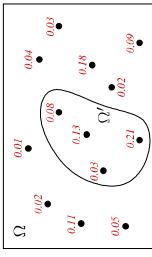
- Ω is finite, and
- a probability distribution is defined by assigning a probability value

$$P(\omega)$$

to each individual possible world $\omega \in \Omega$.

For any $\Omega' \subseteq \Omega$ then

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



$$P(\Omega') = 0.08 + 0.13 + 0.03 + 0.21 = 0.45$$

From now on, we will only consider variables with finite domains.

MI Autumn 2019

Basic Probability Calculus

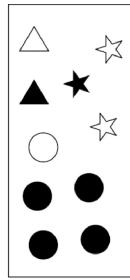
Probabilities of propositions

Probabilities of Propositions

A probability distribution over possible worlds defines probabilities for propositions α :

$$P(\alpha) = P(\{\omega \in \Omega \mid \omega \models \alpha\}) \\ = \sum_{\omega : \alpha \text{ is true in } \omega} P(\omega)$$

Example



Assume probability for each world is 0.1:

- $P(\text{Shape} = \text{circle}) = 0.5$
- $P(\text{Filled} = \text{false}) =$

MI Autumn 2019

7

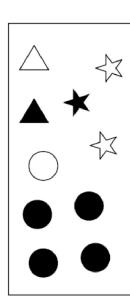
Basic Probability Calculus

Probabilities of Propositions

A probability distribution over possible worlds defines probabilities for propositions α :

$$P(\alpha) = P(\{\omega \in \Omega \mid \omega \models \alpha\}) \\ = \sum_{\omega : \alpha \text{ is true in } \omega} P(\omega)$$

Example



Assume probability for each world is 0.1:

- $P(\text{Shape} = \text{circle}) =$
- $P(\text{Filled} = \text{false}) =$

MI Autumn 2019

7

Basic Probability Calculus

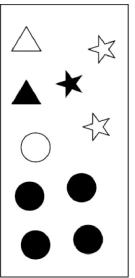
Probabilities of propositions

Probabilities of propositions

A probability distribution over possible worlds defines probabilities for propositions α :

$$P(\alpha) = P(\{\omega \in \Omega \mid \omega \models \alpha\}) \\ = \sum_{\omega : \alpha \text{ is true in } \omega} P(\omega)$$

Example



Assume probability for each world is 0.1:

- $P(\text{Shape} = \text{circle}) = 0.5$
- $P(\text{Filled} = \text{false}) = 0.4$
- $P(\text{Shape} = c \wedge \text{Filled} = f) = 0.1$

[MI Autumn 2019]

Basic Probability Calculus

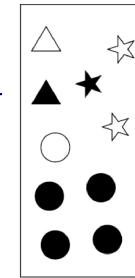
Probabilities of propositions

Probabilities of Propositions

A probability distribution over possible worlds defines probabilities for propositions α :

$$P(\alpha) = P(\{\omega \in \Omega \mid \omega \models \alpha\}) \\ = \sum_{\omega : \alpha \text{ is true in } \omega} P(\omega)$$

Example



Assume probability for each world is 0.1:

- $P(\text{Shape} = \text{circle}) = 0.5$
- $P(\text{Filled} = \text{false}) = 0.4$
- $P(\text{Shape} = c \wedge \text{Filled} = f) = 0.1$

[MI Autumn 2019]

Basic Probability Calculus

Basic probability axioms

Axiom

If \mathcal{A} and \mathcal{B} are disjoint, then $P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B})$.

Consider a deck with 52 cards. If $\mathcal{A} = \{2, 3, 4, 5\}$ and $\mathcal{B} = \{7, 8\}$, then

$$P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B}) = 0.12 + 0.16 = 0.28.$$

[MI Autumn 2019]

Basic Probability Calculus

Basic probability axioms

Axiom

If \mathcal{A} and \mathcal{B} are disjoint, then $P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B})$.

Consider a deck with 52 cards. If $\mathcal{A} = \{2, 3, 4, 5\}$ and $\mathcal{B} = \{7, 8\}$, then

$$P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B}) = 0.12 + 0.16 = 0.28.$$

[MI Autumn 2019]

Basic Probability Calculus

Basic probability axioms

Axiom

If \mathcal{A} and \mathcal{B} are disjoint, then $P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B})$.

Consider a deck with 52 cards. If $\mathcal{A} = \{2, 3, 4, 5\}$ and $\mathcal{B} = \{7, 8\}$, then

$$P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B}) = 0.12 + 0.16 = 0.28.$$

[MI Autumn 2019]

Basic Probability Calculus

Consider a deck with 52 cards. If $\mathcal{A} = \{2, 3, 4, 5\}$ and $\mathcal{B} = \{7, 8\}$, then

$$P(\mathcal{A} \cup \mathcal{B}) = P(\mathcal{A}) + P(\mathcal{B}) = \frac{4}{13} + \frac{2}{13} = \frac{6}{13}.$$

If \mathcal{C} and \mathcal{D} are disjoint, then $P(\mathcal{C} \cup \mathcal{D}) = P(\mathcal{C}) + P(\mathcal{D})$.

Example

If $\mathcal{C} = \{2, 3, 4, 5\}$ and $\mathcal{D} = \{7, 8\}$, then

$$P(\mathcal{C} \cup \mathcal{D}) =$$

[MI Autumn 2019]

Basic Probability Calculus

More generally

If \mathcal{C} and \mathcal{D} are not disjoint, then $P(\mathcal{C} \cup \mathcal{D}) = P(\mathcal{C}) + P(\mathcal{D}) - P(\mathcal{C} \cap \mathcal{D})$.

Example

If $\mathcal{C} = \{2, 3, 4, 5\}$ and $\mathcal{D} = \{7, 8\}$, then

$$P(\mathcal{C} \cup \mathcal{D}) =$$

[MI Autumn 2019]

Basic Probability Calculus

More generally

If \mathcal{C} and \mathcal{D} are disjoint, then $P(\mathcal{C} \cup \mathcal{D}) = P(\mathcal{C}) + P(\mathcal{D})$.

Example

If $\mathcal{C} = \{2, 3, 4, 5\}$ and $\mathcal{D} = \{7, 8\}$, then

$$P(\mathcal{C} \cup \mathcal{D}) =$$

[MI Autumn 2019]

Basic Probability Calculus

More generally

If \mathcal{C} and \mathcal{D} are not disjoint, then $P(\mathcal{C} \cup \mathcal{D}) = P(\mathcal{C}) + P(\mathcal{D}) - P(\mathcal{C} \cap \mathcal{D})$.

Example

If $\mathcal{C} = \{2, 3, 4, 5\}$ and $\mathcal{D} = \{7, 8\}$, then

$$P(\mathcal{C} \cup \mathcal{D}) =$$

[MI Autumn 2019]

Basic Probability Calculus

Basic probability axioms

Probability Updating

Axiom

If A and B are disjoint, then $P(A \cup B) = P(A) + P(B)$.

Example

Consider a deck with 52 cards. If $A = \{2, 3, 4, 5\}$ and $B = \{7, 8\}$, then

$$P(A \cup B) = P(A) + P(B) = \frac{4}{13} + \frac{2}{13} = \frac{6}{13}.$$

How should the probabilities be updated, when I observe $\Omega'?$:

More generally

If C and D are not disjoint, then $P(C \cup D) = P(C) + P(D) - P(C \cap D)$.

Example

If $C = \{2, 3, 4, 5\}$ and $D = \{\spadesuit\}$, then

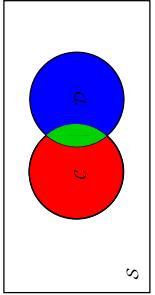
$$P(C \cup D) = \frac{4}{13} + \frac{1}{4} - \frac{4}{52} = \frac{25}{52}.$$

Basic Probability Calculus

MI Autumn 2019

Basic Probability Calculus

9



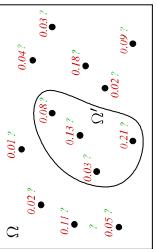
Probability Updating

Given new information (**evidence**), degrees of belief change.

Evidence can consist of:

- learning that a certain proposition p is true ("switch up")
 - measuring the value of some variable ("room_ai = 0.2.90")
 - obtaining partial information on the value of some variable ("room_ai ≠ 0.2.90")
 - ...

- In all cases: evidence can be represented as the set of possible world Ω' not ruled out by the observation.
- How should the probabilities be updated, when I observe $\Omega'?$:



- worlds that are not consistent with evidence have probability 0

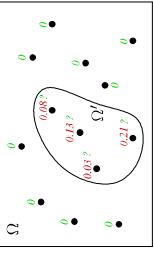
Probability Updating

Given new information (**evidence**), degrees of belief change.

Evidence can consist of:

- learning that a certain proposition p is true ("switch up")
 - measuring the value of some variable ("room_ai = 0.2.90")
 - obtaining partial information on the value of some variable ("room_ai ≠ 0.2.90")
 - ...

- In all cases: evidence can be represented as the set of possible world Ω' not ruled out by the observation.
- How should the probabilities be updated, when I observe $\Omega'?$:



- worlds that are not consistent with evidence have probability 0

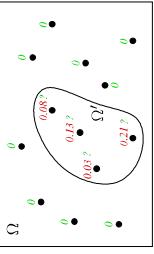
Probability Updating

Given new information (**evidence**), degrees of belief change.

Evidence can consist of:

- learning that a certain proposition p is true ("switch up")
 - measuring the value of some variable ("room_ai = 0.2.90")
 - obtaining partial information on the value of some variable ("room_ai ≠ 0.2.90")
 - ...

- In all cases: evidence can be represented as the set of possible world Ω' not ruled out by the observation.
- How should the probabilities be updated, when I observe $\Omega'?$:



- worlds that are not consistent with evidence have probability 0

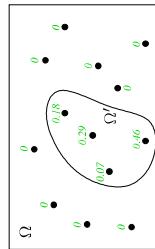
Probability Updating

Given new information (**evidence**), degrees of belief change.

Evidence can consist of:

- learning that a certain proposition p is true ("switch up")
 - measuring the value of some variable ("room_ai = 0.2.90")
 - obtaining partial information on the value of some variable ("room_ai ≠ 0.2.90")
 - ...

- In all cases: evidence can be represented as the set of possible world Ω' not ruled out by the observation.
- How should the probabilities be updated, when I observe $\Omega'?$:



- worlds that are not consistent with evidence have probability 0
- the probabilities of worlds consistent with evidence are proportional to their probability before observation, and they must sum to 1

MI Autumn 2019

MI Autumn 2019

9

$$P(S=\text{circ.} \mid F\text{ill} = f) = \frac{P(S=\text{circ.} \wedge F\text{ill} = f)}{P(F\text{ill} = f)}$$

$$= \frac{0.1}{0.4} = 0.25$$

(probability for each world is 0.1)

Basic Probability Calculus

MI Autumn 2019

10



Conditional probability

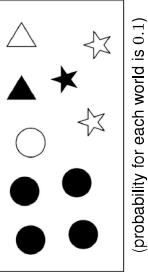
Conditional probability

Definition

The **conditional probability** of proposition p given e is

$$P(p \mid e) = \frac{P(p \wedge e)}{P(e)}$$

Example



(probability for each world is 0.1)

$$P(S=circ. \mid Fill = f) = \frac{P(S=circ. \wedge Fill = f)}{P(Fill = f)}$$

$$= \frac{0.1}{0.4} = 0.25$$

What is the probability of $P(S=\text{star} \mid Fill = f)$?

[M] Autumn 2019

Basic Probability Calculus

[M] Autumn 2019 Basic Probability Calculus 10

Two important rules

Bayes rule

For propositions p, e :

$$P(p \mid e) = \frac{P(e \wedge p)}{P(e)} = \frac{P(e \mid p)P(p)}{P(e)}$$

Bayes rule

For propositions p, e :

$$P(p \mid e) = \frac{P(e \wedge p)}{P(e)} = \frac{P(e \mid p)P(p)}{P(e)}$$

Example

A doctor observes symptoms and wishes to find the probability of a disease:

$$P(\text{disease} \mid \text{symp.}) = \frac{P(\text{symp.} \mid \text{disease})P(\text{disease})}{P(\text{symp.})}$$

[M] Autumn 2019 Basic Probability Calculus 11

Two important rules

Bayes rule

For propositions p, e :

$$P(p \mid e) = \frac{P(e \wedge p)}{P(e)} = \frac{P(e \mid p)P(p)}{P(e)}$$

[M] Autumn 2019 Basic Probability Calculus 11

Bayes rule

For propositions p, e :

$$P(p \mid e) = \frac{P(e \wedge p)}{P(e)} = \frac{P(e \mid p)P(p)}{P(e \wedge p) + P(e \wedge \neg p)}$$

Bayes rule

For propositions p, e :

$$P(p \mid e) = \frac{P(e \wedge p)}{P(e)} = \frac{P(e \mid p)P(p)}{P(e \wedge p) + P(e \wedge \neg p)}$$

[M] Autumn 2019 Basic Probability Calculus 11

Chain rule

A doctor observes symptoms and wishes to find the probability of a disease:

$$P(\text{disease} \mid \text{symp.}) = \frac{P(\text{symp.} \mid \text{disease})P(\text{disease})}{P(\text{symp.})}$$

Both rules are immediate consequences of the definition of conditional probability!

[M] Autumn 2019 Basic Probability Calculus 11

[M] Autumn 2019 Basic Probability Calculus 11

Random Variables

Variables defining possible worlds on which probabilities are defined are called **random variables**.

Distributions

For a random variable A , and $a \in D_A$ we have the probability

$$P(A = a) = P(\{\omega \in \Omega \mid A = a \text{ in } \omega\})$$

The **probability distribution of A** is the function on D_A that maps a to $P(A = a)$. The distribution of A is denoted

$$P(A)$$

Joint Distributions

Extension to several random variables:

$$P(A_1, \dots, A_k)$$

is the **joint distribution of A_1, \dots, A_k** . The joint distribution maps tuples (a_1, \dots, a_k) with

$$a_i \in D_{A_i}, \quad P(A_1 = a_1, \dots, A_k = a_k)$$

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

Basic Probability Calculus

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

Bayes' rule for variables

Bayes' rule can also be written for variables:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Bayes' rule for variables

Bayes' rule can also be written for variables:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{\sum_A P(A|B)}$$

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

Bayes' rule for variables

Bayes' rule can also be written for variables:

$$P(A|B) = \frac{P(B|A)P(A)}{\sum_A P(B|A)P(A)} = \frac{P(B|A)P(A)}{\sum_A P(A|B)} = \frac{P(B|A)P(A)}{\sum_A P(B|A)P(A)}$$

Example

Consider the variables

- Temp : sp(Temp) = $\{l, m, h\}$
- Sensor : sp(Sensor) = $\{l, m, h\}$

$$P(\text{Temp}) = (0.1, 0.6, 0.3)$$

	l	m	h
Usage	0.8	0.15	0.05
Sensor	0.15	0.8	0.1
Temp	0.05	0.1	0.85

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

[MI Autumn 2019]

Basic Probability Calculus

Basic Probability Calculus

Example Continued

Conditionally Independent Variables

$P(\text{Hair length}, \text{Height} \mid \text{Sex} = \text{female})$, $P(\text{Height} \mid \text{Sex} = \text{female})$,
 $P(\text{Height} \mid \text{Hair length}, \text{Sex} = \text{female})$:

Height	Hair length	long	short
tall	0.14	0.06	0.2
medium	0.56	0.24	0.8

\rightsquigarrow Hair length and Height are independent given Sex=female.

Also: Hair length and Height are independent given Sex=male.

\rightsquigarrow Hair length and Height are independent given Sex.

Definition of Conditional Independence

The variables A_1, \dots, A_n are **conditionally independent** of the variables B_1, \dots, B_m given C_1, \dots, C_k , if

$$P(A_1, \dots, A_n \mid B_1, \dots, B_m, C_1, \dots, C_k) = P(A_1, \dots, A_n \mid C_1, \dots, C_k)$$

Agenda: use conditional independence to facilitate specification of probability distributions on complex state spaces

Mi Autumn 2019 20 Independence

Tentative course overview

Machine Intelligence

Lecture 5: Bayesian networks

Thomas Dyhre Nielsen
Aalborg University

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- **Bayesian networks**
- Machine learning
- Planning
- Multi-agent systems

Mi Autumn 2019 21 Independence

Tentative course overview

Machine Intelligence

Lecture 5: Bayesian networks

Thomas Dyhre Nielsen
Aalborg University

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- **Bayesian networks**
- Machine learning
- Planning
- Multi-agent systems

Mi Autumn 2019 1 Independence

Example

Random variables (all Boolean):

Tampering fire alarm has been tampered with
Fire fire in the building
Alarm fire alarm ringing
Smoke smoke in the building
Leaving people leaving the building
Report report of people leaving the building

Bayesian Networks

$P(\text{Tampering}, \text{Fire}, \text{Alarm}, \text{Smoke}, \text{Leaving}, \text{Report}) =$
 $P(\text{Tampering}) \cdot$
 $P(\text{Fire} \mid \text{Tampering}) \cdot$
 $P(\text{Alarm} \mid \text{Tampering}, \text{Fire}) \cdot$
 $P(\text{Smoke} \mid \text{Tampering}, \text{Fire}, \text{Alarm}) \cdot$
 $P(\text{Leaving} \mid \text{Tampering}, \text{Fire}, \text{Alarm}, \text{Smoke}) \cdot$
 $P(\text{Report} \mid \text{Tampering}, \text{Fire}, \text{Alarm}, \text{Smoke}, \text{Leaving})$

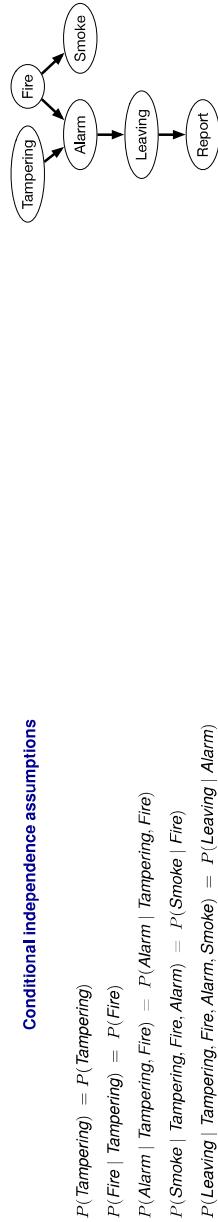
Joint distribution according to chain rule:

Mi Autumn 2019 3 Bayesian Networks

Example continued

Graphical Representation

Representation of conditional dependencies in a graph:

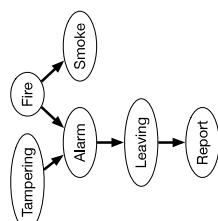


$$P(\text{Tampering}, \text{Fire}, \text{Alarm}, \text{Smoke}, \text{Leaving}, \text{Report}) = \\ P(\text{Tampering}) \cdot P(\text{Fire}) \cdot P(\text{Alarm} \mid \text{Tampering}, \text{Fire}) \cdot P(\text{Smoke} \mid \text{Fire}) \cdot \\ P(\text{Leaving} \mid \text{Report})$$

MI Autumn 2019 Bayesian Networks 4

Graphical Representation

Representation of conditional dependencies in a graph:



The graph is **directed** and **acyclic**.

Bayesian Network

A **Bayesian Network** for variables A_1, \dots, A_k consists of

- a directed acyclic graph with nodes A_1, \dots, A_k
- for each node a **conditional probability table** specifying the conditional distribution $P(A_i \mid \text{parents}(A_i))$ [parents(A_i) denotes the parents of A_i in the graph])

and through the chain rule provides a compact representation of a joint probability distribution.

MI Autumn 2019 Bayesian Networks 5

MI Autumn 2019 Bayesian Networks 5

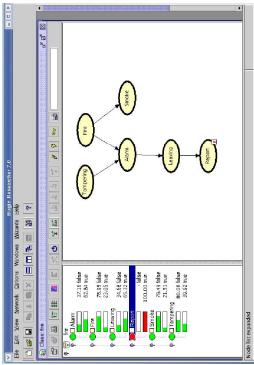
Observations and queries

A Bayesian network specifies a joint distribution from which arbitrary conditional probabilities can be derived.

Inference

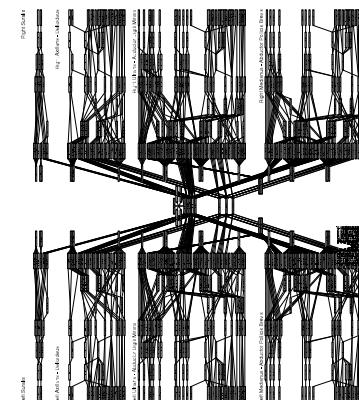
The most common task is to compute the posterior distribution over a query variable A given the observed values of some evidence nodes $E_i = e_i$, for $i = 1, \dots, l$:

$$P(A \mid E_1 = e_1, \dots, E_l = e_l).$$



MI Autumn 2019 Bayesian Networks 6

The Munin network



Characteristics:

- Approximately 1100 variables.
- Each variable has between 2 and 20 states.
- 10⁶⁰⁰ possible state configurations!

Construction via chain rule

- put the random variables in some order
- write the joint distribution using chain rule
- simplify conditional probability factors by conditional independence assumptions. That determines the parents of each node, i.e. the graph structure
- specify the conditional probability tables

Note: the structure of the resulting network strongly depends on the chosen order of the variables.

Construction via causality

- Draw and edge from variable A to variable B if A has a direct causal influence on A .
- Note: this may not always be possible:
- $\text{Inflation} \rightarrow \text{salaries}$ or $\text{salaries} \rightarrow \text{Inflation}$?
 - Rain doesn't cause Sun, and Sun doesn't cause Rain, but they are not independent either!

MI Autumn 2019 Bayesian Networks 7

MI Autumn 2019 Bayesian Networks 8

Reasoning under uncertainty 1

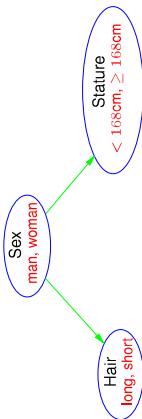
Transmission of evidence



- If there has been a flooding does that tell me something about the amount of rain that has fallen?
- The water level is high: If there has been a flooding does that tell me anything new about the amount of rain that has fallen?

MI Autumn 2019 Transmission of evidence 9

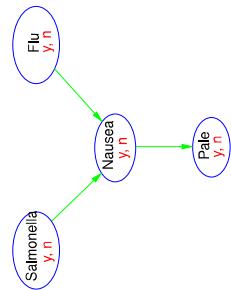
Reasoning under uncertainty 2



- If a person has long hair does that say something about his/her stature?
- It is a woman: If she has long hair does that say something about her stature?

MI Autumn 2019 Transmission of evidence 10

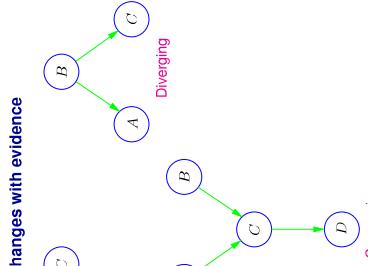
Reasoning under uncertainty 3



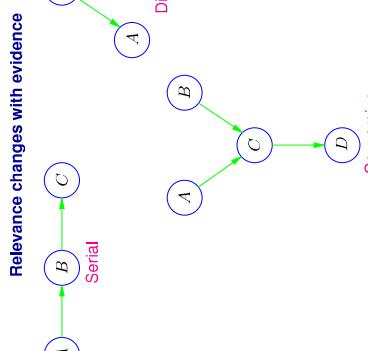
- Does salmonella have an impact on Flu?
- If a person is Pale, does salmonella then have an impact on Flu?

MI Autumn 2019 Transmission of evidence 11

Transmission of evidence 1



Relevance changes with evidence



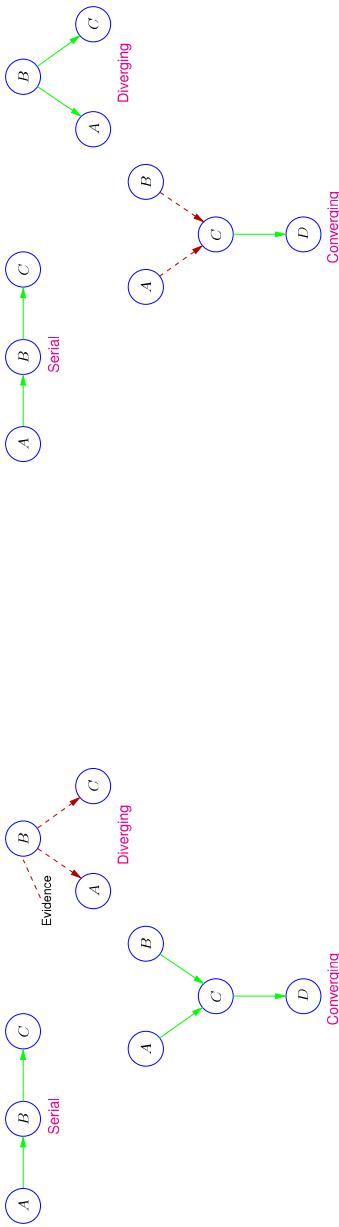
MI Autumn 2019 Transmission of evidence 12

MI Autumn 2019 Transmission of evidence 13

Transmission of evidence 1

Transmission of evidence 1

Relevance changes with evidence



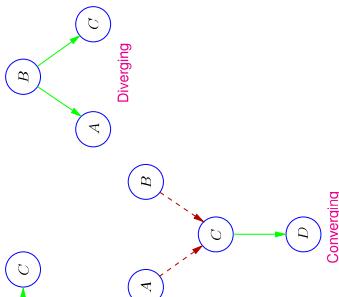
MI Autumn 2019

12

Transmission of evidence

12

Relevance changes with evidence



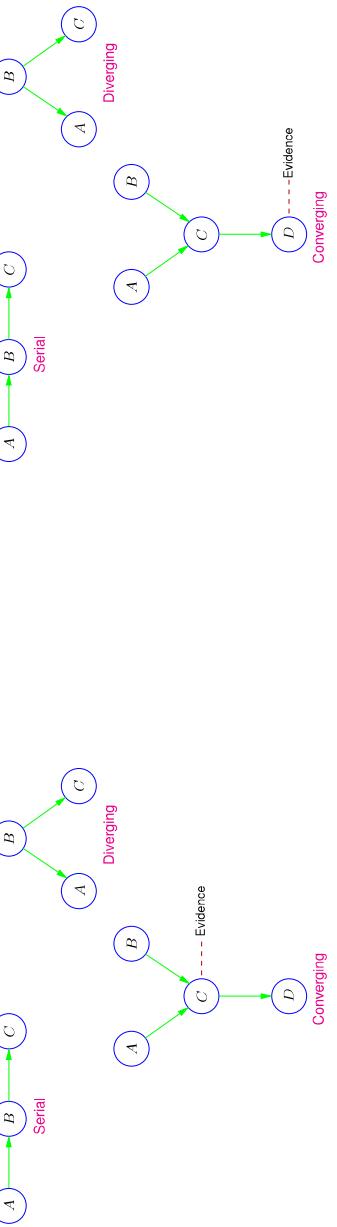
Transmission of evidence

12

Transmission of evidence 1

Transmission of evidence 1

Relevance changes with evidence



MI Autumn 2019

12

Transmission of evidence

12

Transmission of evidence 1

Transmission of evidence 1

Summary of transmission rules (d-separation rules)



Rules for transmission of evidence

- Evidence may be transmitted through a serial or diverging connection unless it is instantiated.
- Evidence may be transmitted through a converging connection only if either the variable in the connection or one of its descendants has received evidence.

Can knowledge of *A* have an impact on our knowledge of *J*?

MI Autumn 2019

12

Transmission of evidence

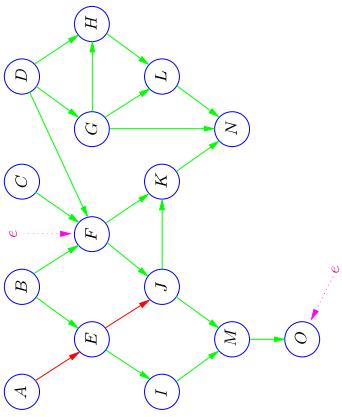
12

Transmission of evidence

13

Transmission of evidence 2

Transmission of evidence 2



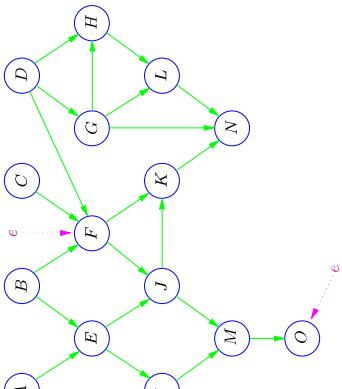
Can knowledge of A have an impact on our knowledge of J ? yes!

MI Autumn 2019

Transmission of evidence

13

13



Can knowledge of A have an impact on our knowledge of B ? yes!

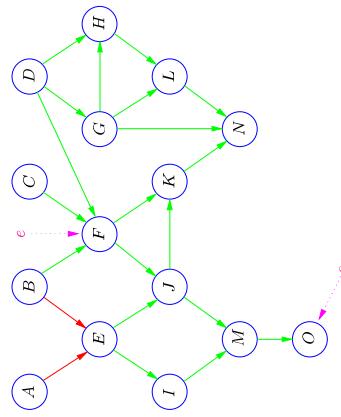
MI Autumn 2019

Transmission of evidence

13

Transmission of evidence 2

Transmission of evidence 2



Can knowledge of A have an impact on our knowledge of B ? yes!

MI Autumn 2019

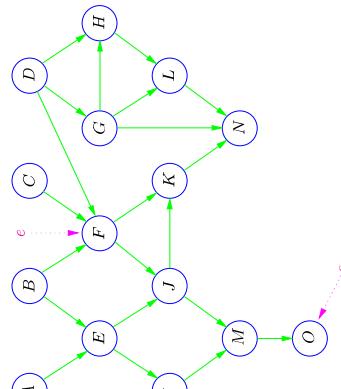
Transmission of evidence

13

13

Transmission of evidence 2

Transmission of evidence 2



Can knowledge of A have an impact on our knowledge of G ? ?

MI Autumn 2019

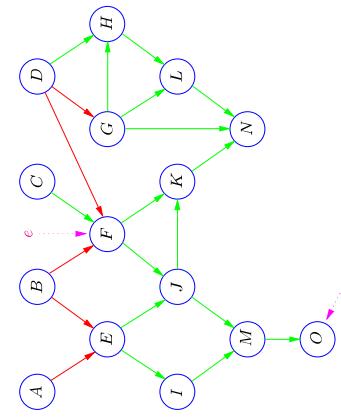
Transmission of evidence

13

13

Transmission of evidence 2

Transmission of evidence 3



Can knowledge of A have an impact on our knowledge of G ? yes!

MI Autumn 2019

Transmission of evidence

13

Is E d-separated from A ?

MI Autumn 2019

Transmission of evidence

14

The d-separation theorem

Use of d-Separation

Theorem

For all pairwise disjoint sets A, B, C of nodes in a Bayesian network:

$$\text{If } C \text{ d-separates } A \text{ from } B, \text{ then } P(A | B, C) = P(A | C).$$

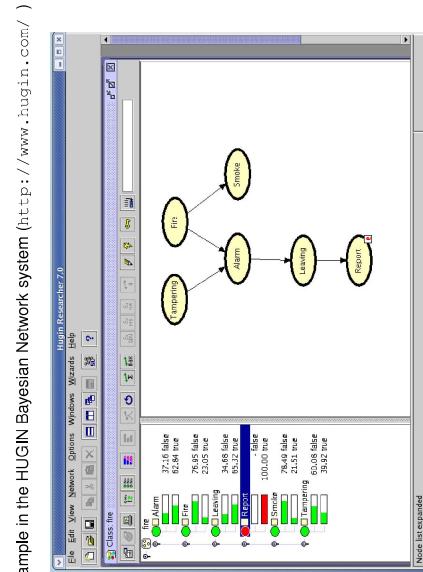
There are no more general graphical conditions than d-separation for which such a result holds.

Why is d-separation important?

- Gaining insight: given a (correct) Bayesian network model, can derive insight into the dependencies among the variables
- Debugging a model: given a Bayesian network model, check whether entailed independence relations are plausible
- Correctness of algorithms: certain computational procedures depend on validity of special independence relations

MI Autumn 2019 15 Transmission of evidence

Fire Example



MI Autumn 2019 16 Transmission of evidence

Specifying the structure of a Bayesian network

Specifying the structure of a Bayesian network

MI Autumn 2019 17 Transmission of evidence

Specifying the structure of a Bayesian network

Building models

Building models

Milk from a cow may be infected. To detect whether or not the milk is infected, you can apply a test which may either give a positive or a negative test result. The test is not perfect: It may give **false positives** as well as **false negatives**.

Milk from a cow may be infected. To detect whether or not the milk is infected, you can apply a test which may either give a positive or a negative test result. The test is not perfect: It may give **false positives** as well as **false negatives**.



MI Autumn 2019 18 Specifying the structure of a Bayesian network

Specifying the structure of a Bayesian network

Specifying the structure of a Bayesian network

MI Autumn 2019 19 Specifying the structure of a Bayesian network

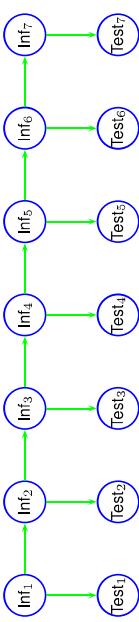
7-day model I

Infections develop over time:



7-day model I

Infections develop over time:

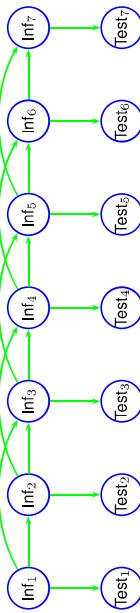


But what if yesterday's Inf-state has an impact on tomorrow's Inf-state?

MI Autumn 2019 Specifying the structure of a Bayesian network 19

7-day model II

Non-Markov relations



Yesterday's Inf-state has an impact on tomorrow's Inf-state.

MI Autumn 2019 Specifying the structure of a Bayesian network 20

Sore throat

I wake up one morning with a sore throat. It may be the beginning of a cold or I may suffer from angina. If it is a severe angina, then I will not go to work. To gain more insight, I can take my temperature and look down my throat for yellow spots.

I wake up one morning with a sore throat. It may be the beginning of a cold or I may suffer from angina. If it is a severe angina, then I will not go to work. To gain more insight, I can take my temperature and look down my throat for yellow spots.

Hypothesis variables

Cold? - {n, y}
Angina? - {no, mild, severe}

Information variables

Sore throat? - {n, y}
See spots? - {n, y}
Fever? - {no, low, high}

MI Autumn 2019 Specifying the structure of a Bayesian network 22

Specifying the structure of a Bayesian network

Specifying the structure of a Bayesian network

Model for sore throat

Model for sore throat



MI Autumn 2019 Specifying the structure of a Bayesian network 23

Insemination of a cow

Six weeks after the insemination of a cow, there are two tests: a **Blood test** and a **Urine test**.



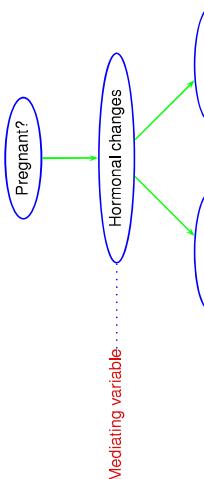
Check the conditional independences

If we know that the cow is pregnant, will a negative blood test then change our expectation for the urine test?

If it **will**, then the model does not reflect reality!

MI Autumn 2019 Specifying the structure of a Bayesian network 24

Insemination of a cow: A more correct model



MI Autumn 2019 Specifying the structure of a Bayesian network 25

But does this actually make a difference?

But does this actually make a difference?

Assume that both **tests** are **negative** in the **incorrect model**:
This will overestimate the probability for **Pregnant? = n**.
Pregnant? --> Blood test
Pregnant? --> Urine test

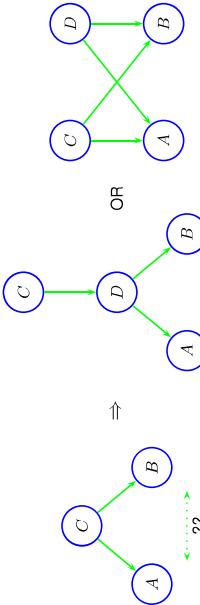
MI Autumn 2019 Specifying the structure of a Bayesian network 26

Why mediating variables?

Why do we introduce **mediating variables**:

- Necessary to catch the correct conditional independences.
- Can ease the specification of the probabilities in the model.

For example: If you find that there is a dependence between two variables *A* and *B*, but cannot determine a causal relation: Try with a **mediating variable**!



MI Autumn 2019 Specifying the structure of a Bayesian network 26

A simplified poker game

The game consists of:

- Two players.
- Three cards to each player.
- Two rounds of changing cards (max two cards in the second round)

What kind of hand does my opponent have?

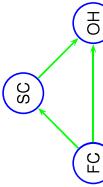
Hypothesis variable:

$$OH - \{no, 1a, 2v, fl, st, 3v, sf\}$$

Information variables:

$$FC - \{0, 1, 2, 3\}$$

and $SC - \{0, 1, 2\}$



But how do we find:

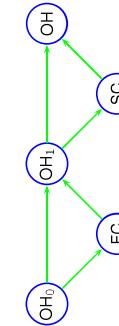
$$P(FC), P(SC|FC), \text{ and } P(OH|SC, FC)??$$

MI Autumn 2019 Specifying the structure of a Bayesian network 27

Naive Bayes Model

Introduce mediating variables:

- The opponent's initial hand, OH_0 .
- The opponent's hand after the first change of cards, OH_1 .



Note

The states of OH_0 and OH_1 are different from OH .



MI Autumn 2019 Specifying the structure of a Bayesian network 28

Specifying the structure of a Bayesian network 29

A simplified poker game

The game consists of:

- Two players.
- Three cards to each player.
- Two rounds of changing cards (max two cards in the second round)

What kind of hand does my opponent have?

Example: Spam filter

- A single query variable: *Spam*
- Many observable features (e.g. words appearing in the body of the message): *abacus*, ..., *informatics*, *pills*, ..., *watch*, ..., *zytogenic*

Network Structure:



Specifying the structure of a Bayesian network 29



We want the posterior probability of the hypothesis variable **Hyp** given the observations $\{\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n\}$:

$$P(\text{Hyp} | \text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n) = \frac{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n | \text{Hyp})P(\text{Hyp})}{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n)}$$

Note: The model assumes that the information variables are independent given the hypothesis variable.

M1 Autumn 2019

Specifying the structure of a Bayesian network

We want the posterior probability of the hypothesis variable **Hyp** given the observations $\{\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n\}$:

$$P(\text{Hyp} | \text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n) = \frac{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n | \text{Hyp})P(\text{Hyp})}{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n)}$$

$$= \mu \cdot P(\text{Inf}_1 = e_1 | \text{Hyp}) \cdot \dots \cdot P(\text{Inf}_n = e_n | \text{Hyp})P(\text{Hyp})$$

Note: The model assumes that the information variables are independent given the hypothesis variable.

Tentative course overview

Machine Intelligence

Lecture 6: Inference in Bayesian networks

Thomas Dyrre Nielsen

Aalborg University

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Inference in Bayesian networks
- Machine learning
- Planning
- Multiagent systems

M1 Autumn 2019

1

Inference Problem:
Given: a Bayesian network
Given: an assignment of values to some of the variables in the network: $E_i = e_i$ ($i = 1, \dots, l$)

30

Exact Inference

Posterior Marginals

- Inference Problem:
- Given: a Bayesian network
 - Given: an assignment of values to some of the variables in the network: $E_i = e_i$ ($i = 1, \dots, l$)

- "Instantiation of the nodes **E**"
- "Evidence **E** = **e** entered"
- "Findings entered"
- ...
- Want: for variables $A \notin E$ the *posterior marginal* $P(A | E = e)$.

M1 Autumn 2019

3

Exact Inference

3

According to the definition of conditional probability:

$$P(A \mid \mathbf{E} = \mathbf{e}) = \frac{P(A, \mathbf{E} = \mathbf{e})}{P(\mathbf{E} = \mathbf{e})}$$

Let A be the variable of interest, \mathbf{E} the evidence variables, and $\mathbf{Y} = Y_1, \dots, Y_l$ the remaining variables in the network not belonging to $A \cup \mathbf{E}$. Then

$$P(A = a, \mathbf{E} = \mathbf{e}) = \sum_{y_1 \in D_{Y_1}} \dots \sum_{y_l \in D_{Y_l}} P(A = a, \mathbf{E} = \mathbf{e}, Y_1 = y_1, \dots, Y_l = y_l).$$

Together with

$$P(\mathbf{E} = \mathbf{e}) = \sum_{a \in D_A} P(A = a, \mathbf{E} = \mathbf{e})$$

this gives the desired posterior distribution.

Note:

- For each \mathbf{y} the probability $P(A = a, \mathbf{E} = \mathbf{e}, \mathbf{Y} = \mathbf{y})$ can be computed from the network (in time linear in the number of random variables).
- There number of configurations over \mathbf{Y} is exponential in l .



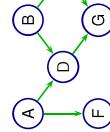
We can if we have access to $P(A, B, C, D, E, F, G, H)$:

$$P(A, B, C, D, E, F, G, H) = P(A)P(B)P(C)P(D|A, B) \dots P(H|E)$$



We can if we have access to $P(A, B, C, D, E, F, G, H)$:

$$P(A, B, C, D, E, F, G, H) = P(A)P(B)P(C)P(D|A, B) \dots P(H|E)$$



We can if we have access to:

$$P(B, \mathbf{a}, \mathbf{f}, \mathbf{g}, \mathbf{h}) = \sum_{C, D, E} P(\mathbf{a}, B, C, D, E, F, \mathbf{g}, \mathbf{h})$$

Inserting evidence we get:

$$P(\mathbf{a}, \mathbf{f}, \mathbf{g}, \mathbf{h}) = \sum_B P(\mathbf{B}, \mathbf{a}, \mathbf{f}, \mathbf{g}, \mathbf{h})$$

and

$$P(\mathbf{Sex} \mid \text{Hair length} = \text{long}) = ?$$

$$P(\text{Sex} \mid \text{Hair length} = \text{short}, \text{Stature} \leq 1.68) = ?$$



We can if we have access to:

$$P(\text{Sex} \mid \text{Hair length} = \text{long}) = ?$$

Posterior probability inference: Given the value of some observed variables (the evidence) compute the conditional distribution of some other variable:

$$P(\text{Sex} \mid \text{Hair length} = \text{short}, \text{Stature} \leq 1.68) = ?$$

Naive Solution Step 1

Naive Solution Step 1

Naive Solution Step 1

		$P(\text{Sex})$		$P(\text{Hair length} \mid \text{Sex})$		$P(\text{Stature} \mid \text{Sex})$	
		Sex		Hair length	Sex	Stature	Sex
male	0.49	long	male	≤ 1.68	male	0.47	
female	0.51	short	female	0.6	female	0.47	
			0.95	0.4	≥ 1.68	0.92	0.53

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

		$P(\text{Sex})$				$P(\text{Stature} \mid \text{Sex})$			
		$P(\text{Hair length} \mid \text{Sex})$		$P(\text{Sex})$		$P(\text{Stature} \mid \text{Sex})$		$P(\text{Sex})$	
Sex	male	0.49	long	0.05	male	0.6	male	0.47	
	female	0.51	short	0.95	female	0.4	female	0.53	
						</th			

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

$P(\text{Sex} \text{Stature})$		$P(\text{Stature} \text{Sex})$	
		male	female
		male	female
Hair length	male	0.05	0.6
	female	0.95	0.4
long			
short			

Query: $P(\text{Stature} \mid \text{Hair length} = \text{long}) = ?$

Naive Solution Step 1

Naive Solution Step 1

Naive Solution Step 1

$P(\text{Sex})$		$P(\text{Hair length} \mid \text{Sex})$		$P(\text{Stature} \mid \text{Sex})$		$P(\text{Hair length} \mid \text{Sex})$		$P(\text{Stature} \mid \text{Sex})$	
Sex		Hair length	male female	Stature	male female	Hair length	male female	Stature	male female
male	0.49	long	0.05	male	0.98	long	0.49	male	0.98
female	0.51	long	0.6	female	0.47	long	0.51	female	0.47
				$P(\text{Sex})$		$P(\text{Stature} \mid \text{Sex})$		$P(\text{Stature} \mid \text{Sex})$	
				Sex		Sex		Sex	
				male		male		male	
				female		female		female	

SOMMAIRE

8..... D / Cn+4..... | H -> H | I -> I | J -> J | K -> K | L -> L | M -> M | N -> N | O -> O

Step 1: Construct joint distribution

Step 1: Construct joint distribution

Step 1: Construct joint distribution

	male		female		male		female	
	Hair length	Hair length	Hair length	Hair length	Stature	Hair length	Stature	Hair length
Stature	long	short	long	short	≤ 1.68	0.00196	0.00196	0.0324
					> 1.68		> 1.68	

Native Solution Stan 1

Naive Solution Step 2

$P(\text{Sex})$		$P(\text{Hair length} \mid \text{Sex})$				$P(\text{Stature} \mid \text{Sex})$			
		Sex		Sex		Sex		Sex	
Sex	Sex	Hair length	male	Hair length	male	Stature	male	Stature	male
male	0.49	long	0.05	long	0.6	≤ 1.68	0.08	≤ 1.68	0.47
female	0.51	long	0.95	short	0.4	> 1.68	0.92	> 1.68	0.53

Answers: $D(\text{Statutes})$ | $H(\text{air length})$ | $I(\text{age})$ | ?

THE JOURNAL OF CLIMATE

Joint distribution $P(Sex, Hair\ length, Stature)$

		Sex	
		male	female
Stature	long	short	long
	short	long	short

Forest plot showing the effect of Hair length on Statute (cm) for males and females across different height distributions. The y-axis represents the effect size (cm) and the x-axis represents the height distribution.

Height distribution	$P(\text{Sex}, \text{Hair length})$		$P(\text{Sex}, \text{Hair length}, \text{Statute})$	
	male	female	male	female
≤ 1.68	0.00136	0.03124	0.4382	0.04956
> 1.68	0.02254	0.04206	0.16218	0.07082

Note: The table on the right shows neither a joint nor a conditional distribution!

P: Sex, Hair length=long, Signature			
		Sex	
		male	female
		≤ 1.68	0.1482
		> 1.68	0.16218

Exhibit Finance
M1 Autumn 2019 6 M1 Autumn 2019

MI Autumn 2019
Exact inference 9

Step 3 Marginalize (sum out Sex variable):

$P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$		\sum_{Sex}		$P(\text{Hair length}=\text{long}, \text{Stature})$
		male	female	
Stature	≤ 1.68	0.00196	0.14592	0.14578
	> 1.68	0.02254	0.16218	0.18472

Step 3 Marginalize (sum out Sex variable):

$P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$		\sum_{Sex}		$P(\text{Hair length}=\text{long}, \text{Stature})$
		male	female	
Stature	≤ 1.68	0.00196	0.14592	0.14578
	> 1.68	0.02254	0.16218	0.18472

Step 4 Normalize

$P(\text{Hair length}=\text{long}, \text{Stature})$		\sum_{Sex}		$P(\text{Hair length}=\text{long})$
		male	female	
Stature	≤ 1.68	0.14578	0.14578	0.14578
	> 1.68	0.18472	0.18472	0.18472

Step 4 Normalize

$P(\text{Hair length}=\text{long}, \text{Stature})$		\sum_{Sex}		$P(\text{Stature} \text{Hair length}=\text{long})$
		male	female	
Stature	≤ 1.68	0.14578	0.14578	0.441
	> 1.68	0.18472	0.18472	0.559

Exact inference

M1 Autumn 2019

10

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} | \text{Sex})P(\text{Stature} | \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Exact inference

M1 Autumn 2019

10

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} | \text{Sex})P(\text{Stature} | \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Exact inference

M1 Autumn 2019

11

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} | \text{Sex})P(\text{Stature} | \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Exact inference

M1 Autumn 2019

11

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} | \text{Sex})P(\text{Stature} | \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Exact inference

M1 Autumn 2019

11

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} | \text{Sex})P(\text{Stature} | \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Exact inference

M1 Autumn 2019

11

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} | \text{Sex})P(\text{Stature} | \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Exact inference

M1 Autumn 2019

11

Naive Solution: Summary

Construct Joint: $P(\text{Sex}, \text{Hair length}, \text{Stature}) = P(\text{Sex})P(\text{Hair length} | \text{Sex})P(\text{Stature} | \text{Sex})$

Insert Evidence: $P(\text{Sex}, \text{Hair length}=\text{long}, \text{Stature})$

Marginalize: $P(\text{Hair length}=\text{long}, \text{Stature}) = \sum_{\text{sex} \in \{\text{male, female}\}} P(\text{Sex}=\text{sex}, \text{Hair length}=\text{long}, \text{Stature})$

We can use

- the form of the joint distribution P
- the law of distributivity

to make the computation of the sum more efficient.

$$P(\text{Hair length}=\text{long}, \text{Stature}) \leq 1.68 + P(\text{Hair length}=\text{long}, \text{Stature} > 1.68)$$

Complexity

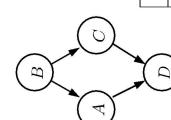
Complexity dominated by initial table $P(\text{Sex}, \text{Hair length}, \text{Stature})$ (size 2^3).

For model with n binary variables:

$$O(2^n)$$

MI Autumn 2019 Exact Inference 11

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

MI Autumn 2019 Exact Inference 13

Example



A	C	t	f
B	t	.9	.1
C	t	.7	.3
\vdots	\vdots	\vdots	\vdots
B	t	.3	.7
C	t	.2	.8
\vdots	\vdots	\vdots	\vdots
B	t	.5	.5
C	t	.4	.6

$$P(A, D = f) =$$

$$\sum_{b \in \{t, f\}} \sum_{c \in \{t, f\}} P(B = b, A, C = c, D = f) =$$

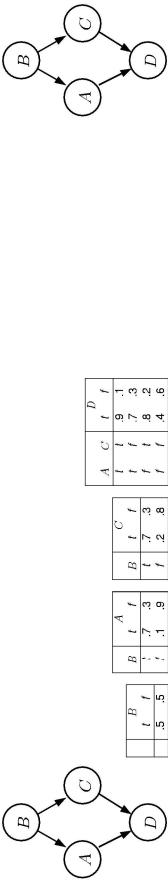
MI Autumn 2019 Exact Inference 13

Example



Example continued

Example continued



	A	C	t	f
B	t	t	$.9$	$.1$
t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$
$.5$	$.1$	$.9$	$.4$	$.6$

$$\sum_b P(B = b)P(A \mid B = b) \sum_c P(C = c \mid B = b)P(D = f \mid A, C = c) =$$

$$\sum_b P(B = b)P(A \mid B = b)F_1(B = b, A) = F_2(A)$$

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

$$\sum_b P(B = b)P(A \mid B = b) \sum_c P(C = c \mid B = b)P(D = f \mid A, C = c) =$$

$$\sum_b P(B = b)P(A \mid B = b) \sum_c P(C = c \mid B = b)F_1(B = b, A) = F_2(A)$$

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

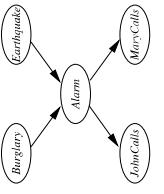
	B	t	A	t	C	t	f
B	t	t	t	t	t	$.9$	$.1$
t	t	t	t	t	t	$.7$	$.3$
$.5$	$.7$	$.3$	$.2$	$.8$	$.2$	$.8$	$.2$
$.5$	$.1$	$.9$	$.4$	$.6$			

</th

Alarm Example

Alarm Example

Example



Bad ordering for computing $P(MC, B = t)$:

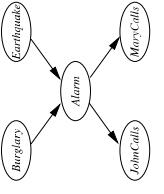
$$\sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(EQ = eq | A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) = \\ \sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) = \\ \sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) =$$

MI Autumn 2019

16

Exact Inference

Example



Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) = \\ \sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) = \\ \sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) = \\ \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq | f_2(eq, MC)) = \\ \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq | f_2(eq, MC)) =$$

MI Autumn 2019

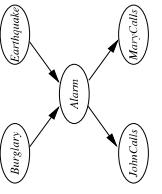
16

Exact Inference

Alarm Example

Alarm Example

Example



Bad ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) = \\ \sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) = \\ \sum_{eq \in \{t, f\}} \sum_{f \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) = \\ \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq | f_2(eq, MC)) = \\ P(B = t)P_3(MC)$$

MI Autumn 2019

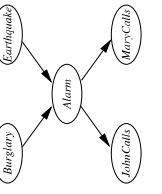
16

Exact Inference

Alarm Example

Alarm Example continued

Example



Good ordering for computing $P(MC, B = t)$:

$$\sum_{eq \in \{t, f\}} \sum_{eq \in \{t, f\}} P(B = t)P(A = a | B = t, EQ = eq)P(JC = jc | A = a)P(MC | A = a) = \\ \sum_{eq \in \{t, f\}} \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) = \\ \sum_{eq \in \{t, f\}} \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq | f_1(eq, jc, MC)) = \\ \sum_{eq \in \{t, f\}} P(B = t)P(EQ = eq | f_2(eq, MC)) = \\ P(B = t)P_3(MC)$$

Largest factor (f_1) is function of 3 variables.

MI Autumn 2019

16

Exact Inference

17



We want the posterior probability of the hypothesis variable Hyp given the observations $\{\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n\}$:

$$P(\text{Hyp}|\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n) = \frac{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n|\text{Hyp})P(\text{Hyp})}{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n)}$$

Note: The model assumes that the information variables are independent given the hypothesis variable.

MI Autumn 2019

Exact Inference

20

We want the posterior probability of the hypothesis variable Hyp given the observations $\{\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n\}$:

$$P(\text{Hyp}|\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n) = \frac{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n|\text{Hyp})P(\text{Hyp})}{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n)} = \frac{P(\text{Inf}_1 = e_1, \dots, \text{Inf}_n = e_n|\text{Hyp})P(\text{Hyp})}{P(\text{Inf}_1 = e_1|\text{Hyp}) \cdot \dots \cdot P(\text{Inf}_n = e_n|\text{Hyp})P(\text{Hyp})}$$

Note: The model assumes that the information variables are independent given the hypothesis variable.

MI Autumn 2019

Exact Inference

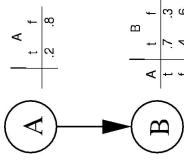
20

Sampling Inference

Sample Generator

Observation: can use Bayesian network as random generator that produces states $X = x$ according to distribution P defined by the network.

Example:



- Generate random numbers r_1, r_2 uniformly from [0,1].
- Set $A = t$ if $r_1 \leq .2$ and $A = f$ else
- Depending on the value of A and r_2 set B to t or f .

Random generation of one state: linear in size of network.

MI Autumn 2019

Approximate Inference

21

MI Autumn 2019

Approximate Inference

21

Sampling Inference

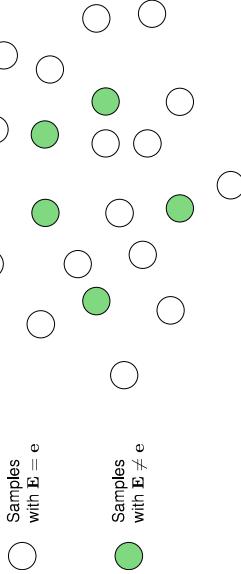
Approximate Inference from Samples

To compute an approximation of $P(E = e)$ (E a subset of the variables in the Bayesian network):

- generate a (large) number of random states
- count the frequency of states in which $E = e$.

Then

$$P(|s - p| > \epsilon) \leq 2e^{-2n\epsilon^2}$$



MI Autumn 2019

Approximate Inference

22

MI Autumn 2019

Approximate Inference

22

Accuracy

Hoeffding Bound

- p : true probability $P(E = e)$
- s : estimate for p from sample of size n
- ϵ : an error bound > 0 .

Then

MI Autumn 2019

Approximate Inference

23

MI Autumn 2019

Approximate Inference

23

Likelihood weighting I

Likelihood weighting I

First idea (not to be followed)

- Fix evidence variables to their observed states.
- Sample from non-evidence variables.

Problem: This gives a sampling distribution

$$\prod_{X \in \mathbf{X} \setminus \mathbf{E}} P(X \mid \text{pa}(X) \setminus \mathbf{E}, \text{pa}(X) \cap \mathbf{E})$$

somewhere between $P(\mathbf{X})$ and $P(\mathbf{X} \mid \mathbf{e})$.

MI Autumn 2019

Approximate Inference

26

Approximate Inference

26

Likelihood weighting I

First idea (not to be followed)

- Fix evidence variables to their observed states.
- Sample from non-evidence variables.

Problem: This gives a sampling distribution

$$\underbrace{\prod_{X \in \mathbf{X} \setminus \mathbf{E}} P(X \mid \text{pa}(X) \setminus \mathbf{E}, \text{pa}(X) \cap \mathbf{E})}_{\text{Part 1}} \cdot \underbrace{\prod_{E \in \mathbf{E}} P(E \mid \text{pa}(E) \setminus \mathbf{E}, \text{pa}(E) \cap \mathbf{E})}_{\text{Part 2}}$$

somewhere between $P(\mathbf{X})$ and $P(\mathbf{X} \mid \mathbf{e})$.

Likelihood weighting

We would like to sample from

$$P(\mathbf{X}, \mathbf{e}) = \underbrace{\prod_{X \in \mathbf{X} \setminus \mathbf{E}} P(X \mid \text{pa}(X) \setminus \mathbf{E}, \text{pa}(X) \cap \mathbf{E})}_{\text{Part 1}} \cdot \underbrace{\prod_{E \in \mathbf{E}} P(E \mid \text{pa}(E) \setminus \mathbf{E}, \text{pa}(E) \cap \mathbf{E})}_{\text{Part 2}}$$

So instead weigh each generated sample with a weight corresponding to Part 2.

MI Autumn 2019

Approximate Inference

26

Approximate Inference

26

Likelihood weighting II

Estimate $P(\mathbf{X} = \mathbf{e} \mid \mathbf{e})$ as

$$\hat{P}(\mathbf{X} = \mathbf{e} \mid \mathbf{e}) = \frac{\sum_{\text{sample: } X=x} w(\text{sample})}{\sum_{\text{sample}} w(\text{sample})},$$

where

$$w(\text{sample}) = \prod_{E \in \mathbf{E}} P(E = e \mid \text{pa}(E) = \pi)$$

(Part 2 on the previous slide)

and π is the values of $\text{pa}(E)$ under sample and e .

Importance sampling

Likelihood weighting is an instance of importance sampling, where

- samples are weighted and can come from (almost) any proposal distribution.

MI Autumn 2019

Approximate Inference

27

Approximate Inference

27

Machine Intelligence

Lecture 7: Learning - Introduction and decision trees

Thomas Dyre Nielsen

Aalborg University

1

Tentative course overview

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty
- Bayesian networks
- Inference in Bayesian networks
- **Machine Learning**
- Planning
- Reinforcement learning
- Multi-agent systems

Supervised Learning

MI Autumn 2019 2 Supervised Learning

So far

The general pattern so far:

Problem/Domain description	State Space Problem	Variables, Constraints	Probabilistic Model
Inference Algorithms	Search (A^*, \dots)	Arc Consistency, Variable Elimination	
Solutions		Goal states, plans, diagnoses, predictions,...	

The general pattern so far:

Problem/Domain description	State Space Problem	Variables, Constraints	Probabilistic Model
Inference Algorithms	Search (A^*, \dots)	Arc Consistency, Variable Elimination	Variable Elimination
Solutions		Goal states, plans, diagnoses, predictions,...	

- Problem/Domain description and algorithms designed by human
- Agent will always act the same in the same situation

- Problem/Domain description and algorithms designed by human
 - Agent will always act the same in the same situation
- Objective of (machine) **learning**:
- Agent can learn by experience: improve performance over time
 - Agent (program) can be automatically constructed from examples (rather than designed by expert)

Tasks and Model Types

The models constructed by machine learning algorithms are used for several kinds of tasks:

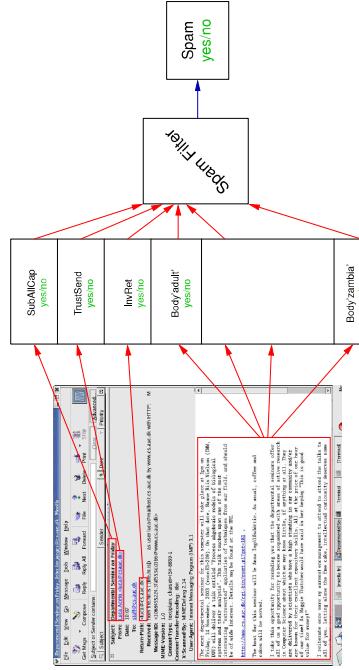
Predictive tasks/models

- Task: predict some (unobserved) **target** or **class** variable based on observed values of **(predictor) attributes**
- **Regression**, if target is continuous
- **Classification**, if target is discrete
- Examples: Spam filtering, Character recognition, ...
- Model types e.g.: Decision trees, k nearest neighbors, Neural networks, Naive Bayes, ...

Descriptive tasks/models

- Task: **Clustering**: identify coherent subgroups in data
- Examples: Recommender systems,
- Model types e.g.: k -means, hierarchical clustering, Self-organizing maps, probabilistic clustering, ...

Example: Spam filter



MI Autumn 2019 4 Supervised Learning

Supervised Learning

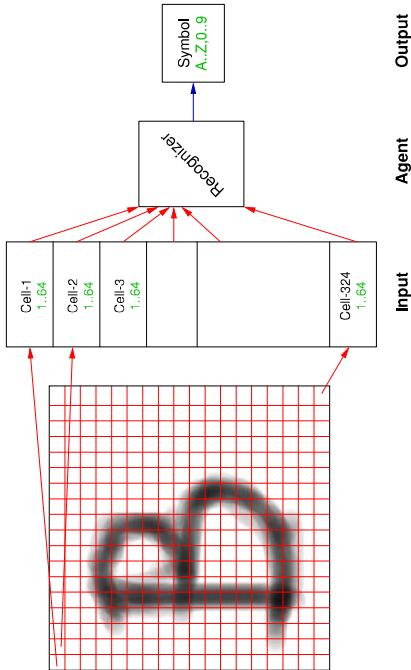
MI Autumn 2019 5 Supervised Learning

Output

Agent

Example: Character Recognition

Experience, Data



MI Autumn 2019 Supervised Learning 6

Data

Continuous Attributes

We assume that data (experience) consists of an **attribute-value table**:

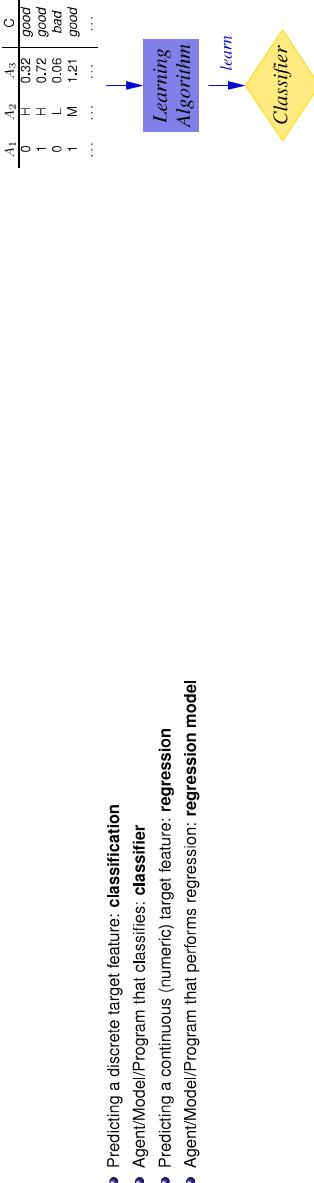
Input Features		Target Feature (Class variable)	
(Attributes, Predictor Variables)	SubAllCap	Spam	
n	n	n	
n	n	n	
n	y	n	
n	n	n	
n	n	n	
...

- Columns correspond to **attributes** given by a **name** and a **state space** (attributes are basically the same as (chance) variables).
 - Rows correspond to **examples** (also called **cases** or **instances**): observations of joint occurrences of values of the attributes.
 - In prediction problems, there is a distinguished **target attribute**. When the target attribute is discrete, it is usually called the **class variable**. The attributes used for prediction are then called **predictor attributes**.
- In table above: **Spam** is class variable for the prediction problem: predict whether a mail is spam, given characteristics of the mail.

MI Autumn 2019 Supervised Learning 7

Classification, Regression

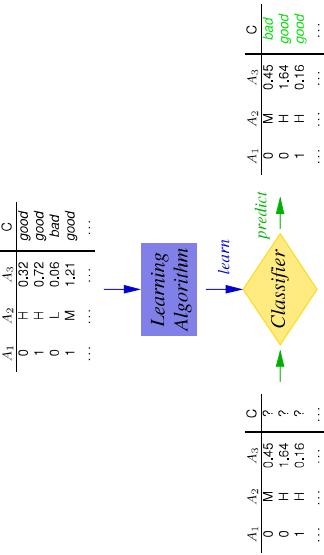
Learning and Using a Classifier



MI Autumn 2019 Supervised Learning 10

Supervised Learning

11



- Key ingredients of a learning method are:
- a **Hypothesis Space**: set of all possible classifiers that could be learned based on a given **Representation**
- an **Evaluation Measure** that is used to decide how good a candidate hypothesis is
- a **Search or Optimization** method used to find a hypothesis that scores high according to the evaluation measure.

- Key ingredients of a learning method are:
- a **Hypothesis Space**: set of all possible classifiers that could be learned based on a given **Representation**
- an **Evaluation Measure** that is used to decide how good a candidate hypothesis is
- a **Search or Optimization** method used to find a hypothesis that scores high according to the evaluation measure.

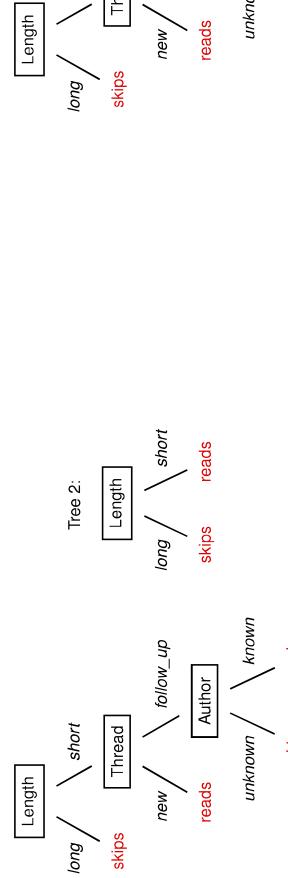
Example: Training Data

User preference data:

Example	Author	Thread	Length	WhereRead	UserAction
e1	known	new	long	home	skips
e2	unknown	new	short	work	reads
e3	unknown	follow up	long	home	skips
e4	known	new	long	home	reads
e5	known	follow up	long	work	skips
e6	unknown	follow up	short	work	reads
e7	unknown	new	short	work	reads
e8	known	follow up	long	home	skips
e9	known	new	long	work	skips
e10	unknown	follow up	short	home	skips
e11	known	new	long	work	skips
e12	known	follow up	short	home	reads
e13	known	new	short	work	reads
e14	known	new	short	home	reads
e15	known	follow up	short	work	reads
e16	unknown	new	short	home	reads
e17	unknown	new	short	work	reads
e18	unknown	new	short	work	reads
e19	unknown	follow up	long	work	?
e20	unknown	follow up	long	home	?
e21	unknown	follow up	short	home	?

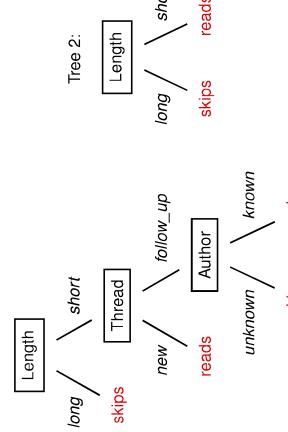
Example: Decision Trees

Tree 1:



Example: Decision Trees

Tree 2:



Example: Decision Trees

Tree 1 is equivalent to the following logic program:

```
skips ← long
reads ← short ∧ new
```

```
reads ← short ∧ follow_up ∧ known
```

```
skips ← short ∧ follow_up ∧ unknown
```

Example: Classifications

Learning Decision Trees

Example	Author	Thread	Length	WhereRead	UserAction	Tree 1	Tree 2
e ₁	known	new	long	home	skips	skips	Top-down construction:
e ₂	unknown	new	short	work	reads	reads	
e ₃	known	follow Up	long	home	skips	skips	
e ₄	known	new	short	home	reads	skips	
e ₅	known	follow Up	short	work	skips	reads	
e ₆	known	follow Up	short	work	skips	skips	
e ₇	unknown	new	short	work	reads	reads	
e ₈	known	follow Up	short	work	skips	skips	
e ₉	known	new	long	work	skips	skips	
e ₁₀	unknown	follow Up	short	home	reads	reads	
e ₁₁	known	new	long	work	skips	skips	
e ₁₂	known	follow Up	short	home	reads	skips	
e ₁₃	known	new	short	work	skips	reads	
e ₁₄	known	new	short	work	reads	reads	
e ₁₅	known	follow Up	short	home	reads	reads	
e ₁₆	known	new	short	home	reads	reads	
e ₁₇	unknown	new	short	work	reads	reads	
e ₁₈	unknown	new	long	work	?	?	
e ₁₉	unknown	follow Up	long	home	?	skips	
e ₂₀	unknown	follow Up	short	home	?	skips	
e ₂₁	unknown	new	short	work	?	skips	

Decision Trees

MI Autumn 2019

15

Choosing X_i

Key question: which X_i to choose in line 4?

Approach: choose the feature that would provide the best classifier if construction would terminate with that feature.



Showing

- Number of examples with class labels skips, reads belonging to different sub-trees
- Green: predicted class label (possibly a tie between two labels)

Decision Trees

MI Autumn 2019

17

Class Purity

Principle: Prefer features that split the examples into *class pure* subsets:

pure:	nearly pure:	impure:
skips: 7, reads: 0	skips: 6, reads: 1	skips: 5, reads: 5
skips: 0, reads: 5	skips: 2, reads: 15	skips: 7, reads: 7
skips: 11, reads: 0	skips: 11, reads: 1	skips: 13, reads: 12

Normalized to probabilities:

pure:	nearly pure:	impure:
skips: 1, reads: 0	skips: 0.85, reads: 0.15	skips: 0.54, reads: 0.46
skips: 0, reads: 1	skips: 0.12, reads: 0.88	skips: 0.5, reads: 0.5
skips: 1, reads: 0	skips: 0.91, reads: 0.09	skips: 0.52, reads: 0.48

Decision Trees

MI Autumn 2019

18

Entropy

Entropy Measure

For a probability distribution $(p, 1-p)$ of a two-valued class label, define

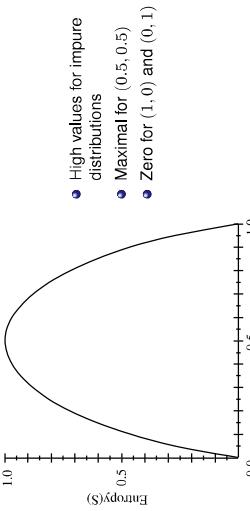
$$h(p, 1-p) = -p \cdot \log(p) - (1-p) \cdot \log(1-p)$$

Decision Trees

MI Autumn 2019

19

Generalization to larger domains



Again:

- Maximal for $\mathbf{p} = (1/n, \dots, 1/n)$
- Zero for $\mathbf{p} = (1, 0, \dots, 0)$, $\mathbf{p} = (0, 1, 0, \dots, 0)$, ..., $\mathbf{p} = (0, \dots, 0, 1)$

Examples:

- Entropy(0.5, 0.5) = $-0.5 \cdot \log_2(0.5) - 0.5 \cdot \log_2(0.5) = 0.5 \cdot 1 + 0.5 \cdot 1 = 1$
- Entropy(0.35, 0.65) = $-0.35 \cdot \log_2(0.35) - 0.65 \cdot \log_2(0.65) = 0.93$
- Entropy(0, 1) = $-0 \cdot \log_2(0) - 1 \cdot \log_2(1) = -0 - 0 = 0$
- Entropy(1, 0) = $-1 \cdot \log_2(1) - 0 \cdot \log_2(0) = -0 - 0 = 0$

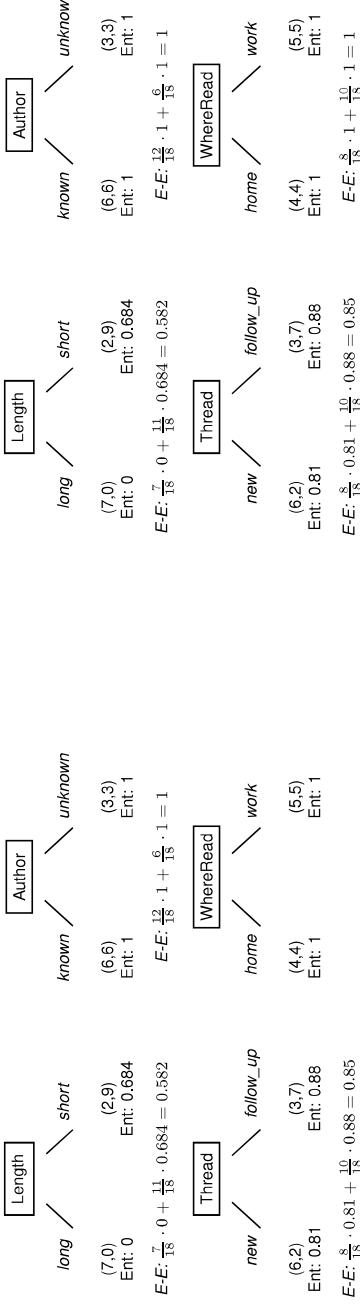
Decision Trees

MI Autumn 2019

20

Information gain

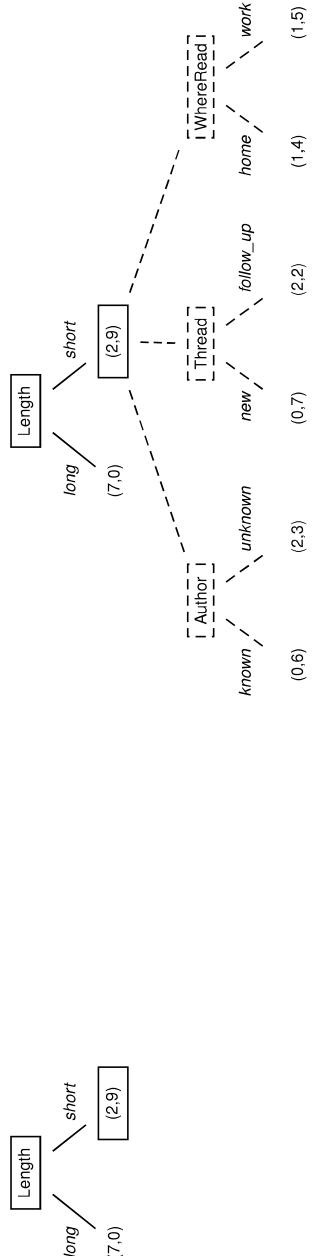
Information gain



Mi Autumn 2019 Decision Trees 23

Constructing the decision tree

Constructing the decision tree



Mi Autumn 2019 Decision Trees 24

Constructing the decision tree

- The **information gain** measure favors attributes with many values:
 - For example, the attribute **Date** (with the possible dates as states) will have a very high **information gain** but is unable to generalize!
- One approach for avoiding this problem is to select **attributes** based on **GainRatio**:
- $$\text{GainRatio}(S, A) = \frac{\text{Gain}(S, A)}{\text{SplitInformation}(S, A)}$$
- $$\text{SplitInformation}(S, A) = - \sum_{i=1}^c \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|},$$
- where S_i is the subset of examples produced by splitting on the i th value of A .
- Note that **SplitInformation** is the entropy of S w.r.t. the values of A .

Mi Autumn 2019 Decision Trees 25

- We require that the attributes being tested are discrete valued. So in order to test a continuous valued attribute we need to 'discretize' it.
- Suppose that the training examples are associated with the attribute **Temperature**:
- | | | | | | | |
|-----------------------|-----|----|----|----|----|-----|
| Temperature: | 40 | 48 | 60 | 72 | 80 | 90 |
| Rain tomorrow: | yes | no | no | no | no | yes |

Mi Autumn 2019 Decision Trees 26

Continuous/many-valued attributes |

Overfitting

We require that the attributes being tested are discrete valued. So in order to test a continuous valued attribute we need to "discretize" it.

Suppose that the training examples are associated with the attribute **Temperature**:

Temperature:	40	48	60	72	80	90
Rain tomorrow:	yes	no	no	no	yes	yes

Create a new boolean valued attribute by first testing the two candidate thresholds:

- $(48+60)/2$

- $(80+90)/2$

Next, pick the one with highest information gain (i.e., $\text{Temperature}_{>54}$)

MI Autumn 2019 | Decision Trees | 26

MI Autumn 2019 | Decision Trees | 27

Tentative course overview

Machine Intelligence Lecture 8: Learning II

Thomas Dyrhe Nielsen
Aalborg University

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Inference in Bayesian networks
- **Machine learning**
- Planning
- Reinforcement learning
- Multi-agent systems

MI E19 | 1

MI E19 | 2

Example applications

Neural Networks



Colorization of images

MI E19 | Neural Networks | 3

Zhang, Isola, Efros. Colorful Image Colorization. In ECCV 2016.

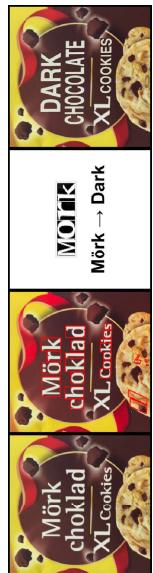
Neural Networks

MI E19 | Neural Networks | 5

Example applications

Example applications

Text-image translation



<https://research.googleblog.com/2015/07/cow-vs-train-are-squeeze-deep.html>

MI E19

Neural Networks

Regression: Example

Cereals Data

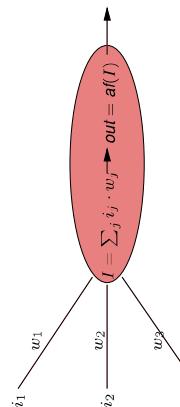
Name	Calories	Protein	Fat	Fiber	Carbo	Sugars	Rating
All-Bran	70	4	1	9	7	5	59.42
All-Bran_with_Extra_Fiber	50	4	0	14	8	0	93.70
Almond_Delight	110	2	2	1	14	8	34.38
Apple_Cinnamon_Cheerios	110	2	2	15	10.5	10	29.50
Apple_Jacks	110	2	0	1	11	14	33.17
Basic_4	130	3	2	2	18	8	37.03
Bran_Cheex	90	2	1	4	15	6	49.12
Bran_Flakes	90	3	0	5	13	5	53.31
Cap_n_Crunch	120	1	2	0	12	12	18.04
Cheerios	110	6	2	2	17	1	50.76
...

Problem: predict nutritional rating of cereal.

MI E19

Neural Networks

Single Neuron



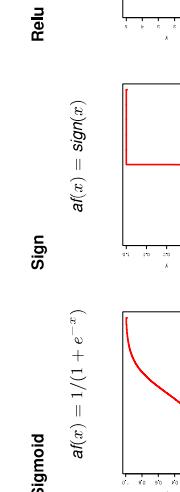
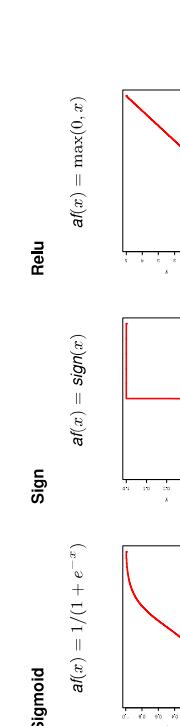
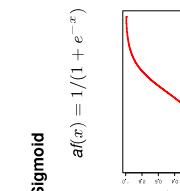
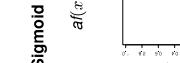
Two step computation:

- Combine inputs as **weighted sum**
- Compute output by **activation function** of combined inputs

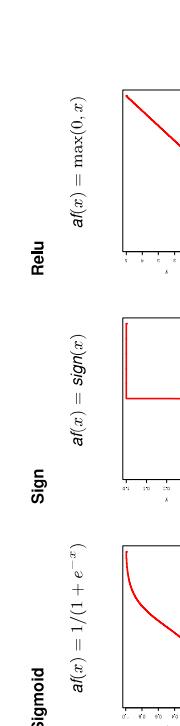
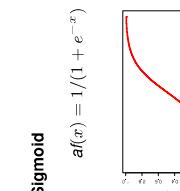
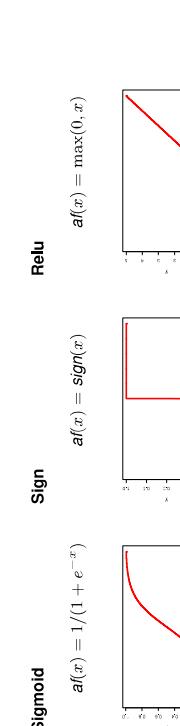
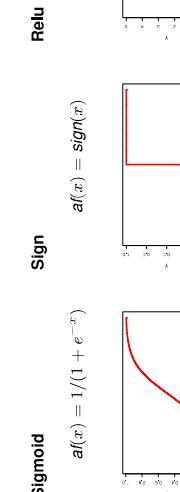
MI E19

Neural Networks

Activation Functions



The most common activation functions are:



- If activation function is sigmoid, i.e. $out = \sigma(\sum_j i_j \cdot w_j)$, we also talk of **squashed linear function**.
- For the output neuron also the **identity function** is used: $af(x) = id(x) = x$

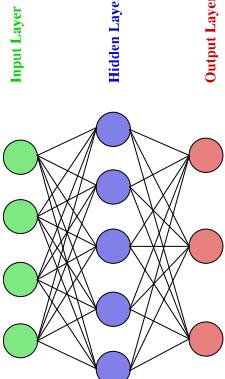
MI E19

Neural Networks

MI E19

Neural Networks

MI E19



Given

- the network structure,
 - the weights associated with links/nodes,
 - the activation function (usually the same for all hidden/output nodes),
- a neural network with n input and k output nodes defines k real-valued functions on continuous input attributes:

$$o_i(a_1, \dots, a_n) \in \mathbb{R} \quad (i = 1, \dots, k).$$

Discrete Inputs

Discrete Features in general

Neural networks can also handle discrete features as inputs or outputs:

Indicator variables

If the regression should also use discrete predictor attributes, e.g.:

Calories	Protein	Sugars	Vitamins	Manufacturer	Rating
70	105	8	135	Kellogg	59.3
110	80	23	99	Nabisco	43.6
...

replace discrete attributes with 0-1-valued indicator variables for their possible values:

Calories	Protein	Sugars	Vitamins	M_Kellogg	M_Nabisco	M_xxx	Rating
70	105	8	135	0	1	...	59.3
110	80	23	99	1	0	...	43.6
...

Discrete Features in general

Neural networks can also handle discrete features as inputs or outputs:

Indicator variables

For each value x_i of X with domain $\{x_1, \dots, x_k\}$ introduce a binary feature $X_{is_x_i}$ with values 0, 1.Encode input $X = x_i$ by inputs

$$X_{is_x_0} = 0, \dots, X_{is_x_{i-1}} = 0, X_{is_x_i} = 1, X_{is_x_{i+1}} = 0, \dots, X_{is_x_k} = 0$$

Numerical Encoding

Translate values into numbers, e.g.:

$$\text{true}, \text{false} \rightarrow 1, 0$$

$$\text{low}, \text{medium}, \text{high} \rightarrow 0, 1, 2$$

Probably not sensible:

$$\text{red}, \text{green}, \text{blue} \rightarrow 0, 1, 2$$

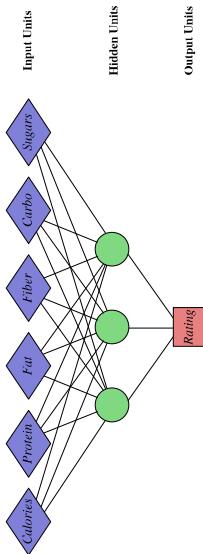
because blue is not "two times green"

Neural Networks

Neural Networks

Neural Networks for Regression

Task: predict the (health-)rating of breakfast cereals based on their contents.



10

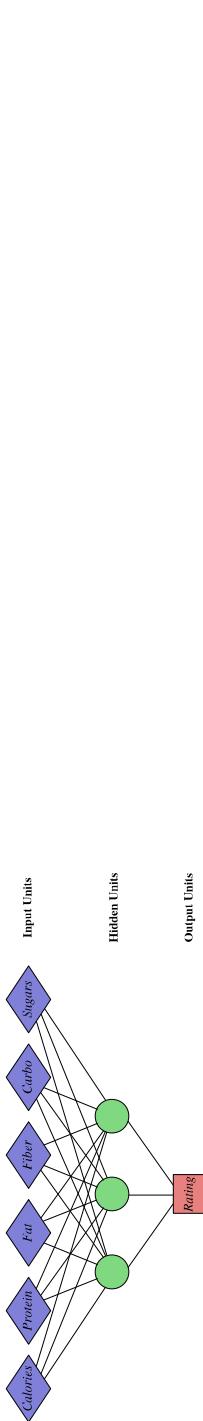
ML E19

Neural Networks

Neural Networks for Classification

Task: hand-written character recognition. Predictor attributes: (continuous) grey-scale values for 18×18 grid cells. Class label: one of $A, \dots, Z, 0, \dots, 9$.

NN regression model, all predictor attributes are continuous:



11

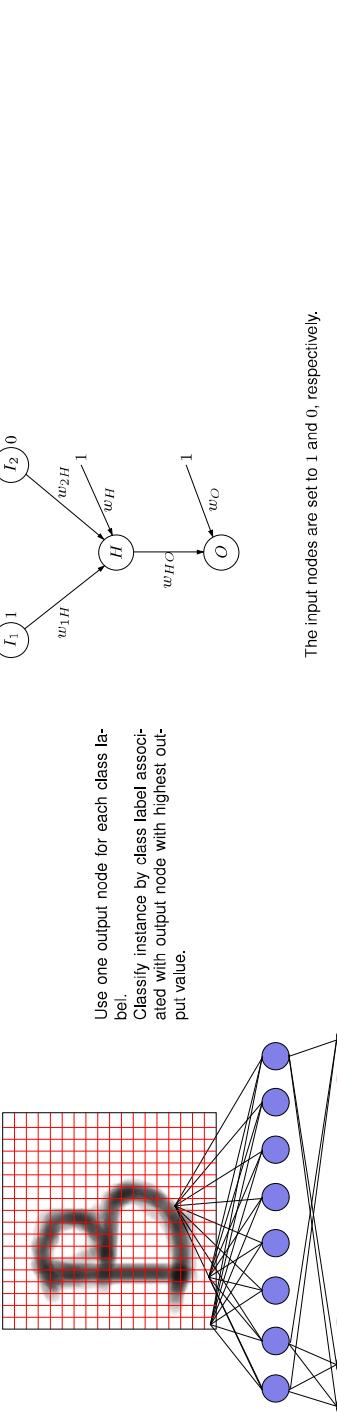
ML E19

Neural Networks

Neural Networks

Neural Networks for Classification

Task: hand-written character recognition. Predictor attributes: (continuous) grey-scale values for 18×18 grid cells. Class label: one of $A, \dots, Z, 0, \dots, 9$.



12

ML E19

Neural Networks

Neural Networks

Propagation in Neural Networks

Use one output node for each class label.
Classify instance by class label associated with output node with highest output value.



The input nodes are set to 1 and 0, respectively.

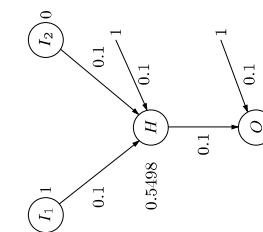
13

ML E19

Neural Networks

Propagation in Neural Networks

The output of neuron H is:
$$o_H = \sigma(1 \cdot 0.1 + 0 \cdot 1 + 1 \cdot 1) = 0.5498.$$



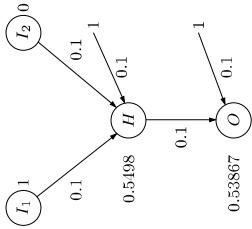
The input nodes are set to 1 and 0, respectively.

13

ML E19

Neural Networks

Propagation in Neural Networks



The input nodes are set to 1 and 0, respectively.

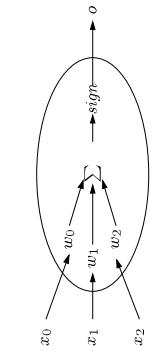
The output of neuron O is:

$$o_H = \sigma(0.5498 \cdot 1 + 1 \cdot 0.1) = 0.53867.$$

Expressive power

The perceptron

Expressive power



The decision surface of a two-input perceptron $a(x_1, x_2) = \text{sign}(x_1 \cdot w_1 + x_2 \cdot w_2 + w_0)$ is given by a straight line, separating positive and negative examples.

- No hidden layer

- One output neuron o

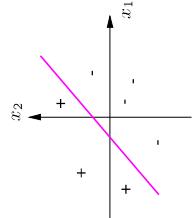
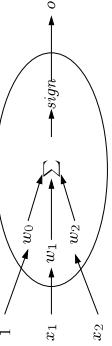
- sign activation function

Function computed:

$$O(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

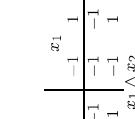
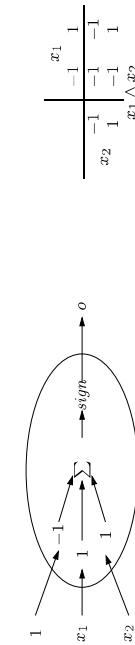
Convention: from now on assume that x_0 is an input neuron with constant input value 1. Then write $w \cdot \mathbf{x}$ for $w_0 x_0 + w_1 x_1 + \dots + w_n x_n$.

Expressive power

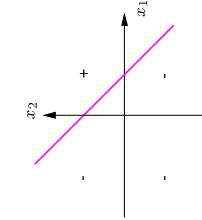
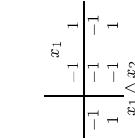
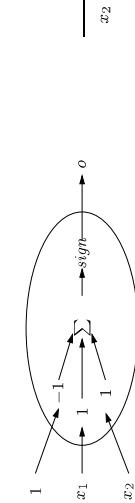


This perceptron specifies the decision surface:

$$-1 + 1 \cdot x_1 + 1 \cdot x_2 = 0$$

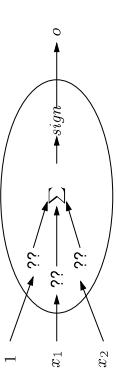


Can the perceptron represent the Boolean function?

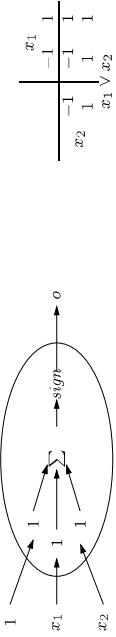
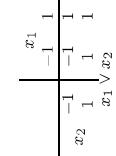


Perceptron examples

Perceptron examples

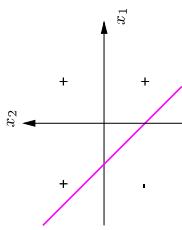


Can the perceptron represent the Boolean function?



This perceptron specifies the decision surface:

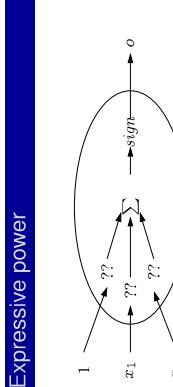
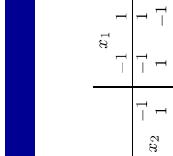
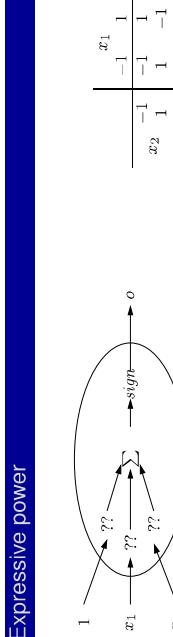
$$1 + 1 \cdot x_1 + 1 \cdot x_2 = 0$$



MI E 19 Expressive power 17

Expressive power

Expressive power



Can the perceptron represent the Boolean function?

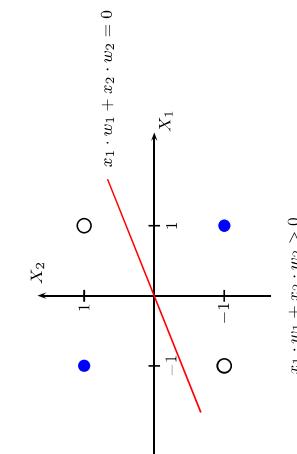
We cannot specify any values for the weights so that the perceptron can represent the "xor" function:

The examples are not linear separable!

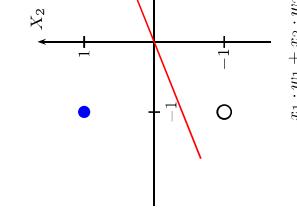
MI E 19 Expressive power 18

Linear separable

Linear separable



$$x_1 \cdot w_1 + x_2 \cdot w_2 > 0$$



$$x_1 \cdot w_1 + x_2 \cdot w_2 > 0$$

MI E 19 Expressive power 19

MI E 19 Expressive power

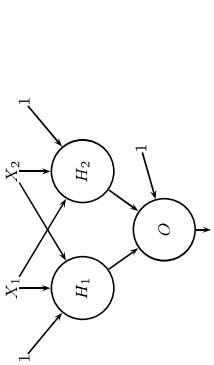
19

Only linearly separable functions can be computed by a single perceptron.

XOR again

XOR again

Can multiple neurons help?

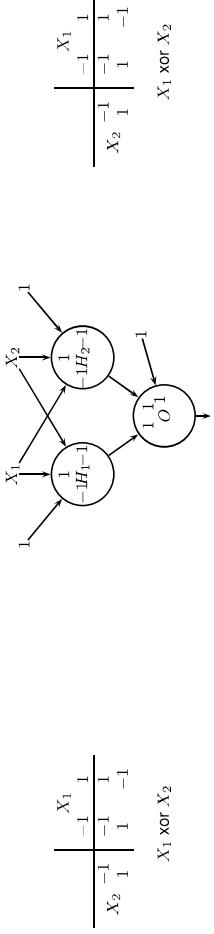


MI E 19

Expressive power

20

Can multiple neurons help?



MI E 19

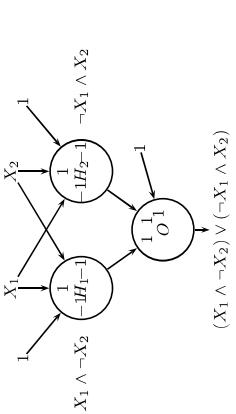
Expressive power

20

XOR again

XOR again

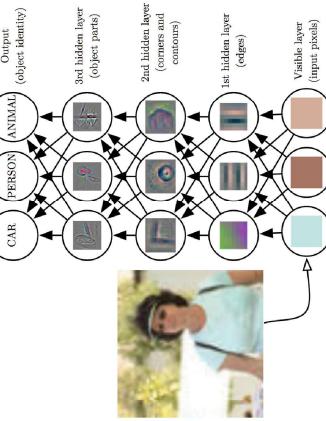
Can multiple neurons help?



MI E 19

Expressive power

20



MI E 19

Expressive power

21

The task of learning



MI E 19

Expressive power

22

playground.tensorflow.org

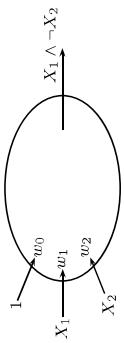
Expressive power

23

The task of learning

Learning weights and threshold |

Learning weights and threshold |



$$X_1 \wedge \neg X_2$$

$$\begin{array}{c|ccccc} & & X_1 & & \\ & & -1 & 1 & \\ \hline X_2 & -1 & -1 & 1 & -1 \\ 1 & 1 & & & \end{array}$$

$$X_1 \wedge \neg X_2$$

We have:

- $\mathcal{D} = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4)$ input vectors (cases).

- $\mathbf{t} = (-1, 1, -1, -1)$ vector of target outputs.

- $\mathbf{w} = (w_0, w_1, w_2)$ vector of current parameters.

- $\mathbf{o} = (o_1, o_2, o_3, o_4)$ vector of current outputs.

We request:

- \mathbf{w}^* parameters yielding $\mathbf{o} = \mathbf{t}$.

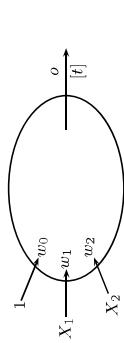
ML E 19

The task of learning

23

Learning weights and threshold |

Example: $X_1 \wedge X_2$



$$\text{Error: } E = t - o$$

$$X_1 \wedge X_2$$

Weight updating procedure

The weights are updated as follows:

- $E > 0 \Rightarrow o$ shall be increased $\Rightarrow \mathbf{x} \cdot \mathbf{w}$ up $\Rightarrow \mathbf{w} := \mathbf{w} + \alpha E \mathbf{x}$
- $E < 0 \Rightarrow o$ shall be decreased $\Rightarrow \mathbf{x} \cdot \mathbf{w}$ down $\Rightarrow \mathbf{w} := \mathbf{w} - \alpha E \mathbf{x}$

α is called the learning rate.

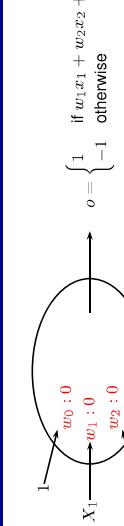
With \mathcal{D} linearly separable and α not too large this procedure will converge to a correct set of parameters.

ML E 19

The task of learning

23

Example: $X_1 \wedge X_2$



$$X_1 \wedge \neg X_2$$

$$\begin{array}{c|ccccc} & & X_1 & & \\ & & -1 & 1 & \\ \hline X_2 & -1 & -1 & 1 & -1 \\ 1 & 1 & & & \end{array}$$

$$X_1 \wedge \neg X_2$$

Cases:

$$\mathbf{t}: \quad (-1, 1, -1, -1)$$

$$\mathbf{o}: \quad \left(\begin{array}{c} 1 \\ -1 \\ 1 \\ -1 \end{array} \right)$$

$$E: \quad \left(\begin{array}{c} 1 \\ -1 \\ 1 \\ -1 \end{array} \right)$$

$$\mathbf{w} := \left(\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right) + \frac{1}{4} \cdot 2 \cdot \left(\begin{array}{c} 1 \\ -1 \\ 1 \\ -1 \end{array} \right) = \left(\begin{array}{c} \frac{1}{2} \\ -\frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \end{array} \right)$$

$\alpha = \frac{1}{4}$

$$X_1 \wedge \neg X_2$$

Cases:

$$\mathbf{t}: \quad (1, 1, -1, 1)$$

$$\mathbf{o}: \quad \left(\begin{array}{c} 1 \\ -1 \\ 1 \\ -1 \end{array} \right)$$

$$E: \quad \left(\begin{array}{c} 1 \\ -1 \\ 1 \\ -1 \end{array} \right)$$

$$\mathbf{w} := \left(\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right) + \frac{1}{4} \cdot 2 \cdot \left(\begin{array}{c} 1 \\ -1 \\ 1 \\ -1 \end{array} \right) = \left(\begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right)$$

$$X_1 \wedge \neg X_2$$

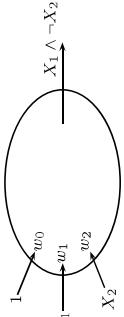
ML E 19

The task of learning

24

Learning weights and threshold |

Learning weights and threshold |



$$X_1 \wedge \neg X_2$$

$$\begin{array}{c|ccccc} & & X_1 & & \\ & & -1 & 1 & \\ \hline X_2 & -1 & -1 & 1 & -1 \\ 1 & 1 & & & \end{array}$$

$$X_1 \wedge \neg X_2$$

We have:

- $\mathcal{D} = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4)$ input vectors (cases).

- $\mathbf{t} = (-1, 1, -1, -1)$ vector of target outputs.

- $\mathbf{w} = (w_0, w_1, w_2)$ vector of current parameters.

- $\mathbf{o} = (o_1, o_2, o_3, o_4)$ vector of current outputs.

We request:

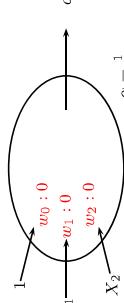
- \mathbf{w}^* parameters yielding $\mathbf{o} = \mathbf{t}$.

ML E 19

The task of learning

23

Example: $X_1 \wedge X_2$



$$X_1 \wedge X_2$$

$\alpha = \frac{1}{4}$

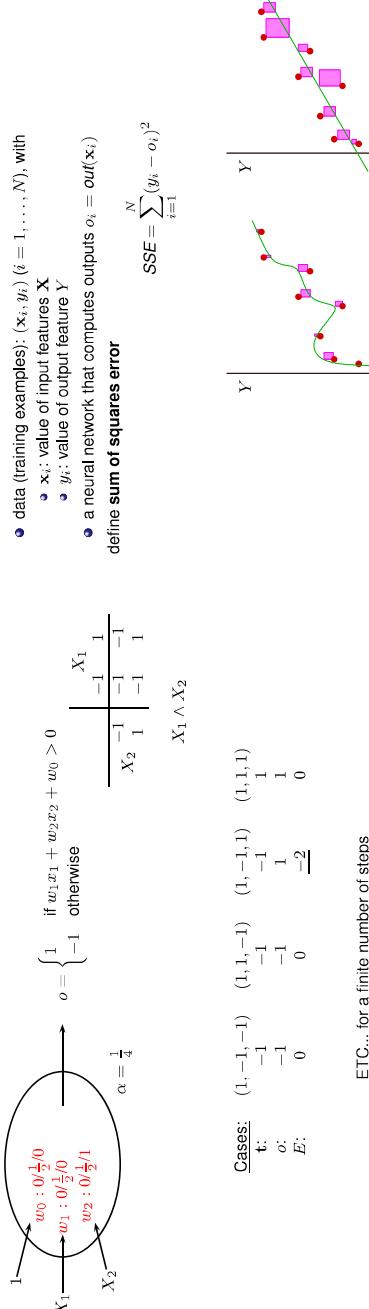
$$X_1 \wedge X_2$$

ML E 19

The task of learning

24

Example: $X_1 \wedge X_2$



ETC... for a finite number of steps

Set of training examples (red dots), function defined by neural network (green), squared errors for each training example (area of magenta squares).

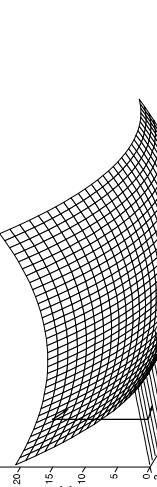
[M1 E 19](#) [The task of learning](#) 25

Minimizing SSE

For a given dataset, the SSE is a smooth, convex function of the weights.

Example for $n = 2$ (and a linear activation function):

$$\frac{\partial E}{\partial w_i} = \sum_{k=1}^N (t_k - \mathbf{w} \cdot \mathbf{x}_k) (-x_{k,i})$$



~~ weights \mathbf{w} that minimize $E(\mathbf{w})$ can be found by gradient descent.

Properties:

- The procedure converges to the weights \mathbf{w} that minimize the SSE (if η is small enough).

[M1 E 19](#) [The task of learning](#) 26

Stochastic Gradient Descent

Variation of gradient descent: instead of following the gradient computed from the whole dataset:

$$\frac{\partial E}{\partial w_i} = \sum_{k=1}^N (t_k - \mathbf{w} \cdot \mathbf{x}_k) (-x_{k,i}),$$

iterate through the data instances one by one, and in one iteration follow the gradient defined by a single data instance (\mathbf{x}_k, t_k) :

$$\frac{\partial E}{\partial w_i} = (t_k - \mathbf{w} \cdot \mathbf{x}_k) (-x_{k,i}),$$

Find the weights that minimize the *sum of squared errors (SSE)*

$$\sum_{i=1}^N \sum_{j=1}^m (t_{j,i} - o_{j,i})^2,$$

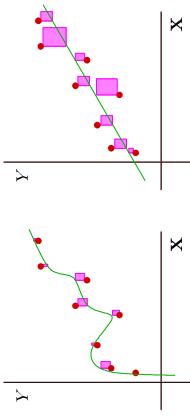
where $o_{j,i}$ is the value of the j 'th output neuron for the i 'th data instance.

[M1 E 19](#) [The task of learning](#) 29

Sum of Squared Errors

Given:

- data (training examples): (\mathbf{x}_i, y_i) ($i = 1, \dots, N$), with
 - \mathbf{x}_i : value of input features \mathbf{X}
 - y_i : value of output feature Y
- a neural network that computes outputs $o_i = \text{out}(\mathbf{x}_i)$



Set of training examples (red dots), function defined by neural network (green), squared errors for each training example (area of magenta squares).

[M1 E 19](#) [The task of learning](#) 25

Gradient Descent Learning

The gradient is the vector of partial derivatives:

$$\nabla E[\mathbf{w}] = \left(\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right)$$

The partial derivatives are (with a linear activation function):

$$\frac{\partial E}{\partial w_k} = \sum_{i=1}^N (t_i - \mathbf{w} \cdot \mathbf{x}_i) (-x_{i,k}).$$

Gradient descent rule:

```
Initialize w with random values
repeat
    w := w - η ∇E(w)
until ∇E(w) ≈ 0
(η is again a small constant,
the learning rate).
```

Properties:

- The procedure converges to the weights \mathbf{w} that minimize the SSE (if η is small enough).

[M1 E 19](#) [The task of learning](#) 26

Neural Networks

Given: structure and activation functions. To be learned: weights.

Goal: given the training examples

X_1	Input	X_2	...	X_n	X_1	Output	X_2	...	X_m
$x_{1,1}$	$x_{2,1}$...		$x_{n,1}$	$t_{1,1}$	$t_{2,1}$...		$t_{m,1}$
$x_{1,2}$	$x_{2,2}$...		$x_{n,2}$	$t_{1,2}$	$t_{2,2}$...		$t_{m,2}$
:	:	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
$x_{1,N}$	$x_{2,N}$...		$x_{n,N}$	$t_{1,N}$	$t_{2,N}$...		$t_{m,N}$

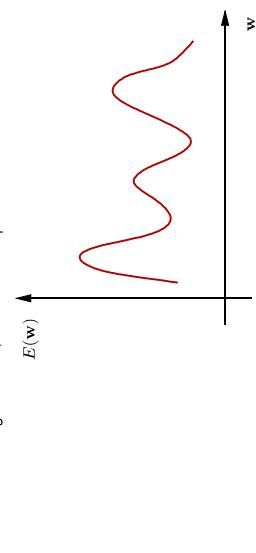
Find the weights that minimize the *sum of squared errors (SSE)*

$$\sum_{i=1}^N \sum_{j=1}^m (t_{j,i} - o_{j,i})^2,$$

where $o_{j,i}$ is the value of the j 'th output neuron for the i 'th data instance.

[M1 E 19](#) [The task of learning](#) 29

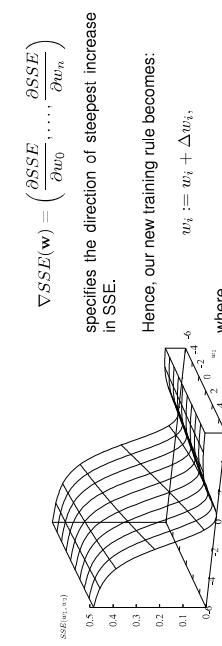
- As for perceptron with SSE error:
- Error is smooth function of weights \mathbf{w}
 - Can use gradient descent to optimize weights
 - Error no longer convex, can have multiple local minima:



- Partial derivatives more difficult to compute

Learning

- Basic principle:** SSE is a differentiable function of the weights (for differentiable activation functions such as the sigmoid function!). Use gradient descent to optimize SSE:



specifies the direction of steepest increase in SSE.

Hence, our new training rule becomes:

$$w_i := w_i + \Delta w_i,$$

$$\Delta w_i = -\eta \frac{\partial SSE}{\partial w_i}$$

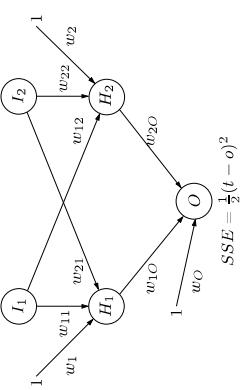
- In practice: use the back propagation algorithm (approximation of gradient descent)

MI E 19 | The task of learning | 30

Neural Networks

The Principle of Back Propagation

Training examples provide target values for only network outputs, so no target values are directly available for indicating the error of the hidden units values.



$$SSE = \frac{1}{2}(t - o)^2$$

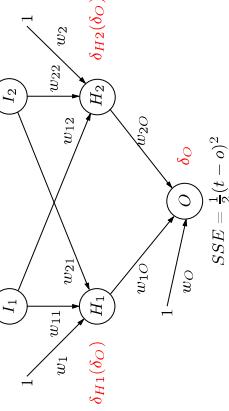
Idea: Calculate an error term δ_h for a hidden unit by taking the weighted sum of the error terms, δ_i , for each output unit it influences.

MI E 19 | The task of learning | 32

Neural Networks

The Principle of Back Propagation

Training examples provide target values for only network outputs, so no target values are directly available for indicating the error of the hidden units values.

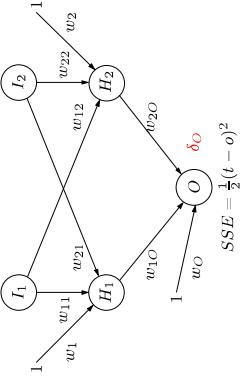


Idea: Calculate an error term δ_h for a hidden unit by taking the weighted sum of the error terms, δ_k , for each output unit it influences.

MI E 19 | The task of learning | 32

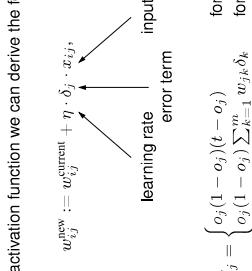
The Principle of Back Propagation

Training examples provide target values for only network outputs, so no target values are directly available for indicating the error of the hidden units values.



$$SSE = \frac{1}{2}(t - o)^2$$

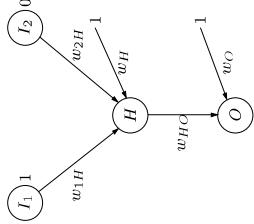
When using a sigmoid activation function we can derive the following updating rule:



$$\delta_j = \begin{cases} o_j(1 - o_j)(t - o_j) & \text{for output nodes,} \\ o_j(1 - o_j) \sum_{k=1}^m w_{jk} \delta_k & \text{for hidden nodes.} \end{cases}$$

MI E 19 | The task of learning | 33

Back Propagation Example



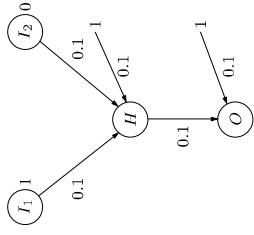
Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

M1 E19

The task of learning

34

Back Propagation Example



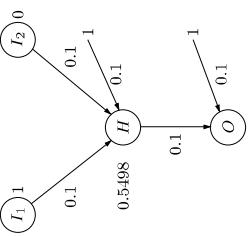
Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

M1 E19

The task of learning

34

Back Propagation Example



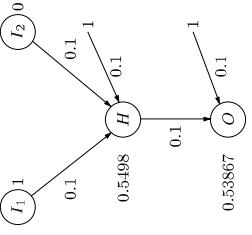
Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

M1 E19

The task of learning

34

Back Propagation Example



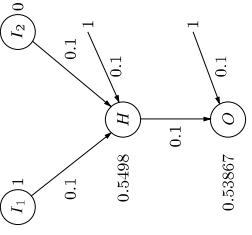
Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

M1 E19

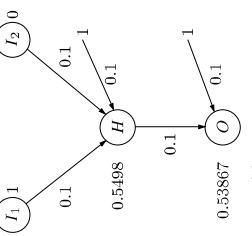
The task of learning

34

Back Propagation Example



Back Propagation Example



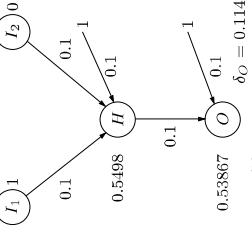
Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

M1 E19

The task of learning

34

Back Propagation Example



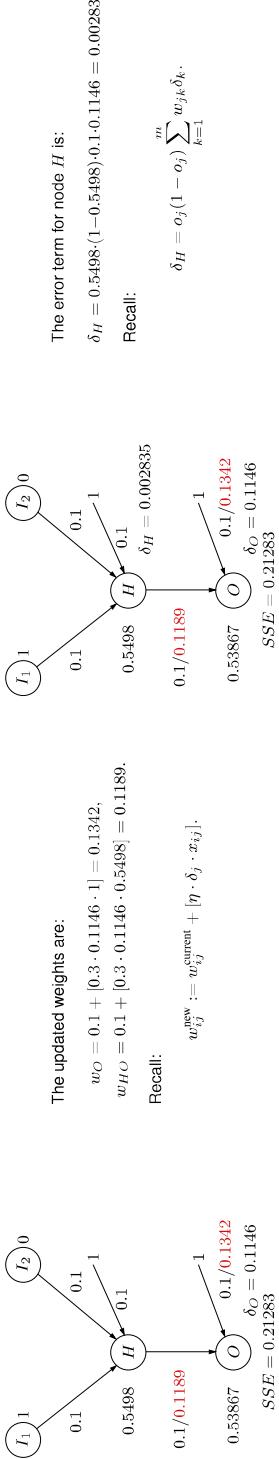
Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

M1 E19

The task of learning

34

Back Propagation Example



Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

MI E 19

The task of learning

MI E 19

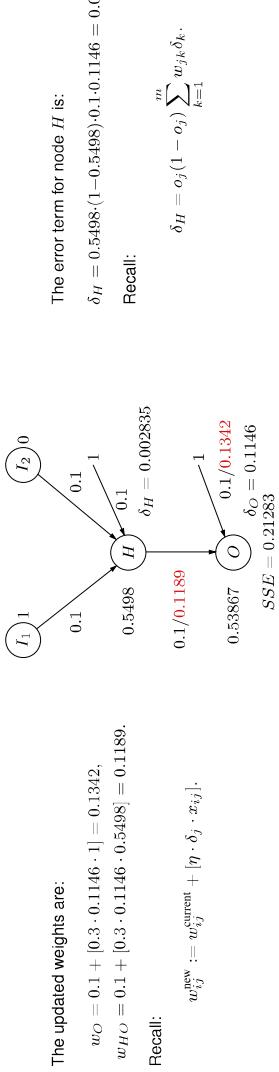
The task for node H is:

$$\delta_H = 0.5498 \cdot (1 - 0.5498) \cdot 0.1 \cdot 0.1146 = 0.002836.$$

Recall:

$$\delta_H = o_j(1 - o_j) \sum_{k=1}^m w_{jk} \delta_k.$$

Back Propagation Example



Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

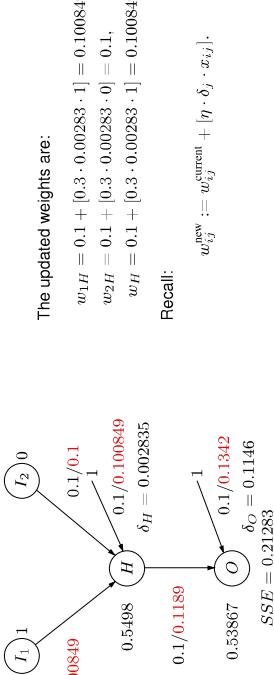
MI E 19

The task of learning

MI E 19

Neural Networks

Back Propagation Example



Assume that we have the training example ($I_1 = 1, I_2 = 0, O = 1$).

MI E 19

The task of learning

MI E 19

Tentative course overview

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Inference in Bayesian networks
- Machine learning
- Planning
- Clustering
- Multi-agent systems

Probabilistic Models



Word occurrence in emails (thousands of input features):

Mail	abacus	...	informatics	...	pills	...	watch	...	zytogenic	Spam
m_1	n	...	y	...	n	...	n	...	n	no
m_2	n	...	n	...	n	...	n	...	y	yes
m_3	n	...	n	...	n	...	n	...	y	no
m_4	n	...	n	...	n	...	n	...	y	yes
m_5	n	...	n	...	n	...	n	...	y	yes
m_6	n	...	n	...	n	...	n	...	n	no
m_7	n	...	n	...	n	...	n	...	n	yes
m_8	n	...	y	...	n	...	n	...	y	yes
m_9	n	...	y	...	n	...	y	...	n	yes

- Need to learn entries in conditional probability tables.
- Simplest approach: use empirical frequencies, e.g.:

$$P(\text{pills} = y \mid \text{Spam} = \text{yes}) = \frac{\#\text{mails with pills} = y \text{ and } \text{Spam} = \text{yes}}{\#\text{mails with } \text{Spam} = \text{yes}}$$

Example: Learning a Naive Bayes

Example	Author	Thread	Length	WhereRead	UserAction
e1	known	new	long	home	skips
e2	unknown	new	short	work	reads
e3	unknown	follow Up	long	home	skips
e4	known	new	long	home	reads
e5	known	follow Up	long	work	skips
e6	known	follow Up	short	work	reads
e7	unknown	new	short	work	reads
e8	known	follow Up	long	home	skips
e9	known	new	long	work	skips
e10	unknown	follow Up	short	home	skips
e11	known	new	long	work	reads
e12	known	follow Up	short	home	skips
e13	known	new	long	work	reads
e14	known	new	short	work	reads
e15	known	new	short	home	reads
e16	known	follow Up	short	work	reads
e17	known	new	short	home	reads
e18	unknown	new	short	work	reads

Example: Learning a Naive Bayes

$$\text{Probabilities to estimate:}$$

$$P(\text{reads}) = \frac{9}{18}$$

$$P(\text{known}|\text{reads}) = \frac{2}{3}$$

$$P(\text{known}|\text{skips}) = \frac{2}{3}$$

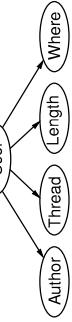
$$P(\text{new}|\text{reads}) = \frac{2}{3}$$

$$P(\text{new}|\text{skips}) = \frac{2}{3}$$

$$P(\text{new}|\text{home}) = \frac{2}{3}$$

$$P(\text{new}|\text{work}) = \frac{2}{3}$$

NB model:



Example: Learning a Naive Bayes

$$\text{Probabilities to estimate:}$$

$$P(\text{reads}) = \frac{9}{18}$$

$$P(\text{known}|\text{reads}) = \frac{2}{3}$$

$$P(\text{known}|\text{skips}) = \frac{2}{3}$$

$$P(\text{new}|\text{reads}) = \frac{2}{3}$$

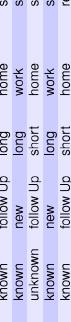
$$P(\text{new}|\text{skips}) = \frac{2}{3}$$

$$P(\text{new}|\text{home}) = \frac{2}{3}$$

$$P(\text{new}|\text{work}) = \frac{2}{3}$$

$$P(\text{new}|\text{skips}) = \frac{2}{3}$$

NB model:



Example: Learning a Naive Bayes

$$\text{Probabilities to estimate:}$$

$$P(\text{reads}) = \frac{9}{18}$$

$$P(\text{known}|\text{reads}) = \frac{2}{3}$$

$$P(\text{known}|\text{skips}) = \frac{2}{3}$$

$$P(\text{new}|\text{reads}) = \frac{2}{3}$$

$$P(\text{new}|\text{skips}) = \frac{2}{3}$$

NB model:



The naive Bayes assumption I

Target: $Symbol \in \{A, \dots, Z, 0, \dots, 9\}$

Predictors: $Cell-1, \dots, Cell-9 \in \{b, w\}$.

1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9

1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9

For example:

$$P(Cell - 2 = b | Cell - 5 = b, Symbol = 1) > P(Cell - 2 = b | Symbol = 1)$$

Attributes not independent given $Symbol = 1$

The naive Bayes assumption II

Target: $Spam \in \{y, n\}$

Predictors: $Subject-all-caps, Known-spam-server, \dots, ContainsMoney' \in \{y, n\}$.

The naive Bayes assumption II

Target: $Spam \in \{y, n\}$

Predictors: $Subject-all-caps, Known-spam-server, \dots, ContainsMoney' \in \{y, n\}$.

For example:

$$\begin{aligned} P(\text{Body-nigeria} = y | \text{Body-confidential} = y, \text{Spam} = y) \\ \gg \\ P(\text{Body-nigeria} = y | \text{Spam} = y) \end{aligned}$$

Attributes not independent given $Spam = yes$

Naive Bayes assumption often not realistic. Nevertheless, Naive Bayes often successful.

The paradoxical success of Naive Bayes

One explanation for the surprisingly good performance of Naive Bayes in many domains: do not require exact distribution for classification, only the right decision boundaries [Domingos, Pazzani 97]

$$\oplus : P(C = \oplus | a_1, \dots, a_n) \text{ (real)}$$

$$\oplus : P(C = \oplus | a_1, \dots, a_n) \text{ (Naive Bayes)}$$



The naive Bayes assumption I

Target: $Symbol \in \{A, \dots, Z, 0, \dots, 9\}$

Predictors: $Cell-1, \dots, Cell-9 \in \{b, w\}$.

1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9

1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9
1	2	3	4	5	6	7	8	9

For example:

$$P(Cell - 2 = b | Cell - 5 = b, Symbol = 1) > P(Cell - 2 = b | Symbol = 1)$$

Attributes not independent given $Symbol = 1$

When Naive Bayes must fail

No Naive Bayes Classifier can produce the following classification:

A	B	Class
yes	⊕ no	⊕ ⊖
yes	no	⊖ ⊕
no	no	⊕

because assume it did, then:

1. $P(A = \text{y} | \oplus)P(B = \text{y} | \oplus)P(\oplus) > P(A = \text{y} | \ominus)P(B = \text{y} | \ominus)P(\ominus)$
2. $P(A = \text{y} | \ominus)P(B = \text{n} | \ominus)P(\ominus) > P(A = \text{y} | \oplus)P(B = \text{n} | \oplus)P(\oplus)$
3. $P(A = \text{n} | \ominus)P(B = \text{y} | \ominus)P(\ominus) > P(A = \text{n} | \oplus)P(B = \text{y} | \oplus)P(\oplus)$
4. $P(A = \text{n} | \oplus)P(B = \text{n} | \oplus)P(\oplus) > P(A = \text{n} | \ominus)P(B = \text{n} | \ominus)P(\ominus)$

MI E19

Probabilistic Models

11

Probabilistic Models

12

Multiplying the four left sides and the four right sides of these inequalities:

$$\prod_{i=1}^4 (\text{left side of } i.) > \prod_{i=1}^4 (\text{right side of } i.)$$

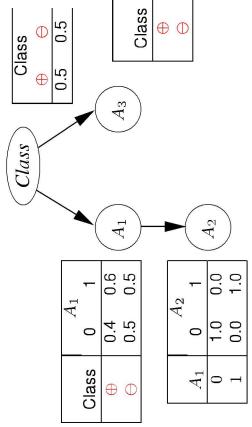
But this is false, because both products are actually equal.

MI E19

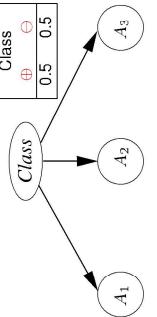
A word of caution

When features don't help

Data generated by process described by Bayesian network:



The Naive Bayes model learned from data:



In Naive Bayes model:

$$P(\oplus | A_1 = 1, A_2 = 1, A_3 = 0) = 0.507$$

Intuitively, the NB model double counts the information provided by A_1, A_2 . Recall our previous discussion on the use of intermediate variables!

MI E19

Probabilistic Models

13

14

Attribute A_2 is just a duplicate of A_1 . Conditional class probability for example:

$$P(\oplus | A_1 = 1, A_2 = 1, A_3 = 0) = 0.461$$

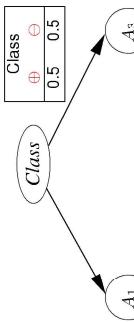
MI E19

Probabilistic Models

15

A word of caution

The Naive Bayes model with selected features A_1 and A_3 :



Case-based Reasoning

Class	0	1
⊕	0.5 0.5	0.5 0.5
⊖	0.7 1.0	0.3 0.3

In this Naive Bayes model:

$$P(\oplus | A_1 = 1, A_3 = 0) = 0.461$$

(and all other posterior class probabilities also are the same as in the true model).

MI E19

Probabilistic Models

16

Case-based Reasoning

Distances for numeric features

- To predict the output feature of a new example \mathbf{x}' :
- find among the training examples the one (several) most similar to \mathbf{x}'
 - predict the output for \mathbf{x}' from the known output values of the similar cases

Several names for this approach:

- Instance based learning
- Lazy learning
- Case-based reasoning

Required : **distance measure** on values of input features.

If all features \mathbf{X} are numeric:

- Euclidean distance: $d(\mathbf{x}, \mathbf{x}') = \sqrt{\sum_i (x_i - x'_i)^2}$
- Manhattan distance: $d(\mathbf{x}, \mathbf{x}') = \sum_i |x_i - x'_i|$

Distances for discrete features

For a single feature X with domain $\{x_1, \dots, x_k\}$:

- Zero-One distance: $d(x_i, x_j) = 0$ if $j = i$, $d(x_i, x_j) = 1$ if $j \neq i$
- Distance matrix: specify for each pair x_i, x_j the distance $d(x_i, x_j)$ in a $k \times k$ -matrix. Example:

	low	medium	high
low	0	2	5
medium	2	0	1
high	5	1	0

For a set of discrete features \mathbf{X} :

- Define distance d_i and weight w_i for each $X_i \in \mathbf{X}$
- Define $d(\mathbf{x}, \mathbf{x}') = \sum_i w_i d_i(x_i, x'_i)$

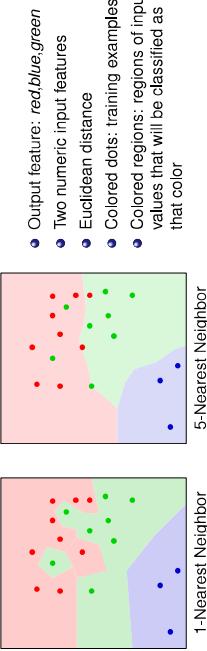
MI E 19 Case-based Reasoning 17

 K -Nearest-Neighbor Classifier

Given

- training examples (\mathbf{x}_i, y_i) ($i = 1, \dots, N$)
 - a new case \mathbf{x} to be classified
 - a distance measure $d(\mathbf{x}, \mathbf{x}')$
- classify \mathbf{x} as follows:
- find the K training examples $\mathbf{x}_{i1}, \dots, \mathbf{x}_{iK}$ with minimal distance to \mathbf{x}
 - predict for \mathbf{x} the most frequent value in y_{i1}, \dots, y_{iK} .

Example



MI E 19

Case-based Reasoning 18

Overfitting

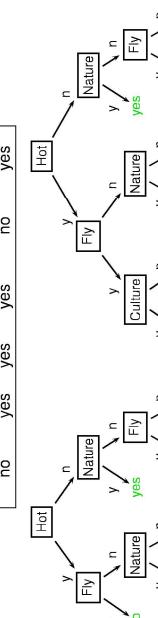
Holiday Data

Culture	Fly	Hot	Music	Nature	Likes
no	no	yes	no	no	no
no	yes	yes	no	no	no
yes	yes	yes	yes	no	no
no	yes	yes	no	yes	no
yes	no	no	yes	yes	yes
no	no	no	yes	yes	yes
yes	yes	yes	no	no	no
yes	yes	yes	yes	yes	yes
yes	yes	yes	no	yes	no
yes	no	no	yes	yes	yes
yes	yes	yes	no	no	no
yes	yes	yes	yes	no	no
yes	yes	yes	yes	yes	yes
yes	yes	yes	no	yes	no
no	no	no	yes	yes	yes
yes	no	yes	no	no	no
yes	yes	yes	no	no	no
yes	no	no	yes	no	no
yes	yes	yes	yes	no	no
yes	yes	yes	yes	yes	yes
no	no	no	no	no	yes

Overfitting: Decision Trees

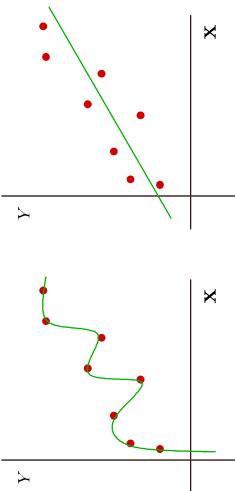
Decision tree learned from holiday data (left) and holiday data augmented with one more example

Culture	Fly	Hot	Music	Nature	Likes
no	no	yes	yes	yes	no
no	yes	yes	yes	yes	no



Overfitting: Neural Networks

Overfitting: Neural Networks



Left model: minimizes SSE on training examples

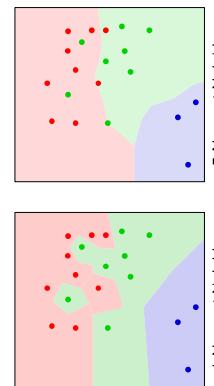
Right model: may have smaller SSE on **future observations**

MI E19

21

Oversimplifying

Overfitting: Nearest Neighbor



1-Nearest Neighbor correctly classifies all training examples

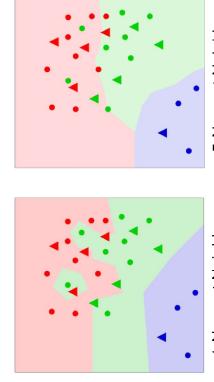
5-Nearest Neighbor may be better for future observations (triangles)

MI E19

22

Oversimplifying

Overfitting: Nearest Neighbor



1-Nearest Neighbor correctly classifies all training examples

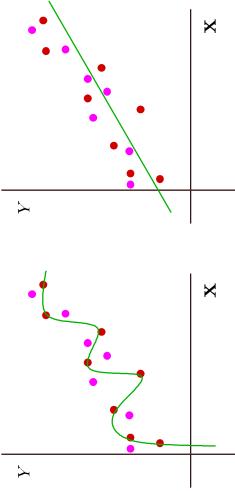
5-Nearest Neighbor may be better for future observations (triangles)

MI E19

21

Oversimplifying

Overfitting: Nearest Neighbor



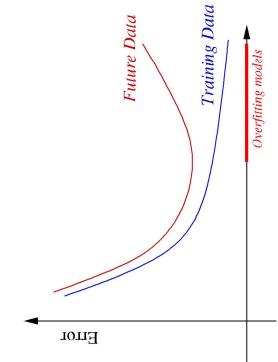
Left model: minimizes SSE on training examples

Right model: may have smaller SSE on **future observations**

MI E19

21

Oversimplifying



Model Complexity

Future Data

Training Data

Error

Oversimplifying model

Model	Complexity	Error
Decision Tree	Depth, # Nodes	Classification error
Neural Network	# hidden nodes/layers	SSE
Nearest Neighbor	Decision regions	Classification error

A model **overfits** the training data, if a smaller error on future data would be obtained by a less complex model.

MI E19

23

Oversimplifying

Reduced hypothesis space

Do not allow overly complex models:

- Decision Trees: use an early stopping criterion, e.g. stop construction when (sub-) set of training examples contains fewer than k elements (for not too small k).
- Add to evaluation measure a **penalty term** for complexity. This penalty term can have an interpretation as a **prior model probability**, or **model description length**. These techniques will usually lead to more complex models only when the data strongly supports it (especially: large number of examples).

- Do not allow to learn models that are too complex (relative to the available data):
- Decision Trees: use an early stopping criterion, e.g. stop construction when (sub-) set of training examples contains fewer than k elements (for not too small k).
 - Add to evaluation measure a **penalty term** for complexity. This penalty term can have an interpretation as a **prior model probability**, or **model description length**. These techniques will usually lead to more complex models only when the data strongly supports it (especially: large number of examples).

Modified Search/Scoring

MI E19

24

Oversimplifying

Basic idea

Reserve part of the training examples as a **validation set**:

- not used during model search
- used as proxy for future data in model evaluation

Train/Validation split

Simplest approach: split the available data randomly into a **training** and a **validation** set.
Typically: 50%/50% or 66%/33% split.

Post pruning

Use of validation set in decision tree learning:

- Build decision tree using **training** set
- For each internal node:
 - test whether accuracy on **validation** set is improved by replacing sub-tree rooted at this node by single leaf (labeled with most frequent target feature value of **training** examples in this sub-tree)
 - if yes.: replace sub-tree with leaf (*prune* sub-tree)

MI E 19

Overfitting

26

Example: K -NN and Neural Networks

Selection of K

for $K = 1, 2, 3, \dots$:

- measure accuracy of K -NN based on **training** examples for **validation** examples
- use K with best performance on **validation** examples
- **validation** examples can now be merged with **training** examples for predicting future cases

Selection of K

for $#hl = 1, 2, \dots, \max$, and $#nd = 1, 2, \dots, \max$:

- learn Neural Network with $#hl$ hidden layers and $#nd$ nodes per hidden layer using **training** examples
- use K with best performance on **validation** examples
- evaluate SSE of learned network on **validation** examples
- select network structure with minimal SSE
- re-learn weights using merged training and validation examples

Selection of Neural Network Structure

for $#hl = 1, 2, \dots, \max$, and $#nd = 1, 2, \dots, \max$:

- learn Neural Network with $#hl$ hidden layers and $#nd$ nodes per hidden layer using **training** examples
- evaluate SSE of learned network on **validation** examples
- select network structure with minimal SSE
- re-learn weights using merged training and validation examples

Cross-Validation

Disadvantage of training/validation splits (for small datasets):

- the examples in the validation set are partly "wasted"
- (unrepresentative) patterns in the validation set can bias the model selection

Cross-Validation

- Divide the examples into n equal sized subsets, or *folds* (e.g. $n = 10$)
- for $i = 1, \dots, n$:
 - learn model (with given "complexity setting") using folds $1, \dots, i - 1, i + 1, \dots, n$
 - evaluate using fold i
 - average the evaluation scores from the n folds
- Choose "complexity setting" that obtained highest average evaluation score
- Learn final model with chosen "complexity setting" using all available examples

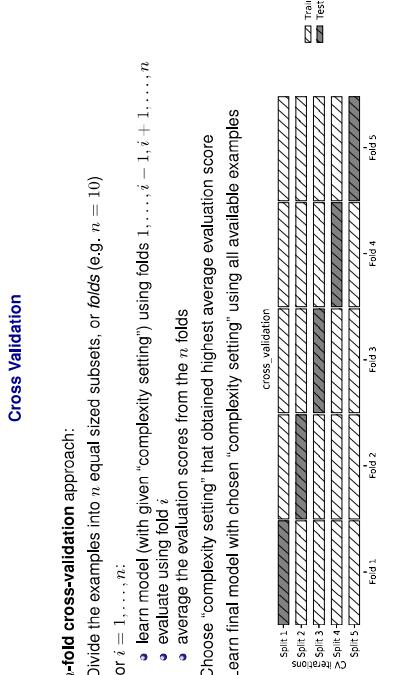
Cross Validation

The *n-fold cross-validation* approach:

- Divide the examples into n equal sized subsets, or *folds* (e.g. $n = 10$)
- for $i = 1, \dots, n$:
 - learn model (with given "complexity setting") using folds $1, \dots, i - 1, i + 1, \dots, n$
 - evaluate using fold i
 - average the evaluation scores from the n folds
- Choose "complexity setting" that obtained highest average evaluation score
- Learn final model with chosen "complexity setting" using all available examples

Cross Validation

cross-validation



Post-pruning

Overfitting

28

Machine Intelligence

Lecture 10: Clustering

Thomas Dyrhe Nielsen

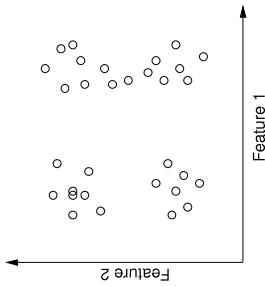
Aalborg University

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty.
- Bayesian networks
- Inference in Bayesian networks
- Machine learning: Classification
- **Machine learning: Clustering**
- Planning
- Multi-agent systems

Clustering: Introduction

The objective of clustering is to find structure in the data.



Clustering

- Examples:**
- Based on customer data, find groups of customers with similar profiles.
 - Based on image data, find groups of images with similar motif.
 - Based on article data, find groups of articles with the same topics.
 - ...

Clustering

Iris Data

Measurement of petal width/length and sepal width/length for 150 flowers of 3 different species of Iris.

first reported in:
Fisher R.A. "The use of multiple measurements in taxonomic problems" Annal Eugenics, 7 (1936).

	Attributes			Class variable
	SL	SW	PL	
5.1	3.5	1.4	0.2	Setosa
4.9	3.0	1.4	0.2	Setosa
6.3	2.9	6.0	2.1	Virginica
6.3	2.5	4.9	1.5	Virginica
...



The Iris data with class labels removed:

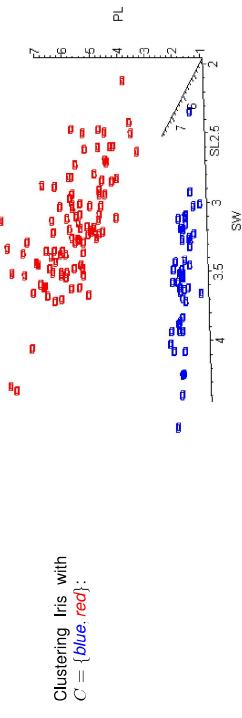
	Attributes			Species
	SL	SW	PL	
5.1	3.5	1.4	0.2	Setosa
4.9	3.0	1.4	0.2	Setosa
6.3	2.9	6.0	2.1	Virginica
6.3	2.5	4.9	1.5	Virginica
...

Clustering

Clustering

Clustered Iris

A clustering of the data $S = \{a_1, \dots, a_N\}$ consists of a set $C = \{c_1, \dots, c_k\}$ of cluster labels, and a cluster assignment $ca : S \rightarrow C$.



MI Autumn 2019

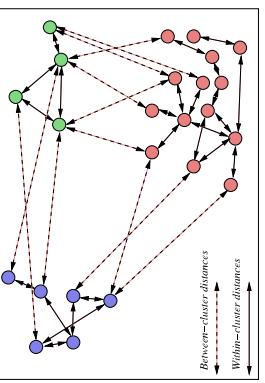
Clustering

MI Autumn 2019

7

Distance Function and Clustering

The k -means algorithm



A candidate clustering (indicated by colors) of data cases in instance space. Arrows indicate between- and within-cluster distances (selected).

General goal: find clustering with

- large between-cluster variation (sum of between-cluster distances)
- small within-cluster variation (sum of within-cluster distances)

MI Autumn 2019

The k -means algorithm

6

The k -means algorithm: Example

$k = 3$:



MI Autumn 2019

The k -means algorithm

MI Autumn 2019

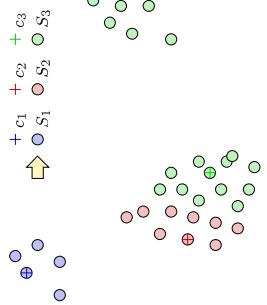
9

The k -means algorithm

9

The k -means algorithm: Example

$k = 3$:



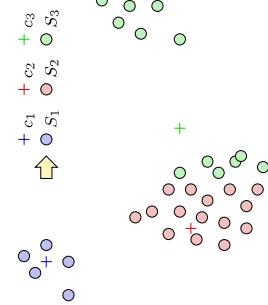
MI Autumn 2019

The k -means algorithm

9

The k -means algorithm: Example

$k = 3$:



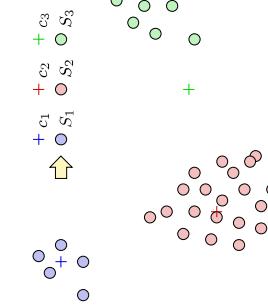
MI Autumn 2019

The k -means algorithm

9

The k -means algorithm: Example

$k = 3$:



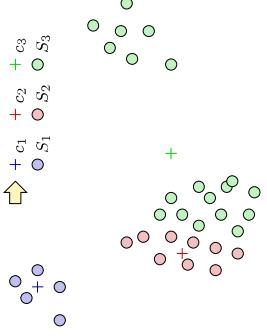
MI Autumn 2019

The k -means algorithm

9

The k -means algorithm: Example

$k = 3$:



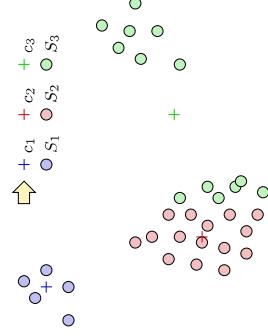
MI Autumn 2019

The k -means algorithm

9

The k -means algorithm: Example

$k = 3$:



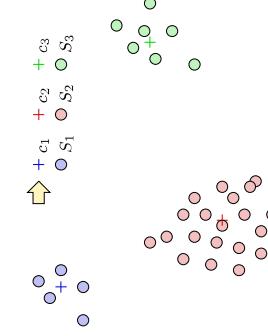
MI Autumn 2019

The k -means algorithm

9

The k -means algorithm: Example

$k = 3$:



MI Autumn 2019

The k -means algorithm

9

The k -means algorithm: Example

$k = 3$:



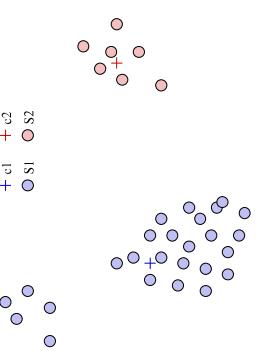
MI Autumn 2019

9

The k -means Algorithm

Different k

Result for clustering the same data with $k = 2$:

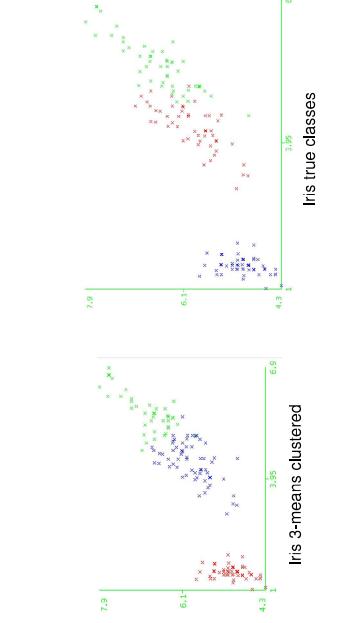


MI Autumn 2019

10

The k -means Algorithm

Clustering Iris



MI Autumn 2019

11

The k -means Algorithm

Result can depend on choice of initial cluster centers!

The k -means algorithm: Background

k -means as an optimization problem

Assume that we use the Euclidean distance d as proximity measure and that the quality of the clustering is measured by the sum of squared errors:

$$SSE = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d(\mathbf{c}_i, \mathbf{x})^2,$$

where:

- \mathbf{c}_i is the i th centroid
- $C_i \subseteq S$ is the points closest to \mathbf{c}_i according to d .

In principle ...

We can minimize the SSE by looking at all possible partitionings \leadsto not feasible!

Instead, k -means

The centroid that minimizes the SSE is the *mean* of the data-points in that cluster:

$$\mathbf{c}_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}$$

Local optimum found by alternating between cluster assignments and centroid estimation.

MI Autumn 2019

12

The k -means Algorithm

13

The k -means algorithm: Background

Convergence

- The k -means algorithm is guaranteed to converge
- Each step reduces the sum of squared errors.
 - There is only a finite number of cluster assignments.

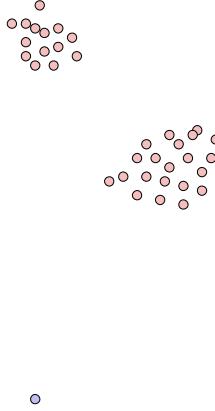
There is no guarantee of reaching the global optimum:

- Improve by running with multiple random restarts.

MI Autumn 2019 13 The k -means Algorithm Some practical issues

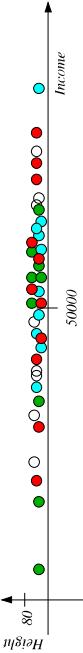
Outliers

The result of partitional clustering can be skewed by outliers. Example with $k = 2$:



Instances defined by attributes

$$A_1 = \text{height in inches} \text{ and } A_2 = \text{annual income in \$} :$$



- all distance functions for continuous attributes dominated by *income* values
- \rightsquigarrow may need to rescale or normalize continuous attributes

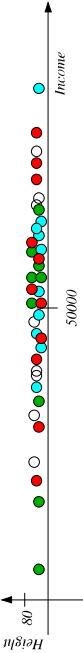
\rightsquigarrow useful preprocessing: outlier detection and elimination.

MI Autumn 2019 14 Some practical issues

Different Measuring Scales

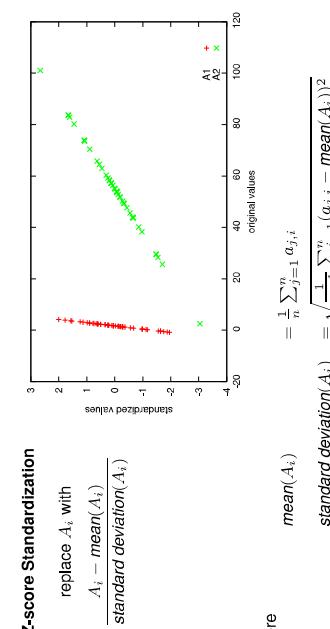
Instances defined by attributes

$$A_1 = \text{height in inches} \text{ and } A_2 = \text{annual income in \$} :$$



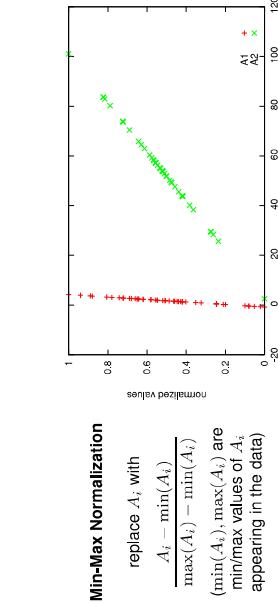
- all distance functions for continuous attributes dominated by *income* values
- \rightsquigarrow may need to rescale or normalize continuous attributes

MI Autumn 2019 15 Some practical issues



$$\begin{aligned} \text{Z-score Standardization} \\ \text{replace } A_i \text{ with} \\ \frac{A_i - \text{mean}(A_i)}{\text{standard deviation}(A_i)} \end{aligned}$$

$$\begin{aligned} &= \frac{1}{n} \sum_{j=1}^n a_{j,i} \\ &\text{standard deviation}(A_i) = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (a_{j,i} - \text{mean}(A_i))^2} \end{aligned}$$



$$\begin{aligned} \text{Min-Max Normalization} \\ \text{replace } A_i \text{ with} \\ \frac{A_i - \text{min}(A_i)}{\text{max}(A_i) - \text{min}(A_i)} \\ (\text{min}(A_i), \text{max}(A_i)) \text{ are} \\ \text{min/max values of } A_i \\ \text{appearing in the data} \end{aligned}$$

MI Autumn 2019 16 Some practical issues

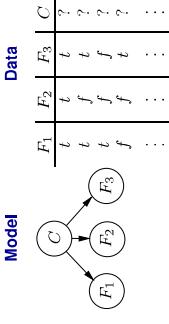
MI Autumn 2019 17 Some practical issues

Soft clustering

The k -means algorithm generates a *hard* clustering; each example is assigned to a single cluster.
Alternatively: In *soft* clustering each example is assigned to a cluster with a certain probability.

Soft clustering

The naive Bayes model for clustering



- C is the hidden cluster variable.
- F_1, F_2 , and F_3 are the feature variables.

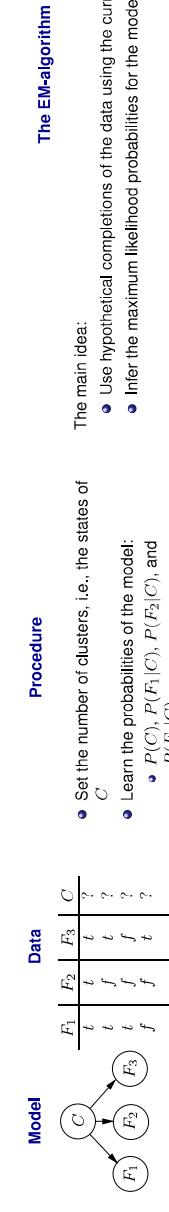
MI Autumn 2019 18 Soft clustering

Soft clustering

The k -means algorithm generates a *hard* clustering; each example is assigned to a single cluster.

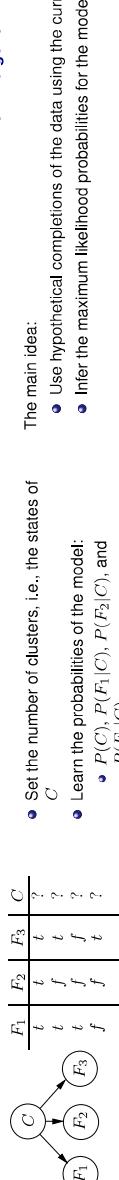
Alternatively: In *soft* clustering each example is assigned to a cluster with a certain probability.

The naive Bayes model for clustering

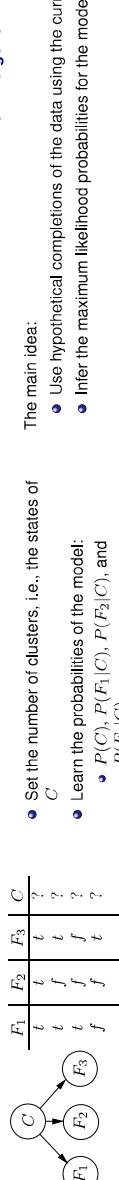
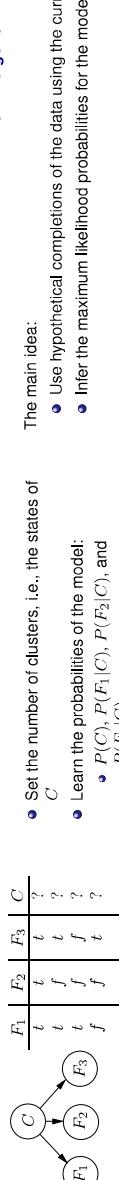


- C is the hidden cluster variable.
- F_1, F_2 , and F_3 are the feature variables.

When learning the probability distributions of the model, the variable C is hidden
 ↵ we cannot directly estimate the probabilities using frequency counts
 Instead we employ the *Expectation-maximization algorithm*

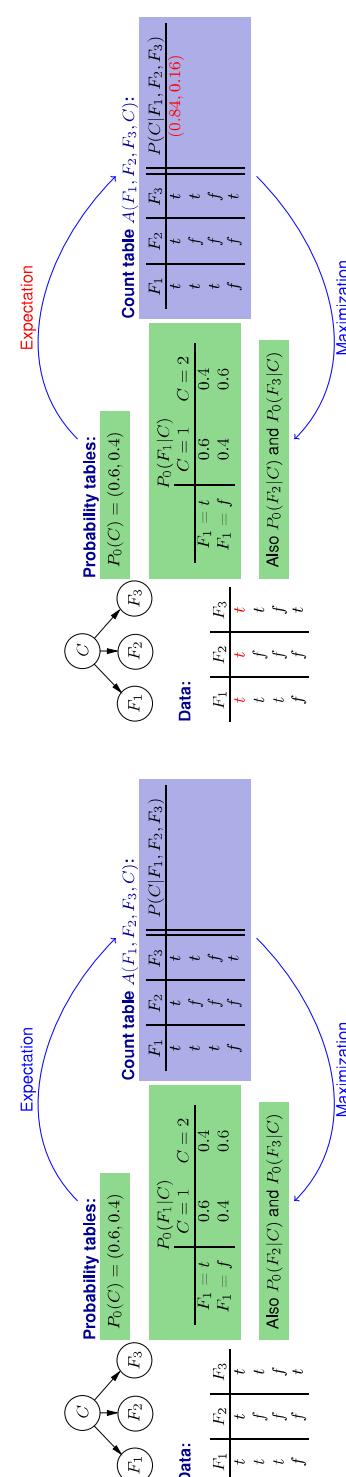


The main idea:
 ↵ we cannot directly estimate the probabilities using frequency counts
 Instead we employ the *Expectation-maximization algorithm*



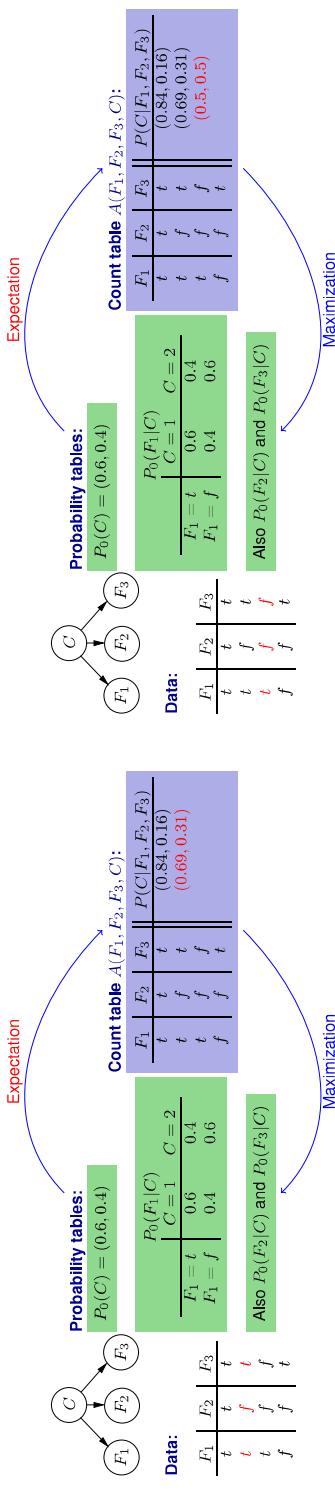
MI Autumn 2019 19 Soft clustering

EM for soft clustering: an example



MI Autumn 2019 20 Soft clustering

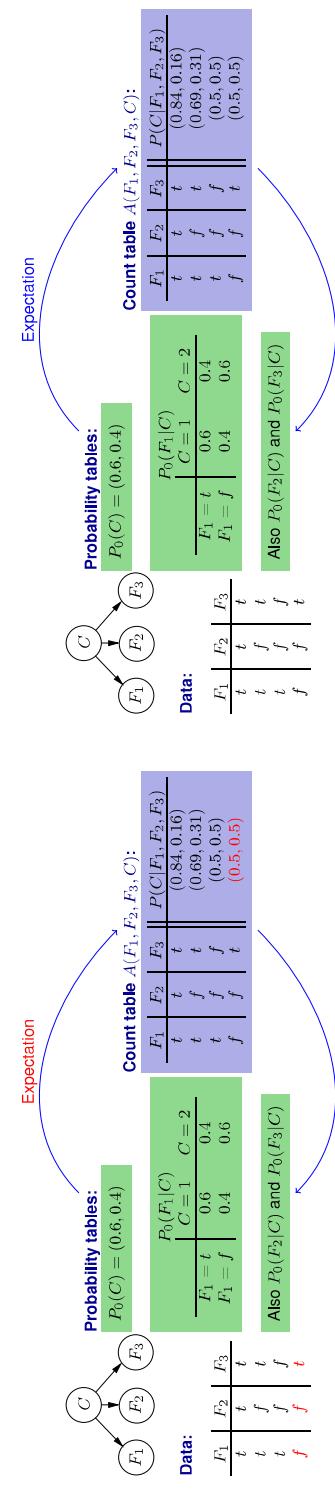
EM for soft clustering: an example



• Fractional counts are being calculated by probability updating.

MI Autumn 2019 | Soft clustering | 20

EM for soft clustering: an example



• Fractional counts are being calculated by probability updating.

MI Autumn 2019 | Soft clustering | 20

Expectation

Probability tables:
 $P_0(C) = (0.6, 0.4)$

Data:

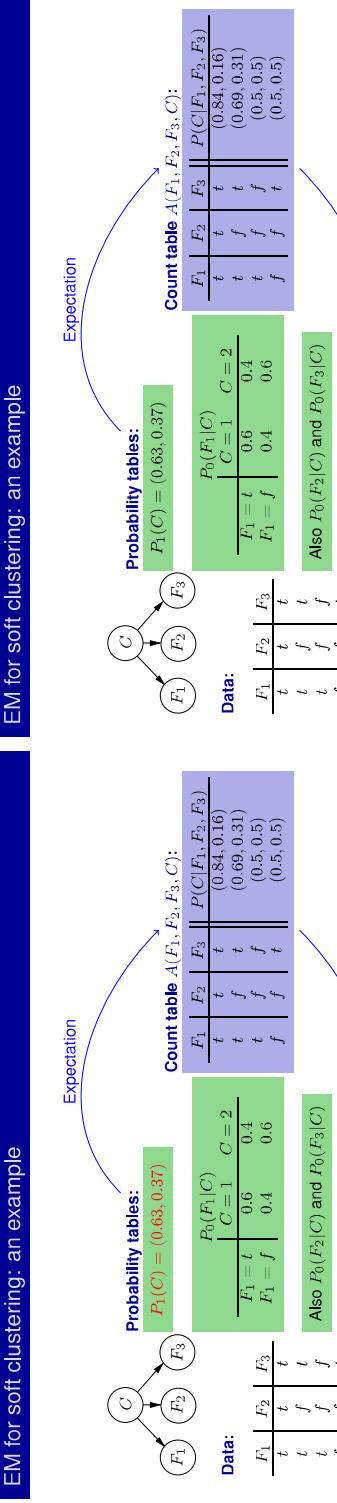
$F_1 = t$	$F_1 = f$	$C = 1$	$C = 2$
0.6	0.4	0.4	0.6

Also $P_0(F_2|C)$ and $P_0(F_3|C)$

Maximization

$$P_1(C) = \frac{1}{4} \sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C) = \frac{1}{4} (0.84 + 0.69 + 0.5 + 0.5, 0.16 + 0.31 + 0.5 + 0.5) \\ = (0.63, 0.37)$$

MI Autumn 2019 | Soft clustering | 20



Expectation

Probability tables:
 $P_1(C) = (0.63, 0.37)$

Data:

$F_1 = t$	$F_1 = f$	$C = 1$	$C = 2$
0.6	0.4	0.4	0.6

Also $P_0(F_2|C)$ and $P_0(F_3|C)$

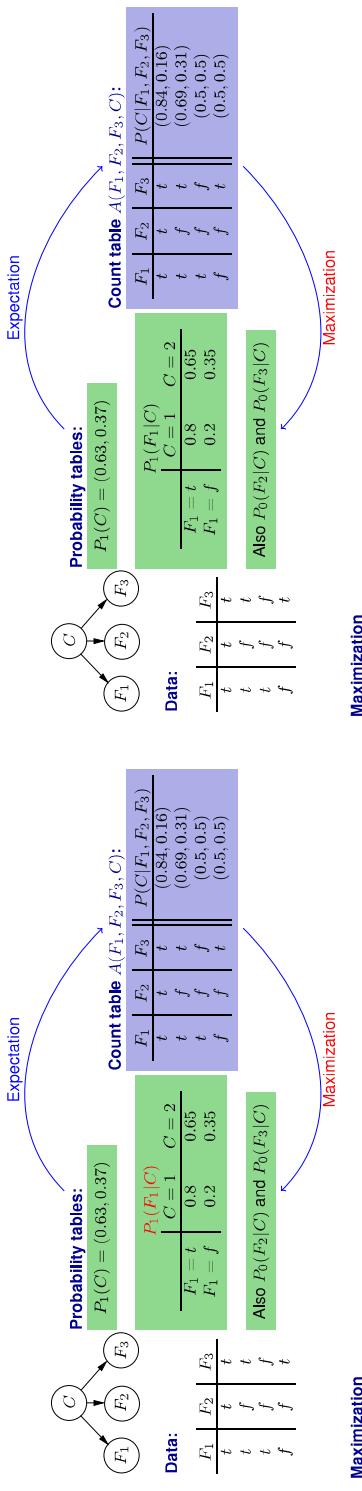
Maximization

$$P_1(F_1|C) = \frac{\sum_{F_2, F_3} A(F_1, F_2, F_3, C)}{\sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C)} = \frac{(0.84 + 0.69 + 0.5 + 0, 0.16 + 0.31 + 0.5 + 0)}{(2.53, 1.47)} \\ = (0.2, 0.35)$$

MI Autumn 2019 | Soft clustering | 20

EM for soft clustering: an example

EM for soft clustering: an example



$$P_1(F_1|C) = \frac{\sum_{F_2, F_3} P_1(F_1, F_2, F_3, C)}{\sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C)} = \frac{(0.84 + 0.69 + 0.5 + 0 \quad 0.16 + 0.31 + 0.5 + 0)}{(2.53, 1.47)}$$

$$= \begin{pmatrix} 0.8 & 0.65 \\ 0.2 & 0.35 \end{pmatrix}$$

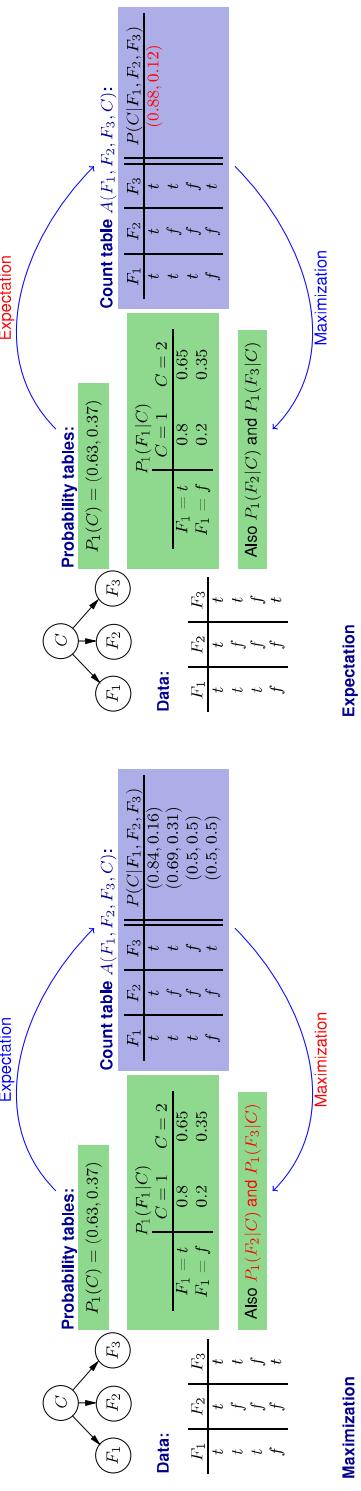
Mi Autumn 2019

Soft clustering

20

EM for soft clustering: an example

EM for soft clustering: an example



$$P_1(F_2|C) = \dots = \begin{pmatrix} 0.33 & 0.11 \\ 0.67 & 0.89 \end{pmatrix}$$

$$P_1(F_3|C) = \dots = \begin{pmatrix} 0.80 & 0.66 \\ 0.20 & 0.34 \end{pmatrix}$$

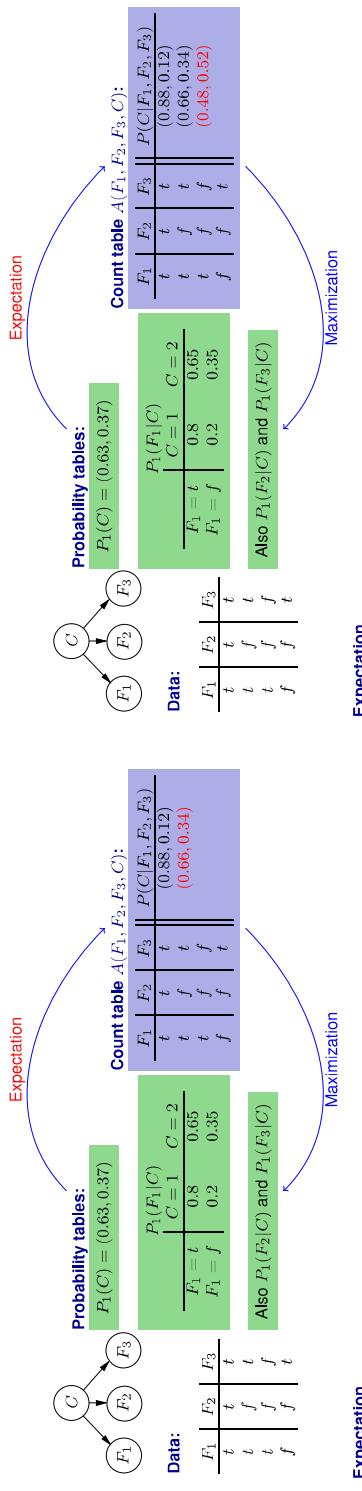
Mi Autumn 2019

Soft clustering

20

EM for soft clustering: an example

EM for soft clustering: an example



$$P_1(F_2|C) = \dots = \begin{pmatrix} 0.8 & 0.65 \\ 0.2 & 0.35 \end{pmatrix}$$

$$P_1(F_3|C) = \dots = \begin{pmatrix} 0.48 & 0.52 \\ 0.52 & 0.48 \end{pmatrix}$$

Mi Autumn 2019

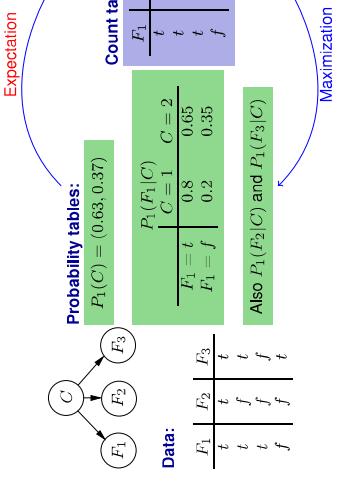
Soft clustering

20

- Fractional counts are being calculated by probability updating.

EM for soft clustering: an example

EM for soft clustering: an example



• Fractional counts are being calculated by probability updating.

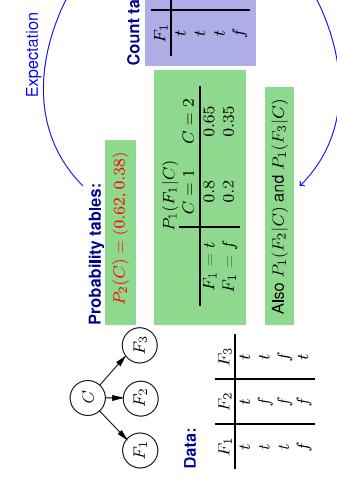
Maximization

$$P_2(C) = \frac{1}{4} \sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C) = \frac{1}{4} (0.88 + 0.66 + 0.48 + 0.47, 0.12 + 0.34 + 0.52 + 0.53) \\ = (0.62, 0.38)$$

MI Autumn 2019 | Seite 20 | Soft clustering

EM for soft clustering: an example

EM for soft clustering: an example

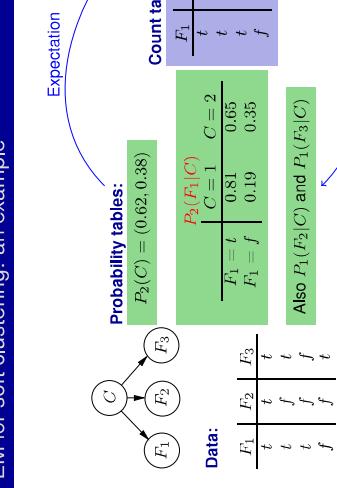


$$P_2(C) = \frac{1}{4} \sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C) = \frac{1}{4} (0.88 + 0.66 + 0.48 + 0.47, 0.12 + 0.34 + 0.52 + 0.53) \\ = (0.62, 0.38)$$

MI Autumn 2019 | Seite 20 | Soft clustering

EM for soft clustering: an example

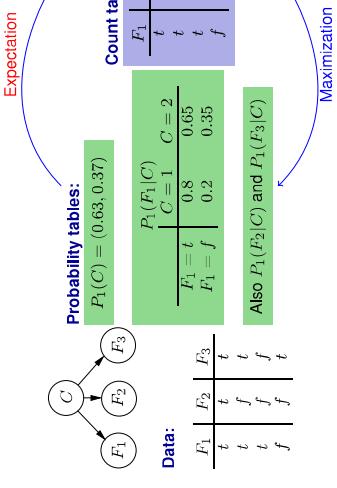
EM for soft clustering: an example



$$P_2(C) = \frac{\sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C)}{\sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C)} = \frac{(0.88 + 0.66 + 0.48 + 0.47, 0.12 + 0.34 + 0.52 + 0)}{(2.49, 1.51)}$$

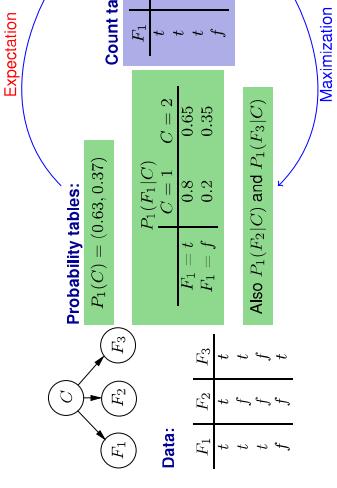
...and so we continue until a termination criterion is reached.

MI Autumn 2019 | Seite 20 | Soft clustering



$$P_2(C) = \frac{1}{4} \sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C) = \frac{1}{4} (0.88 + 0.66 + 0.48 + 0.47, 0.12 + 0.34 + 0.52 + 0.53)$$

MI Autumn 2019 | Seite 20 | Soft clustering



$$P_2(C) = \frac{1}{4} \sum_{F_1, F_2, F_3} A(F_1, F_2, F_3, C) = \frac{1}{4} (0.88 + 0.66 + 0.48 + 0.47, 0.12 + 0.34 + 0.52 + 0.53)$$

MI Autumn 2019 | Seite 20 | Soft clustering

Correctness

- The sequence of probability estimates generated by the EM algorithm converges to a local maximum (in rare cases: a saddle point) of the marginal likelihood given the data.
- To avoid sub-optimal local maxima: run EM several times with different starting points.

Correctness

- The sequence of probability estimates generated by the EM algorithm converges to a local maximum (in rare cases: a saddle point) of the marginal likelihood given the data.
- To avoid sub-optimal local maxima: run EM several times with different starting points.

Notes

- Any permutation of the cluster labels of a local maximum will also be a local maximum.
- Rather than keeping track of a full count table, it suffices to store counts for the variable families, $\text{pa}(X) = \{X\} \cup \text{pa}(X)$. Only one pass through the data is necessary.
- Clustering an existing or new instance x amounts to calculating $P(C|x)$.

Notes

MI Autumn 2019
Soft clustering
21

Cluster evaluation**Cluster evaluation**

A clustering algorithm applied to a dataset will return a clustering - even if there is no meaningful structure in the data!

Question: Do the clusters actually correspond to meaningful groups of data instances?

Question: Are all the clusters relevant, or are there some irrelevant and some meaningless clusters?

MI Autumn 2019
Cluster evaluation
22

MI Autumn 2019
Cluster evaluation
22

Unsupervised

- Uses only the data as given to the clustering algorithm, and the resulting clustering
- The realistic scenario
- Can be guided by considering changes in evaluation score.

Supervised

- Thomas Dyre Nielsen
Aalborg University

Machine Intelligence
Lecture 1: Planning
Thomas Dyre Nielsen
Aalborg University

MI Autumn 2019
Cluster evaluation
23

MI Autumn 2019
Cluster evaluation
23

Preferences and Lotteries

Utilities

Values with certainty

- What do you prefer
- \$100 or \$1000000 ?
 - A 4 or 10 grade in the exam?

Lotteries

A **lottery** is a probability distribution over outcomes. E.g.

$[0.4 : \$100, 0.6 : -\$20]$

means: you win \$100 with probability 0.4, and loose \$20 with probability 0.6.

$[0.3 : 00, 0.5 : 7, 0.2 : 10]$

means: with probability 0.3 you get a 00, with probability 0.5 a 7, and with probability 0.2 a 10.

Values with uncertainty

What do you prefer

• $[1:\$1000000]$ or $[0.5:\$0, 0.5:\$2100000]$?

• $[0.4:00, 0.1:7, 0.5:10]$ or $[0.1:00, 0.8:7, 0.1:10]$?

MI Autumn 2019

Utility

Typically:
Thus: preferences between lotteries with "money outcomes" are not always determined by **expected monetary value**.

Utilities

Utility of Money

Typically:

$[1:\$1000000]$ is preferred over $[0.5:\$0, 0.5:\$2100000]$

Thus: preferences between lotteries with "money outcomes" are not always determined by **expected monetary value**.

Preferences from utilities

The following is a classical result:

If preferences between lotteries obey a certain set of plausible rules, then there exists an assignment of real numbers (utilities) to all outcomes, such that one lottery is preferred over another if and only if it has a higher **expected utility**.

Example

Outcomes:	Expected Utility		
	00	7	10
Utilities:	-5	5	10
Lottery 1:	0.4	0.1	0.5
Lottery 2:	0.1	0.8	0.1

MI Autumn 2019

Utility

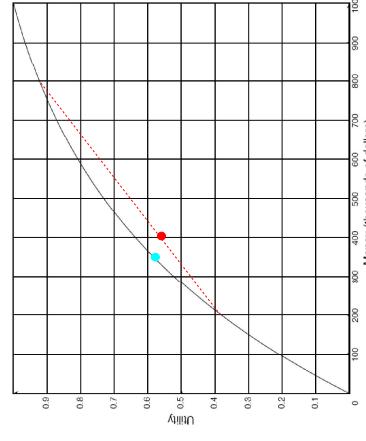
Typically:
Also "money outcomes" have a utility. Utility function of a **risk-averse** agent:



MI Autumn 2019

Utility

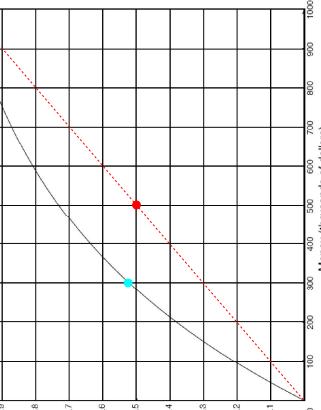
Typically:
Also "money outcomes" have a utility. Utility function of a **risk-averse** agent:



MI Autumn 2019

Utility

Typically:
Also "money outcomes" have a utility. Utility function of a **risk-averse** agent:



MI Autumn 2019

Utility

Typically:
Also "money outcomes" have a utility. Utility function of a **risk-averse** agent:

MI Autumn 2019

Utility

Typically:
Also "money outcomes" have a utility. Utility function of a **risk-averse** agent:

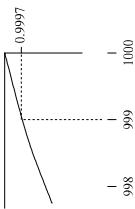
MI Autumn 2019

Utility

Typically:
Also "money outcomes" have a utility. Utility function of a **risk-averse** agent:

Digression: Insurance business

Assume Utility(\$ 999k)=0.9997:



Then agent is indifferent between lotteries

[1:\$ 999k] and [0.0003:0, 0.9997:\$ 1000k]

and prefers

[1:\$ 999.4k] over [0.0003:0, 0.9997:\$ 1000k]

Interpretation:

- right lottery: 0.03 risk of loosing a \$ 1000k property
- left lottery: buying insurance against that risk for \$ 600.

The insurance company prefers

[0.0003:0, 0.9997:\$ 1000k] over [1:\$ 999.4k]

\leadsto the insurance company has a different utility function (near linear).

Mi Autumn 2019

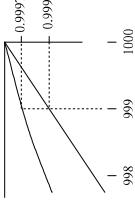
7

Utility

MI Autumn 2019 | Utility

Digression: Insurance business

Assume Utility(\$ 999k)=0.9997:



Then agent is indifferent between lotteries

[1:\$ 999k] and [0.0003:0, 0.9997:\$ 1000k]

and prefers

[1:\$ 999.4k] over [0.0003:0, 0.9997:\$ 1000k]

Interpretation:

- right lottery: 0.03 risk of loosing a \$ 1000k property
- left lottery: buying insurance against that risk for \$ 600.

The insurance company prefers

[0.0003:0, 0.9997:\$ 1000k] over [1:\$ 999.4k]

\leadsto the insurance company has a different utility function (near linear).

Mi Autumn 2019

7

Utility

MI Autumn 2019 | Utility

Do we maximize expected utility?

Recall:

- Lottery $A = [\$1\text{mill.},]$,
- Lottery $B = 0.1[\$5\text{mill.}] + 0.89[\$1\text{mill.}] + 0.01[\$0].$
- Lottery $C = 0.11[\$1\text{mill.}] + 0.89[\$0],$
- Lottery $D = 0.1[\$5\text{mill.}] + 0.9[\$0].$

A preferred over $B \Rightarrow$

$$\begin{aligned} U(\$1m) &> 0.1U(\$5m) + 0.89U(\$1m) + 0.01U(\$0) \\ \Rightarrow \quad & 0.1U(\$1m) > 0.1U(\$5m) + 0.01U(\$0) \end{aligned}$$

Mi Autumn 2019

8

Utility

MI Autumn 2019 | Utility

Do we maximize expected utility?

Recall:

- Lottery $A = [\$1\text{mill.},]$,
- Lottery $B = 0.1[\$5\text{mill.}] + 0.89[\$1\text{mill.}] + 0.01[\$0].$
- Lottery $C = 0.11[\$1\text{mill.}] + 0.89[\$0],$
- Lottery $D = 0.1[\$5\text{mill.}] + 0.9[\$0].$

A preferred over $B \Rightarrow$

$$\begin{aligned} U(\$1m) &> 0.1U(\$5m) + 0.89U(\$1m) + 0.01U(\$0) \\ \Rightarrow \quad & 0.1U(\$1m) > 0.1U(\$5m) + 0.01U(\$0) \end{aligned}$$

Mi Autumn 2019

8

Utility

MI Autumn 2019 | Utility

Do we maximize expected utility?

Recall:

- Lottery $A = [\$1\text{mill.},]$,
- Lottery $B = 0.1[\$5\text{mill.}] + 0.89[\$1\text{mill.}] + 0.01[\$0].$
- Lottery $C = 0.11[\$1\text{mill.}] + 0.89[\$0],$
- Lottery $D = 0.1[\$5\text{mill.}] + 0.9[\$0].$

A preferred over $B \Rightarrow$

$$\begin{aligned} U(\$1m) &> 0.1U(\$5m) + 0.89U(\$1m) + 0.01U(\$0) \\ \Rightarrow \quad & 0.1U(\$1m) > 0.1U(\$5m) + 0.01U(\$0) \end{aligned}$$

Mi Autumn 2019

8

Utility

MI Autumn 2019 | Utility

Do we maximize expected utility?

Recall:

- Lottery $A = [\$1\text{mill.},]$,
- Lottery $B = 0.1[\$5\text{mill.}] + 0.89[\$1\text{mill.}] + 0.01[\$0].$
- Lottery $C = 0.11[\$1\text{mill.}] + 0.89[\$0],$
- Lottery $D = 0.1[\$5\text{mill.}] + 0.9[\$0].$

A preferred over $B \Rightarrow$

$$\begin{aligned} U(\$1m) &> 0.1U(\$5m) + 0.89U(\$1m) + 0.01U(\$0) \\ \Rightarrow \quad & 0.1U(\$1m) > 0.1U(\$5m) + 0.01U(\$0) \end{aligned}$$

Mi Autumn 2019

8

Utility

MI Autumn 2019 | Utility

Factored Utility

So far: outcomes seen as unstructured states. When states are described by features, then overall utility often a combination of utility factors derived from different features.

Example

Two component utility function:

RHC	SWC	Utility
rhc	swc	5
rhc	swc	3
rhc	swc	0
rhc	swc	5
rhc	swc	2
cs	mw	0
off	mw	2
off	mw	2
lab	mw	0
lab	mw	2
mr	mw	5
mr	mw	2

Utility of full outcome (state) is sum of utility factors:

$$U((\text{off}, \text{rhc}, \overline{\text{swc}}, \overline{\text{mw}}, \text{rmw})) = 3 + 2 = 5$$

Assumption: the utility contribution from one factor is independent of the values of other factors.

E.g.: (rhc,swc) should probably be worth less than 5 when at the same time (mw,mrw), because mail needs to be delivered first (the two utility factors are **substitutes**).

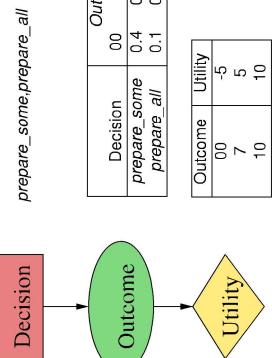
Utility

9

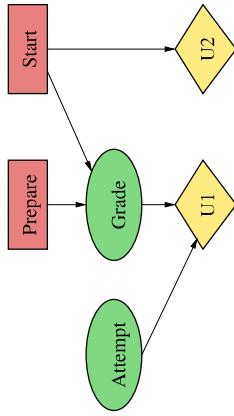
MI Autumn 2019 Single-Stage Decision Networks 10

Decisions and Lotteries

Simple decisions can be seen as choices over lotteries. Graphical representation of study example:



Decisions, outcomes and utilities can all be composed of features or factors:
Two components of decision: prepare *some/all*; start preparations sooner/later
Two utility factors: utility of grade, and utility (cost) of preparation time
Outcome composed of Grade, Attempt



Graph represents:

- One utility factor depends on Attempt and Grade, another only on Start
- Both the Prepare and Start decision influence the probabilities for Grade.

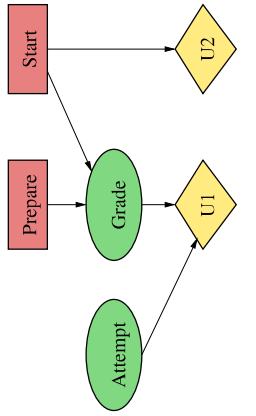
Utility

10

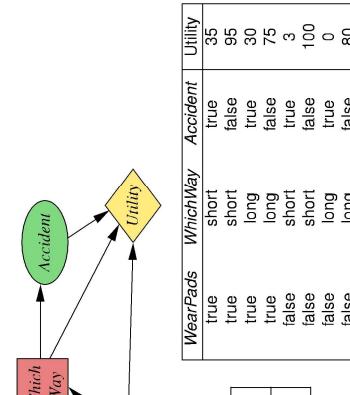
MI Autumn 2019 Single-Stage Decision Networks 11

Feature Based Models

Decisions, outcomes and utilities can all be composed of features or factors:
Two components of decision: prepare *some/all*; start preparations sooner/later
Two utility factors: utility of grade, and utility (cost) of preparation time
Outcome composed of Grade, Attempt



Robot Example



WearPads	WhichWay	Accident	Utility
true	short	true	35
true	short	false	95
true	long	true	30
true	long	false	75
false	short	true	3
false	short	false	100
false	long	true	0
false	long	false	80

WhichWay	Accident
short	0.2
long	0.99

A Single-Stage Decision Network is a directed acyclic graph with three types of nodes:

- Decision Nodes D
- Chance Nodes C
- Utility Nodes U

The graph must have the following structure:

- All decision nodes are connected in one linear sequence (representing the order in which the different sub-decisions are taken)
- The only parent of a decision node is its predecessor in the order
- Chance Nodes can have Decision Node and Chance Node parents
- Utility Nodes can have Decision Node and Chance Node parents

MI Autumn 2019 Single-Stage Decision Networks 12

MI Autumn 2019 Single-Stage Decision Networks 13

SSDN in General

SSDN Semantics and solutions

Structure

- Decision Nodes D
- Chance Nodes C
- Utility Nodes U

The graph must have the following structure:

- All decision nodes are connected in one linear sequence (representing the order in which the different sub-decisions are taken)
- The only parent of a decision node is its predecessor in the order
- Chance Nodes can have Decision Node and Chance Node parents
- Utility Nodes can have Decision Node and Chance Node parents

Tables

- No table is associated with decision nodes (only the list of available decisions)
- A chance node is labeled with a conditional probability table that specifies for each value assignment to its parents (decision and chance nodes) a probability distribution over the domain of the chance node.
- A utility node is labeled with a utility table that specifies for each value assignment to its parents (decision and chance nodes) a utility value.

Mi Autumn 2019

13

Single-Stage Decision Networks

$P(\omega \mid \mathbf{D} = \mathbf{d})$

A possible world ω is an assignment of values to all decision and chance variables.

An SSDN defines:

- For each assignment $\mathbf{D} = \mathbf{d}$ of values to the decision nodes a probability distribution over possible worlds.
- For each possible world ω a utility value

$U(\omega)$

An SSDN defines:

- For each assignment $\mathbf{D} = \mathbf{d}$ of values to the decision nodes a probability distribution over possible worlds.

$P(\omega \mid \mathbf{D} = \mathbf{d})$

over possible worlds.

For each possible world ω a utility value

$U(\omega)$

SSDN Semantics



The graph must have the following structure:

- All decision nodes are connected in one linear sequence (representing the order in which the different sub-decisions are taken)
- The only parent of a decision node is its predecessor in the order
- Chance Nodes can have Decision Node and Chance Node parents
- Utility Nodes can have Decision Node and Chance Node parents

Tables

- No table is associated with decision nodes (only the list of available decisions)
- A chance node is labeled with a conditional probability table that specifies for each value assignment to its parents (decision and chance nodes) a probability distribution over the domain of the chance node.
- A utility node is labeled with a utility table that specifies for each value assignment to its parents (decision and chance nodes) a utility value.

Mi Autumn 2019

13

Single-Stage Decision Networks

$P(\omega \mid \mathbf{D} = \mathbf{d})$

SSDN Semantics and solutions

Robot Example

SSDN Semantics

A possible world ω is an assignment of values to all decision and chance variables.

An SSDN defines:

- For each assignment $\mathbf{D} = \mathbf{d}$ of values to the decision nodes a probability distribution over possible worlds.
- For each possible world ω a utility value

$P(\omega \mid \mathbf{D} = \mathbf{d})$

over possible worlds.

For each possible world ω a utility value

$U(\omega)$

Solving an SSDN

To solve a single-stage decision network (or problem) means to find the decisions \mathbf{d} that maximize the expected utility

$$\mathcal{E}(U \mid \mathbf{D} = \mathbf{d}) = \sum_{\omega} U(\omega) P(\omega \mid \mathbf{D} = \mathbf{d})$$

Mi Autumn 2019

14

Single-Stage Decision Networks

$U(\omega)$

Robot Example

Robot Example

15

Single-Stage Decision Networks

$U(\omega)$

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

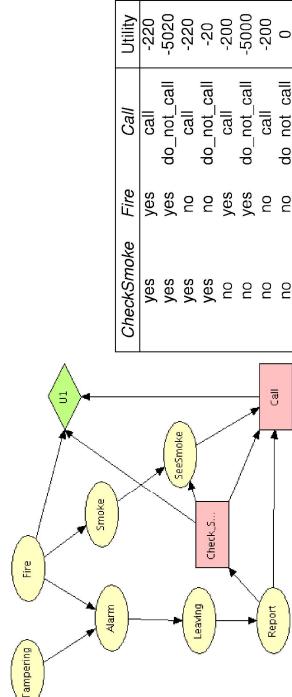
WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
true	short	false	0.2 · 35
true	short	true	0.8 · 95
true	long	false	0.01 · 30
true	long	true	0.99 · 75
false	short	false	0.2 · 3
false	short	true	0.8 · 100
false	long	false	0.01 · 0
false	long	true	0.99 · 80

WearPads	WhichWay	Accident	$P(\omega \mid \mathbf{D} = \mathbf{d}) \cdot U(\omega)$
<

Example: Fire scenario

Fire alarm example extended with

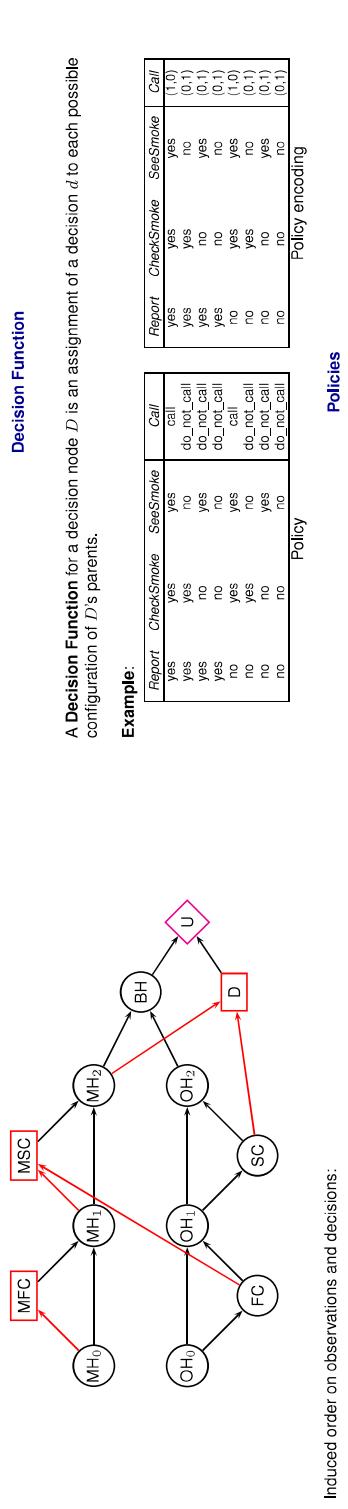
- two decisions: $CheckSmoke \in \{\text{yes}, \text{no}\}$, $Call \in \{\text{call}, \text{do_not_call}\}$
- one more random variable $SeeSmoke \in \{\text{yes}, \text{no}\}$
- a utility function depending on $Fire, CheckSmoke, Call$ (composed of two factors: one depending on $CheckSmoke$, one on $Fire$ and $Call$)



We assume that the decision maker doesn't forget: the *no-forgetting assumption*

Mi Autumn 2019 | 18 | Sequential Decisions

Possible decision model

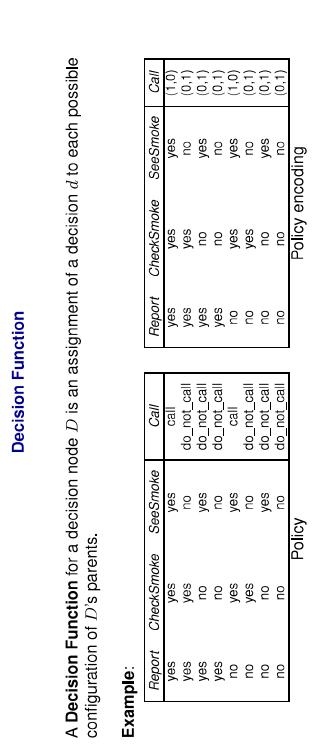


Induced order on observations and decisions:

$$\{MH_0\} \prec \textcolor{red}{MFC} \prec \{MH_1, OFC\} \prec \textcolor{red}{MSC} \prec \{MH_2, OSC\} \prec \textcolor{red}{D} \prec \{OH_0, OH_1, OH_2, BH\}$$

Mi Autumn 2019 | 20 | Sequential Decisions

Policies



- Example:**
- A **Policy** π consists of one decision function for each decision node.
- General strategy for actions (decisions), taking into account the possible (uncertain) effects of previous actions

Mi Autumn 2019 | 21 | Sequential Decisions

Policies

Semantics

Expected Utility

As before: possible worlds ω are assignments for all decision and chance variables.

- A policy π defines a probability distribution

$$P(\omega \mid \pi)$$

over possible worlds:

- if ω contains assignments to a decision node D and its parents which is not consistent with the decision function for D : $P(\omega \mid \pi) = 0$.
- Otherwise: $P(\omega \mid \pi)$ is the product of all conditional probability values for the assignments to chance nodes C , given the assignment to the parents of C .
- Each possible world has a utility
- Obtain expected utility of a policy

$$\mathcal{E}(U \mid \pi) = \sum_{\omega} U(\omega) P(\omega \mid \pi)$$

Optimal Policy

An **optimal policy** is a policy with maximal expected utility (among all possible policies).

Mi Autumn 2019 | 22 | Sequential Decisions

Mi Autumn 2019 | 23 | Sequential Decisions

1. $DFs = \emptyset$ // Set of decision functions
2. Fs = all conditional probability and utility tables
3. **while** there are decision nodes
 - sum out all random variables that are not parents of a decision node
 - // Fs now contains a factor F that depends on one decision node D
 - // and (a subset of) its parents!
4. Add $\max_{\mathcal{V}_D} F$ to Fs
5. Add $\arg \max_{\mathcal{V}_D} F$ to DFs
6. Sum out remaining random variables
7. Sum out DFs and product of remaining factors (expected utility of optimal policy)

Value of Information

MI Autumn 2019
Sequential Decisions
25
Value of Information

Test Decisions

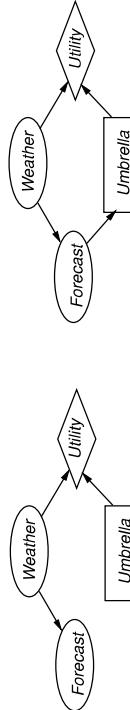
Collecting Information

- *CheckSmoke* is aimed at determining (with some uncertainty) the true state of *Smoke* (or even *Fire*)
- *Test* (medical example) is aimed at determining the true state of *Disease*

Value of Information

Question: what is it worth to know the exact state of a random variable C when making decision D ?

Answer: Compare maximal expected utilities of



MI Autumn 2019
Value of Information
26
Value of Information

Properties of Value of Information

- Value of information is always non-negative
- Value of knowing C for decision D is zero, if no observed value of C can change the decision rule, i.e., for all values p of existing parents of D , and all values c of C , the optimal decision given (p, c) is the same as the optimal decision given (p) .

Machine Intelligence

Lecture 12: Multi-agent systems

Thomas Dyre Nielsen

Aalborg University

Value of Information

1

MI Autumn 2019
Value of Information
27
Value of Information

Tentative course overview

Topics:

- Introduction
- Search-based methods
- Constrained satisfaction problems
- Logic-based knowledge representation
- Representing domains endowed with uncertainty
- Bayesian networks
- Inference in Bayesian networks
- Machine learning: classification
- Machine learning: clustering
- Planning
- **Multi-agent systems**

Multi-Agent Systems

MI Autumn 2019

2

MI-Agent Systems

3

From Single to Multi Agent Systems

Game Trees

So far ...

We have modeled an agent that decides/plans in a world with/without uncertainty.

New Dimension

Agent acts in an environment containing other agents. Other agents might have competing/conflicting objectives.

Using Uncertainty

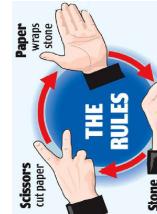
The actions of other agents can partly be represented as uncertainty in effects of own actions (uncertainty of state transitions).

Better: take explicitly into account

- Competing objectives of other agents
- Reasoning about what other agents will do (reduce uncertainty)
- Possibility to collaborate to achieve common objectives

Imperfect Information Games

Representation of game with simultaneous moves:



Collect in an information set the nodes that the agent (Bob) can not distinguish (at all nodes in an information set the same actions must be possible).

Other sources for imperfect information:

- Unobserved: random moves by nature (dealing of cards).
- Hidden moves by other agent

Strategies

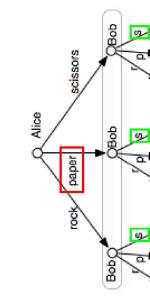
4

MI Autumn 2019

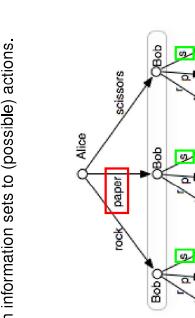
5

A (pure) strategy for one agent is a mapping from information sets to (possible) actions.

Example strategies for A and strategies for B:



(A strategy is essentially the same as a policy)
A strategy profile consists of a strategy for each agent.



(A strategy is essentially the same as a policy)
A strategy profile consists of a strategy for each agent.

MI Autumn 2019

6

MI-Agent Systems

7

MI Autumn 2019

Utility

Solving Perfect Information Gain

- Utility for each agent given a strategy profile:
- each node has the utilities that will be reached at a leaf by following the strategy profile
- the utilities at the nodes represent the outcome of the game (given the strategy profile)
- utilities at a *node* are computed by taking the expectation over the utilities of its successors)

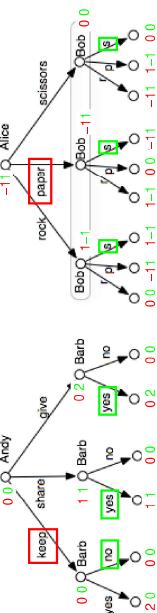
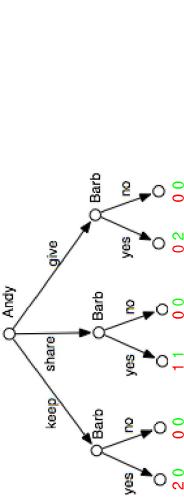


Figure shows the utilities for A and utilities for B at all nodes.

- game is perfect information (no information sets with more than 1 node)
- both agents play rationally (optimize their own utility)

then the optimal strategies for both players are determined by

- bottom-up propagation of utilities under optimal strategies, where
- each player selects the action that leads to the child with the highest utility



○ 00
○ 02
○ 00
○ 11
○ 00
○ 20

Solving Perfect Information Gain

Solving Perfect Information Gain

- game is perfect information (no information sets with more than 1 node)
- both agents play rationally (optimize their own utility)

then the optimal strategies for both players are determined by

- bottom-up propagation of utilities under optimal strategies, where

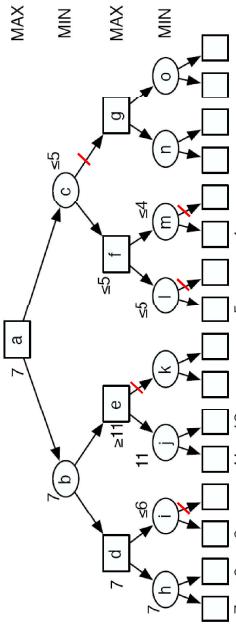


Zero Sum Games

In this case:

- need only one utility value at the leaves
- one player (called *Max*) wants to reach leaf with maximal value
- one player (called *Min*) wants to reach leaf with minimal value

- Schaeffer et al. proved: there is no winning strategy for either player: perfect play by both players will always result in a draw
- checkers has approximately $5 \cdot 10^{20}$ different positions
- in the proof only about 10^{14} positions were explored



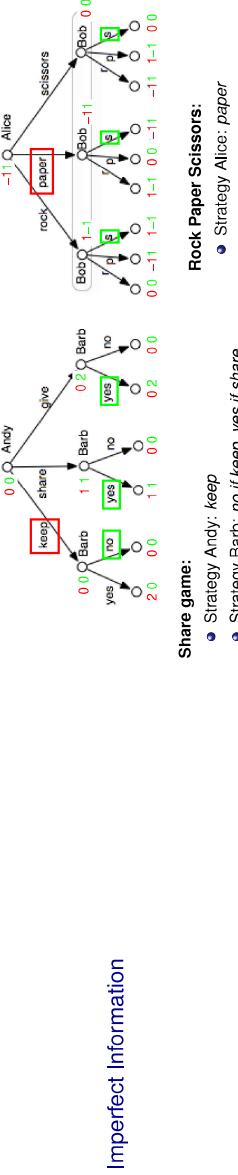
A wooden checkers board with black and orange pieces.

Journal of Oral Rehabilitation 2003; 30: 103–109

10

Normal Form

For each strategy profile, utilities of game are determined:



Share game:

- Strategy Andy: keep
- Strategy Barb: no if keep, yes if share, yes // give
- Utilities: 0 for Andy, 0 for Barb
- Utilities determined by these choices

MI Autumn 2019 11 Imperfect Information

Normal form representations

Share game

		Andy		
		keep	share	give
Barb		2, 0	1, 1	0, 2
$k \rightarrow y, s \rightarrow y, g \rightarrow y$		2, 0	1, 1	0, 2
$k \rightarrow y, s \rightarrow y, g \rightarrow n$		2, 0	0, 0	0, 2
$k \rightarrow y, s \rightarrow n, g \rightarrow y$		2, 0	0, 0	0, 2
$k \rightarrow y, s \rightarrow n, g \rightarrow n$		0, 0	1, 1	0, 2
$k \rightarrow n, s \rightarrow y, g \rightarrow y$		0, 0	1, 1	0, 2
$k \rightarrow n, s \rightarrow y, g \rightarrow n$		0, 0	1, 1	0, 2
$k \rightarrow n, s \rightarrow n, g \rightarrow y$		0, 0	0, 0	0, 2
$k \rightarrow n, s \rightarrow n, g \rightarrow n$		0, 0	0, 0	0, 2

Rock Paper Scissors

		Alice		
		rock	paper	scissors
Bob		0, 0	1, -1	-1, 1
rock		-1, 1	0, 0	1, -1
paper		1, -1	-1, 1	0, 0
scissors		1, -1	-1, 1	0, 0

Difference between perfect and imperfect information not directly visible in normal form representation!

MI Autumn 2019 12 Imperfect Information

Prisoner's Dilemma

No pure strategy Nash equilibrium in Rock Paper Scissors:

		Alice		
		rock	paper	scissors
Bob		0, 0	-1, 1	1, -1
rock		-1, 1	0, 0	1, -1
paper		1, -1	-1, 1	0, 0
scissors		1, -1	-1, 1	0, 0

A mixed strategy is a probability distribution over actions.

$$\begin{aligned} \text{Mixed Strategy for Alice: } r : 1/3 p : 1/3 s : 1/3 \\ \text{Mixed Strategy for Bob: } r : 1/3 p : 1/3 s : 1/3 \end{aligned}$$

Expected utility for Alice = expected utility for Bob =

$$1/9(0 + 1 - 1 - 1 + 0 + 1 + 1 - 1 + 0) = 0$$

The only Nash equilibrium is Alice:testify, Bob:testify

Nash equilibria do not represent cooperative behavior!

MI Autumn 2019 14 Imperfect Information

MI Autumn 2019 15 Imperfect Information

Mixed Strategies

Key Results

No pure strategy Nash equilibrium in Rock-Paper-Scissors:

		Alice		
		Rock	Paper	Scissors
Bob	Rock	0	-1	-1
	Paper	-1	0	1
Scissors	1	-1	0	0
	1	1	-1	0

A mixed strategy is a probability distribution over actions.

Mixed Strategy for Alice: $r : 1/3 p : 1/3 s : 1/3$

Mixed Strategy for Bob: $r : 1/3 p : 1/3 s : 1/3$

Expected utility for Alice = expected utility for Bob =

$$1/9(0 + 1 - 1 - 1 + 0 + 1 + 1 - 1 + 0) = 0$$

Suppose Alice plays some other strategy: $r : p_r : p_s : p_p : p_s$. Expected utility for Alice then:

$$1/3(p_r \cdot 0 + p_p \cdot 1 - p_s \cdot 1 - p_r \cdot 1 + p_p \cdot 0 + p_s \cdot 1p_r + 1 - p_p \cdot 1 + p_s \cdot 0) =$$

If Bob plays $r : 1/3 p : 1/3 s : 1/3$, Alice can not do better than playing $r : 1/3 p : 1/3 s : 1/3$ also.

Same for Bob

Both playing $r : 1/3 p : 1/3 s : 1/3$ is a (the only) Nash equilibrium

MI Autumn 2019

Imprecise Information

- Every (finite) game has a Nash equilibrium (using mixed strategies)
- There can be multiple Nash equilibria
- Playing a Nash equilibrium strategy profile does not necessarily lead to optimal utilities for the agents (prisoner's dilemma)

The exam

Some practical issues

- January 8th, 2019.
- Written exam with internal censor.
- Graded exam.
- Answers should be written in English.
- A "question session" is scheduled for ??.

MI Autumn 2019

The Exam

17

MI Autumn 2019

The Exam

17

What the course has covered

The course has covered the following issues:

- Introduction
- Bayesian networks
- Search-based methods
- Inference in Bayesian networks
- Constrained satisfaction problems
- Machine learning
- Logic-based knowledge representation
- Planning
- Reasoning under uncertainty
- Multi-agent systems

This corresponds to the following literature:

- David L. Poole and Alan K. Mackworth, Artificial Intelligence: Foundations of computational agents (Second edition): Preface, Ch. 1, 3-3.6, 3.7-7.1, 3.7.3, 3.8.2, 3.8.3, 4-4.7.3, 5-5.2, 7-7.5 (except 7.4.2), 7.7, 8-8.4.1, 8.6-8.5, 9-9.4 (except 9.1.3), 10.1.2, 10.2, 11-11.4, A.3
- Finn V. Jensen and Thomas D. Nielsen, Bayesian networks and decision graphs: Sections 2-2.2, 3-3.1.
- The slides from the course.

MI Autumn 2019

The Exam

18