

# TDT4171 - Artificial Intelligence Methods

Martin Bjerke

June 2, 2018

# Contents

<b>1</b>	<b>Book</b>	<b>2</b>
13	Quantifying Uncertainty . . . . .	2
14	Bayesian networks . . . . .	4
15	Probabilistic reasoning over time . . . . .	6
16	Making simple decisions . . . . .	9
17	Making complex decisions . . . . .	12
18	Learning from examples . . . . .	13
21	Reinforcement learning . . . . .	16
22	Philosophical Foundations . . . . .	18
23	AI: The Present and Future . . . . .	19
<b>2</b>	<b>CBR paper</b>	<b>20</b>
1	Background . . . . .	20
2	History . . . . .	20
3	Fundamentals of case-based reasoning methods . . . . .	20
4	Representation of cases/Knowledge representation . . . . .	21
4.1	The Dynamic Memory Model . . . . .	22
4.2	Category & Exemplar Model . . . . .	22
5	Case retrieval . . . . .	23
6	Case Reuse . . . . .	24
7	Case Revision . . . . .	24
8	Case Retainment - Learning . . . . .	24
<b>3</b>	<b>Deep learning</b>	<b>25</b>
1	Supervised learning . . . . .	25
2	Backpropagation . . . . .	25
3	Convolutional neural networks . . . . .	26
4	Distributed representations and language processing . . . . .	26
5	Recurrent neural networks . . . . .	26

In the following summary of the syllabus of TDT4171 in 2018, text that has an additional margin is only mentioned in the syllabus, not in the lectures. The sections and equations have the same numbering as in the book.

# 1 Book

## 13 Quantifying Uncertainty

Agent encounters uncertainty through partial observability, non-determinism or both. To counter this problem-solving agents and logic agents track the world state and store it as a *belief state*. The *qualification problem* arises when an agent tries to store all possibilities (even implausible ones) in its belief state, making it impossible to infer its next action. Therefore *rational decision* takes into consideration both the relative importance of an action and its likelihood.

Using inference to make an action fails on the following ground:

*Laziness* failure to enumerate exceptions, qualifications etc

*Theoretical ignorance* no theory for the domain

*Practical ignorance* not possible to run all tests

Probabilistic assertions can summarise these effects, and an agent can store its *degrees of belief* that something will happen and use *probability theory* to choose the best action. The combination of probability theory and *utility theory* (inference) is called *decision theory* and all of them can be used to choose an action.

Probability gives an assertion about how the world looks (there is a 50% chance of rain today versus it is raining today), and the set of all possible worlds (denoted  $\Omega$ ) is called the sample space. A fully specified probability model gives a probability,  $P(\omega)$ , of each possible world, and by the axioms of probability each world has a probability between 0 and 1 and together the probability must sum to 1:

$$0 \leq P(\omega) \leq 1 \text{ for every } \omega \text{ and } \sum_{\omega \in \Omega} P(\omega) = 1 \quad (13.1)$$

The probability of *propositions* (also called probabilistic assertions or queries, or events) are the sum of the probability of the worlds where the proposition holds:

$$\text{For any proposition } \phi, P(\phi) = \sum_{\omega \in \phi} P(\omega) \quad (13.2)$$

There are different “types” of probability:

**Prior/unconditional probabilities:** Independent probabilities

**Posterior/conditional probabilities:** Dependent probabilities, the probability of rain, given that it's already cloudy. Denoted  $P(a \mid b)$ .

**Probability distributions:** Outputs a probability given a value, most common in settings where the sample space is continuous, like temperature.

### Joint probability distributions: Normalised probability distributions

As more knowledge is gained, conditions arises, like you might not know you have a cavity, but when you get toothache, you might consider it,  $P(\text{Cavity}|\text{Toothache})$ . Note, this means you **ONLY** know that you have a toothache, not anything more relevant information. The conditional probability can be calculated using unconditional probabilities:

$$P(a | b) = \frac{P(a \wedge b)}{P(b)} \text{ or usually presented as the } \textit{product rule} P(a \wedge b) = P(a | b)P(b) \quad (13.3)$$

Some notation, *random variables* are denoted with uppercase names and each variable has a *domain* of possible values, like  $\text{Die} = \{1, 2, 3, 4, 5, 6\}$ , boolean variables has two values  $\{\text{true}, \text{false}\}$  and values with names are denoted with lowercase names. When a bold  $\mathbf{P}$  is used it the probability of all the possible values of the variable,  $\mathbf{P}(\text{Weather}) = \langle 0.6, 0.1, 0.29, 0.01 \rangle^1$  and this is also called a *probability distribution*. With discrete variables the probability distribution can be represented as a *joint probability distribution* as shown below. Going back to the possible worlds talk about above, if each random variable of a world is given a specific value this represents one specific world. If the joint probability distribution contains all possible random variables, then it is *full joint probability distribution*(FJPD).

	toothache		$\neg$ toothache	
	catch	$\neg$ catch	catch	$\neg$ catch
cavity	.108	.012	.072	.008
$\neg$ cavity	.016	.064	.144	.576

Using this table each possible value of each variable can be calculated using probabilistic inference: ex.  $P(\text{toothache}) = .108 + .012 + .016 + .64 = .2$  and  $P(\text{cavity} \vee \text{toothache}) = .108 + .012 + .016 + .64 + .072 + .008 = .28$ .

Together with the product rule (Eq. ??) and the *inclusion-exclusion principle* (Eq. 13.4) we get the *Kolmogorov's axioms*. These build up the rest of probability theory.

$$P(a \vee b) = P(a) + P(b) - P(a \wedge b) \quad (13.4)$$

Summing up all the entries of a FJPD of a specific value, i.e  $P(\text{cavity}) = .108 + .012 + .072 + .008 = 0.2$ , is called *marginalisation* and can be represented as the following general rule, with the variables  $\mathbf{Y}$  and  $\mathbf{Z}$ :

$$\mathbf{P}(\mathbf{Y}) = \sum_{z \in \mathbf{Z}} \mathbf{P}(\mathbf{Y}, z) \quad (13.5)$$

Above this is partially done with *Cavity*,  $\mathbf{P}(\text{Cavity}) = \sum_{z \in \{\text{Catch}, \text{Toothache}\}} \mathbf{P}(\text{Cavity}, z)$ . From Eq. ?? we get the conditioning rule:

$$\mathbf{P}(\mathbf{Y}) = \sum_z \mathbf{P}(\mathbf{Y} | z)P(z) \quad (13.8)$$

Going back to Eq. 13.3 and if you calculate both  $P(\text{cavity} | \text{toothache}) = 0.6$  and  $P(\neg\text{cavity} | \text{toothache}) = 0.4$  these add up to 1, as they should. Also notice that the  $\frac{1}{P(\text{toothache})}$  is constant, and therefore we view it as a *normalisation* constant of the distribution of  $\mathbf{P}(\text{Cavity} | \text{toothache})$  and is therefore swapped with a normalisation variable  $\alpha$ , which is set so that a distribution  $\mathbf{P}(A)$

---

<sup>1</sup>Meaning here that  $P(\text{Weather} = \text{sunny}) = 0.6$ ,  $P(\text{Weather} = \text{rain}) = 0.1$ ,  $P(\text{Weather} = \text{cloudy}) = 0.29$ , and  $P(\text{Weather} = \text{snow}) = 0.01$

has a total probability of 1. This is done since  $\frac{1}{P(\text{toothache})}$  is just a normalisation, so we don't actually need to know it. From this we extract a general inference procedure:

Let  $X$  be all the query variables,  $E$  be the evidence variables and  $e$  be the specific values and let  $Y$  be the remaining hidden variables. This gives the following formula for summing the joint entries:

$$P(X | e) = \alpha P(X, e) = \alpha \sum_y P(X, e, y) \quad (13.9)$$

In more cases, the use of conditional independence reduces the size of the representation of the joint distribution from exponential in  $n$  to linear. For example  $n$  coin flips, if they were dependent there would be  $2^n$  representations in the joint distribution, but since they are independent, only  $n$  is needed.

Identifying independence, two variables are independent iff:

$$P(X | Y) = P(X) \text{ equivalently } P(Y | X) = P(Y)$$

Bayes rule:

$$P(Y | X) = \frac{P(X | Y)P(Y)}{P(X)}$$

or with some background evidence  $e$

$$P(Y | X, e) = \frac{P(X | Y, e)P(Y, e)}{P(X, e)} \quad (13.13)$$

or by using the normalisation trick from before

$$\begin{aligned} P(X | e) &= \alpha \langle P(e | x)P(x), P(e | \neg x)P(\neg x) \rangle \\ &= P(X | Y) = \alpha P(Y | X)P(X) \end{aligned}$$

Variables are said to be conditionally independent if they have the same *cause*, but different *effects*, where the effects are independent, but they are both dependent on the cause. I.e. Toothache and catch are both an effect of having a cavity, but having a toothache is independent of the toothpick to catch and vice versa;  $P(\text{toothache} \wedge \text{catch} | \text{Cavity}) = P(\text{toothache} | \text{Cavity})P(\text{catch} | \text{Cavity})$ . Conditional independence is represented in the following equation:

$$\begin{aligned} P(X, Y | Z) &= P(X | Z)P(Y | Z) \\ &\text{since } Y \text{ and } X \text{ are independent, the following is also true} \\ P(X | Y, Z) &= P(X | Z) \\ P(Y | X, Z) &= P(Y | Z) \end{aligned} \quad (13.19)$$

Using conditional independence to “shorten” a probability distribution is called a *naive Bayes model*, this is because it is used in cases where the effect variables are not actually conditionally independent given the cause variable.

## 14 Bayesian networks

Bayesian networks are graphical notations for conditional independence assertions, giving a compact specification of full joint distributions. The syntax is a set of nodes (one per variable) that makes out a DAG and a conditional distribution for each node given its parents,  $P(X_i | \text{Parents}(X_i))$

All conditional distribution represented as a conditional probability table (CPT) gives the distribution over  $X_i$  for each combination of parent values.

Example: A CPT for boolean  $X_i$  with  $k$  boolean parents has  $2^k$  rows for the combinations of parent values. Since  $X_i$  is boolean,  $P(X_i = \text{true}) = 1 - P(X_i = \text{false})$ , so each row requires only one representative, so if each variable has maximum  $k$  parents, the complete network requires  $O(n \cdot 2^k)$  numbers.

## Representing a full joint distribution and constructing a Bayesian network

A generic entry in a joint distribution is just a conjunction of assignment of values to all variables,  $P(X_1 = x_1 \wedge X_2 = x_2 \wedge \dots \wedge X_n = x_n)$  which is written in shorthand as  $P(x_1, x_2, \dots, x_n)$ . Since  $x_i$  is just dependent of its parents this value is given by:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=0}^n P(x_i \mid \text{parents}(X_i)) \quad (14.1)$$

Example: Let B and E be parents of A, and J and M be children of A,

$$P(j \wedge m \wedge a \wedge \neg b \wedge \neg e) = P(j \mid a)P(m \mid a)P(a \mid \neg b, \neg e)P(\neg b)P(\neg e) \quad (14.2)$$

The chain rule can be extracted from this using the product rule:

$$\begin{aligned} P(x_1, \dots, x_n) &= P(x_n \mid x_{n-1}, \dots, x_1)P(x_{n-1}, \dots, x_1) \\ &= P(x_n \mid x_{n-1}, \dots, x_1)P(x_{n-1} \mid x_{n-2}, \dots, x_1)P(x_{n-2}, \dots, x_1) \\ &= \dots = \prod_{i=1}^n P(x_i \mid x_{i-1}, \dots, x_1) \end{aligned} \quad (14.3)$$

Using the chain rule in Eq. 14.1 and generalising:

$$P(X_i \mid X_1, \dots, X_n) = P(x_i \mid \text{Parents}(X_i)) \quad (14.4)$$

Local semantics: each node is conditionally independent of its non descendants given its parents.  
Markov blanket: Each node is conditionally independent of all other nodes then its parents, children and children's parents.

Construction of a Bayesian network:

1. Choose an ordering of variables  $X_1, \dots, X_n$
2. for  $i = 1$  to  $n$ 
  - (a) add  $X_i$  to the network
  - (b) select parents from  $X_1, \dots, X_{i-1}$  such that  $P(X_i \mid \text{Parents}(X_i)) = P(X_i \mid X_1, \dots, X_{i-1})$

Step 1 - Decide what to model

Step 2 - Define variables

Step 3 - Define the graphical structure that connects the variables(The qualitative part)

Step 4 - Fix parameters to specify each  $P(x_i \mid \text{pa}(x_i))$ (The quantitative part)

Step 5 - Verification

There are other approaches, other than probability, to handle uncertainty:

- Qualitative reasoning: I.e. default reasoning: A conclusion is sound as long as there are no better conclusions to draw.
- Rule-based: Add a “fudge factor” to each rule to accommodate uncertainty. There are 3 desirable properties from rule-based:
  - *Locality*; Given the rule  $A \implies B$ , B can be concluded by only looking at A, not the other evidence.
  - *Detachment*; Once a proof for a proposition is found it can be detached from its justification.
  - *Truth-functionality*; Truth of sentences can be computed by the truth of the components. This can’t be done in probability without independence.

There have been attempts to add uncertainty to RB by assigning a belief to the rules.

- Dempster-Shafer theory: Compute the probability that the evidence supports the proposition (use a *belief function*).
- Fuzzy logic: is fuzzy

## 15 Probabilistic reasoning over time

To handle the coming dynamic models they are viewed discretely using *time slices*, containing sets of observable and unobservable random variables. For notation  $\mathbf{X}_t$  denotes the set of state variables at time  $t$  which are unobservable and  $\mathbf{E}_t$  denotes the set of observable evidence variables, so the observations at time  $t$  are  $\mathbf{E}_t = e_t$ . The notation  $a:b$  denotes (inclusive) the interval from  $a$  to  $b$ , so  $\mathbf{X}_{a:b}$  are the set of variables  $\mathbf{X}_a, \dots, \mathbf{X}_b$ . The examples used are based on the umbrella world of Chapter 15. The notations of *transition model* is how the world evolves and *sensor model* is how the evidence variables get their values.

The *transition model* specifies the probability distribution over the latest state variables, given the previous states,  $P(\mathbf{X}_t \mid \mathbf{X}_{0:t-1})$ . Using this increases the number of previous states need over time. This is solved by making a *Markov assumption*, the next states only depend on a finite set of previous states. This is called a *Markov process* or *Markov chain*, the simplest form being the *first-order Markov process* where the current state only depends on the previous one. In other words, a future state is only dependent on the previous state and conditionally independent of the rest:

$$P(\mathbf{X}_t \mid \mathbf{X}_{0:t-1}) = P(\mathbf{X}_t \mid \mathbf{X}_{t-1}) \quad (15.1)$$

Having bounded the number of transition models using a Markov assumption, we want to bound the possible number of different transition models that arises. We want a *stationary process*, where the rules of changes to the world states does not change. I.e in the umbrella world  $P(R_t \mid R_{t-1})$  is the same for all  $t$ .

Next is bounding the sensor model. This is done by making a *sensor Markov assumption* that works the same way as a Markov assumption, limiting the number of past evidence values needed and giving us the following sensor model (sometimes called the *observation model*):

$$P(\mathbf{E}_t \mid \mathbf{X}_{0:t}, \mathbf{E}_{0:t-1}) = P(\mathbf{E}_t \mid \mathbf{X}_t) \quad (15.2)$$

Lastly we need a prior probability distribution at time  $t = 0$ ,  $P(\mathbf{X}_0)$ . Given these assumptions and equations we have the complete joint distribution for any  $t$ :

$$P(\mathbf{X}_{0:t}, \mathbf{E}_{1:t}) = P(\mathbf{X}_0) \prod_{i=1}^t P(\mathbf{X}_i | \mathbf{X}_{i-1}) P(\mathbf{E}_i | \mathbf{X}_i) \quad (15.3)$$

To use these models, the following inference tasks can be used:

- *Filtering*: Calculate the posterior distribution, aka the belief state, given the evidence observed so up until a time slice  $t$ ;  $P(\mathbf{X}_t | \mathbf{e}_{1:t})$ . A close related calculation can be done to get  $P(\mathbf{e}_{1:t})$
- *Prediction*: Calculate the probability of a future belief state at time slice  $k$ , where  $k > 0$ , given the present evidence values;  $P(\mathbf{X}_{t+k} | \mathbf{e}_{1:t})$ .
- *Smoothing*: Calculate the belief state of a previous time slice  $k$ , given the observed evidence up until  $t$ , where  $0 < k < t$ ;  $P(\mathbf{X}_k | \mathbf{e}_{1:t})$ .
- *Most likely explanation*: Calculate the most likely belief states that would give the observed evidence;  $\text{argmax}_{\mathbf{x}_{1:t}} P(\mathbf{x}_{1:t} | \mathbf{e}_{1:t})$

Together with these tasks learning is also possible. Both the transition model and the sensor model can be learned through observation.

### Filtering/forward message

To calculate the belief state at time slice  $t + 1$ , we can use the following equation;

$$\begin{aligned} P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) &= P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}, \mathbf{e}_{t+1}) && \text{(dividing up the evidence)} \\ &= \alpha P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}, \mathbf{e}_{1:t}) P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}) && \text{(using Bayes rule)} \\ &= \alpha P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t}) && \text{(using sensor Markov assumption)} \end{aligned} \quad (15.4)$$

The terms are as follows;  $\alpha$ , the normalisation constant;  $P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1})$ , the sensor model is used to update the belief state;  $P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t})$ , a prediction of the next belief state. This last prediction is based on the current belief state  $\mathbf{X}_t$  and this is used integrated in the following way:

$$\begin{aligned} P(\mathbf{X}_{t+1} | \mathbf{e}_{1:t+1}) &= \alpha P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1} | \mathbf{x}_t, \mathbf{e}_{1:t}) P(\mathbf{x}_t | \mathbf{e}_{1:t}) && \text{(predicting } \mathbf{X}_{t+1} \text{ using } \mathbf{X}_t) \\ &= \alpha P(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} P(\mathbf{X}_{t+1} | \mathbf{x}_t) P(\mathbf{x}_t | \mathbf{e}_{1:t}) && \text{(using Markov assumption)} \end{aligned} \quad (15.5)$$

Notice that  $P(\mathbf{X}_{t+1} | \mathbf{x}_t)$  comes from the transition model. Further  $P(\mathbf{x}_t | \mathbf{e}_{1:t})$  can be viewed as a *forward message* (as this is calculated at each time step using the previous state) and Eq. 15.5 can be in shorthand as:

$$\mathbf{f}_{1:t+1} = \alpha \text{FORWARD}(\mathbf{f}_{1:t}, \mathbf{e}_{1:t+1})$$

where  $\mathbf{f}_{1:0} = P(\mathbf{X}_0)$



## Prediction

This is done by using Eq. 15.5 using just the evidence observed so far,  $\mathbf{e}_{1:t}$  and the previous predicted state,  $\mathbf{X}_{t+k}$ :

$$P(\mathbf{X}_{t+k+1} \mid \mathbf{e}_{1:t}) = \sum_{\mathbf{x}_{t+k}} P(\mathbf{X}_{t+k+1} \mid \mathbf{x}_{t+k}) P(\mathbf{x}_{t+k} \mid \mathbf{e}_{1:t}) \quad (15.6)$$

These predictions will converge to a fixed state called the *stationary distribution* after a number of time steps called the *mixing time*, dooming any attempt to predicting too far into the future.

## Likelihood

This is useful when reviewing the assumed belief states calculated given the evidence and is done using the *likelihood message*,  $\mathbf{l}_{1:t}(\mathbf{X}_t) = P(\mathbf{X}_t, \mathbf{e}_{1:t})$ , this in turn can be inferred into the following:

$$\begin{aligned} \mathbf{l}_{1:t+1} &= FORWARD(\mathbf{l}_{1:t}, \mathbf{e}_{1:t+1}) \\ \text{and from } \mathbf{l}_{1:t} &\text{ the actual likelihood can be calculated} \end{aligned} \quad (15.7)$$

$$L_{1:t} = P(\mathbf{e}_{1:t}) = \sum_{\mathbf{x}_t} \mathbf{l}_{1:t}(\mathbf{x}_t)$$

## Smoothing

Computing the belief state of a previous state,  $P(\mathbf{X}_k \mid \mathbf{e}_{1:t})$  at time  $0 \leq k \leq t$ , using evidence up to the present,  $t$ , can be done by splitting the computation into two:

$$\begin{aligned} P(\mathbf{X}_k \mid \mathbf{e}_{1:t}) &= P(\mathbf{X}_k \mid \mathbf{e}_{1:k}, \mathbf{e}_{k+1:t}) && \text{split the evidence} \\ &= \alpha P(\mathbf{X}_k \mid \mathbf{e}_{1:k}) P(\mathbf{e}_{k+1:t} \mid \mathbf{X}_k, \mathbf{e}_{1:k}) && \text{using Baye's rule} \\ &= \alpha P(\mathbf{X}_k \mid \mathbf{e}_{1:k}) P(\mathbf{e}_{k+1:t} \mid \mathbf{X}_k) && \text{using Eq. 15.2} \\ &= \alpha \mathbf{f}_{1:k} \times \mathbf{b}_{k+1:t} && \text{element wise multiplication} \end{aligned} \quad (15.8)$$

Where  $\mathbf{b}_{k+1:t} = P(\mathbf{e}_{k+1:t} \mid \mathbf{X}_k)$  is the backward message that is the recursive analogy to the forward message that is calculated backwards from  $t$ :

$$\begin{aligned} P(\mathbf{e}_{k+1:t} \mid \mathbf{X}_k) &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1:t} \mid \mathbf{x}_{k+1}, \mathbf{X}_k) P(\mathbf{x}_{k+1} \mid \mathbf{X}_k) \\ &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1:t} \mid \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1} \mid \mathbf{X}_k) && \text{conditioning on } \mathbf{X}_{k+1} \\ &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1}, \mathbf{e}_{k+2:t} \mid \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1} \mid \mathbf{X}_k) && \text{conditional independence} \\ &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1} \mid \mathbf{x}_{k+1}) P(\mathbf{e}_{k+2:t} \mid \mathbf{x}_{k+1}) P(\mathbf{x}_{k+1} \mid \mathbf{X}_k) \end{aligned} \quad (15.9)$$

Notice that the first and third products are from the models, and the second is the recursive call.

## Most likely sequence

Because of the *Markov assumptions* the transition between states are dependent on the previous state, therefore, if you change the question from which sequence is most likely to what path through the states that takes (see Fig 15.5 in book), the most likely state at time  $t + 1$  is dependent on the most likely state at  $t$ . This relationship can be expressed as:

$$\begin{aligned} & \max_{\mathbf{x}_1 \dots \mathbf{x}_t} P(\mathbf{x}_1, \dots, \mathbf{x}_t, \mathbf{X}_{t+1} \mid \mathbf{e}_{1:t+1}) = \\ & \alpha P(\mathbf{e}_{1:t+1} \mid \mathbf{X}_{t+1}) \max_{\mathbf{x}_t} (P(\mathbf{X}_{t+1} \mid \mathbf{x}_t) \max_{\mathbf{x}_1 \dots \mathbf{x}_{t-1}} P(\mathbf{x}_1, \dots, \mathbf{x}_t \mid \mathbf{e}_{1:t})) \end{aligned} \quad (15.10)$$

The *Hidden Markov Model* (HMM) is a way to represent the world state as a single discrete random variable, where the values are the possible states of the world. This will be the state variable  $X_t$ , which is a vector of the size  $S$ , where  $S$  is the number of possible states of the world. The *transition model*  $P(X_t \mid X_{t-1})$  is written as a  $S \times S$  matrix  $\mathbf{T}$ , where  $\mathbf{T}_{ij} = P(X_t = j \mid X_{t-1} = i)$ , the transition from state  $i$  to state  $j$ . Next the sensor model is written as a diagonal matrix  $\mathbf{O}$ , where the diagonal  $\mathbf{O}_i = P(e_t \mid X_t = i)$  and the rest are 0. The forward and backward messages then becomes:

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^T \mathbf{f}_{1:t} \quad (15.11)$$

$$\mathbf{b}_{k+1:t} = \mathbf{T} \mathbf{O}_{k+1} \mathbf{b}_{k+2:t} \quad (15.12)$$

In the real world there are often failures in the observable variables, like a sensor failing, both temporally and permanently. These event can be handled using *Dynamic Bayesian Networks* (DBN) which incorporates several different models:

- *Transient failure model*: This model handles the event of a temporary error in the observable variable by adding a small probability that the wrong value might be observed.
- *Persistent failure model*: This model handles the event when the the observable variable permanently return the wrong value. This is done by adding a new variable to the DBN that represents the state of the sensor of that observable variable.

## 16 Making simple decisions

When dealing with decisions we define  $RESULT(a)$  as a random variable whose values are possible outcome states. The probability of outcome  $s'$  of an action  $a$  given evidence  $\mathbf{e}$  is written as:  $P(RESULT(a) = s' \mid a, \mathbf{e})$ . An agents preference is captured by the *utility function*,  $U(s)$ , and the *expected utility* of an action given evidence,  $EU(a \mid \mathbf{e})$ , is just the average utility value of the possible outcomes:

$$EU(a \mid \mathbf{e}) = \sum_{s'} P(RESULT(a) = s' \mid a, \mathbf{e}) U(s') \quad (16.1)$$

The principle of *maximum expected utility* (MEU) says that a rational agent should choose the action that maximises the agent's expected utility;  $action = \arg\max_a EU(a \mid \mathbf{e})$  Note that the utility score is not the same as the performance score, as the utility score is the expected value of an action, the performance score is the actual value of the action. For an agent to be *unbiased* the difference of expected value and predicted value must be 0.

Some notation on preference:

$A \prec B$  A is preferred over B

$A \sim B$  A and B are equally preferred

$A \succsim B$  A is preferred over B or they are equal

The probability of the outcomes of an action is represented as a lottery  $L$  with each outcome state,  $S_1, \dots, S_n$ , and their probabilities,  $p_1, \dots, p_n$ ;  $L = [p_1, S_1; p_2, S_2; \dots; p_n, S_n]$ .

For a preference relation to be reasonable it must fulfil the following restraints:

- *Orderability*: Exactly one of the following must hold;  $(A \prec B), (B \prec A), (A \sim B)$   
The states must either be equal or one can order them.
- *Transitivity*:  $(A \prec B) \wedge (B \prec C) \implies (A \prec C)$
- *Continuity*:  $A \prec B \prec C \implies \exists p[p, A; 1 - p, C] \sim B$   
There must be a way to “represent” the outcome of B, by using A and C.
- *Substitutability*:  $A \sim B \implies [p, A; 1 - p, C] \sim [p, B; 1 - p, C]$   
Here  $\sim$  can be swapped with  $\prec$ .
- *Monotonicity*:  $A \prec B \implies (p > q \iff [p, A; 1 - p, B] \prec [q, A; 1 - q, B])$   
If A is preferred over B, then the lottery with the highest probability of getting A is preferred.
- *Decomposability*:  $[p, A; 1 - p, [q, B; 1 - q, C]] \sim [p, A; (1 - p)q, B; (1 - p)(1 - q), C]$   
Two consecutive lotteries can be turned into one.

Some utility function properties:

- $U(A) < U(B) \iff A \prec B$
- $U(A) = U(B) \iff A \sim B$
- The utility of a lottery is defined as:  
 $U(L) = U([p_1, S_1; p_2, S_2; \dots; p_n, S_n]) = \sum_i p_i U(S_i)$

As an agent might only be interested in ordering states, the actual utility value isn't necessary, these functions are called *value functions* or *ordinal utility functions*. *Preference elicitation* is used to discover a utility function and is done by presenting choices to an agent the using its observed preferences.

When an agent prefers to take risks for higher reward they are *risk-seeking*, the reverse is *risk-averse*, preferring safer bets. The value an agent accepts as guaranteed after a bet is the *certainty equivalent*, the difference between this and the *expected monetary value* (EMV) is the *insurance premium*.

*Normative theory* is the theory of how an agent should act, and *descriptive theory* is the theory of how the agent actually acts. Not all actions are completely rational as covered by the *certainty effect*, that a person usually chooses the guaranteed win. Something similar is covered by the *ambiguity aversion*, that people go for the known probabilities, instead of the unknown probabilities.

Problems where the utility is based on several variables are called *multiattribute utility theory*. Here the terms *strict dominance* means that a solution  $S_1$  is better in every utility aspect than an solution  $S_2$ , and *stochastic dominance*, if two actions,  $A_1, A_2$  leads to probability distributions  $p_1(x), p_2(x)$  the  $A_1$  stochastically dominates  $A_2$  if:  $\forall x \int_{-\infty}^x p_1(x') dx' \leq \int_{-\infty}^x p_2(x') dx'$ .

*Preference independence* is when two variables are not dependent on the value of a third variable. *Mutual preferential independence* (MPI) says that each variable is equally important, so they don't

affect the other variables during a trade off. For attributes that are MPI the preferred behaviour can be described by maximising the function *additive value function*,  $V(x_1, \dots, x_n) = \sum_i V_i(x_i)$ . The preference independence can be extended with *utility independence*(UI) to cover lotteries. A set of attributes is *mutually utility independent* (MUI) if each of its subsets is UI of the remaining attributes. MUI can be expressed using a *multiplicative utility function*.

Rational decisions can be done using decision networks, an extension of Bayesian networks. The network is constructed by the following node types:

- *Chance nodes*: same as for Bayesian networks, circles with variables
- *Decision nodes*: rectangles representing where the decision maker can take actions
- *Utility nodes*: diamonds representing the utility function.

These networks can be evaluated in the following way:

1. Set the evidence variables for the current state.
2. For each possible value of the decision node:
  - (a) Set the decision node to that value.
  - (b) Calculate the posterior probabilities for the parent nodes of the utility node, using a standard probabilistic inference algorithm
  - (c) Calculate the resulting utility for the action.
3. Return the action with the highest utility.

In information value theory the *value of perfect information*(VPI) is expressed as the following:

$$VPI_e(E_j) = \left( \sum_k P(E_j = e_{jk} \mid \mathbf{e}) EU(\alpha_{e_{jk}} \mid \mathbf{e}, E_j = e_{jk}) \right) - EU(\alpha \mid \mathbf{e})$$

Where  $E_j$  is the variable we don't know but want to find out,  $\mathbf{e}$  is the observations so far,  $e_{jk}$  is the new evidence and  $\alpha$  is the value of the current best action.

Lastly in decision-theoretic expert systems *decision analysis* studies the application of decision theory on actual decision problems. There are usually talk about two roles, the *decision maker* that states preferences about outcomes, and the *decision analyst* that enumerates the possible actions and outcomes and elicits preferences from the decision maker to determine the best action. A decision-theoretic expert system is created in the following way:

1. *Create a causal model*: determine all possible variables, utility functions and actions.
2. *Simplify to a qualitative decision model*: Remove variables that does not affect the choice of action.
3. *Assign probabilities*
4. *Assign utilities*: These can be enumerated from the probabilities
5. *Verify and refine the model*: Using a *gold standard*(validation set)
6. *Perform sensitivity analysis*: If small changes results in big changes in the chosen action, the system might not be ready.

## 17 Making complex decisions

Focuses on *sequential decision problems* where the utility function depends on a sequence of decisions. First we look at the environment as *fully observable*, that is has a *transition model*,  $P(s' | a, s)$  and that the transitions are *Markovian*. For each state an agent receives a reward  $R(s)$  depending on that state and the utility function is just the sum of these rewards for the *environment history*. This is called a *Markov decision problem* (MDP).

When adding uncertainty to actions, i.e. there is a probability to take a right instead of a left, a sequence of actions does not guarantee success, therefore only the optimal action at each state is stored,  $\pi(s)$ , called a *policy*. The *optimal policy* for a state is the policy that gives the highest expected utility. As the performance of an agent is based on the states it visited this becomes a story in *multiattribute utility theory*. There is also the properties of *finite* or *infinite horizon*, where the agent either must reach a terminate state after  $N$  states, or has no such constriction, on wards the an infinite horizon is assumed. If the optimal policy in a state,  $\pi^*(s)$ , changes over time it is *nonstationary*, and if it does not it is *stationary*. If an agent is stationary for preference it has the same preference-order for  $[s_0, s_1, s_2, \dots]$  and  $[s'_0, s'_1, s'_2, \dots]$ , as for  $[s_1, s_2, \dots]$  and  $[s'_1, s'_2, \dots]$  iff  $s_0 = s'_0$ .

If a system is stationary the following utility function can be assigned:

1. *Additive rewards*:  $U_h([s_0, s_1, s_2, \dots]) = \sum_{i=0} R(s_i)$

2. *Discounted rewards*:  $U_h([s_0, s_1, s_2, \dots]) = \sum_{i=0} \gamma^i R(s_i)$

Here  $0 < \gamma \leq 1$  is the *discount factor*, this also gives an interest rate of  $(1/\gamma) - 1$ .

When assuming the infinite horizon these functions become tricky as they both go on to plus/minus infinity and we then have to order infinities. There are three solutions to this:

1. *Discounted rewards*: If  $\gamma < 1$  this is infinite geometric series, which is finite:

$$U_h([s_0, s_1, s_2, \dots]) = \sum_{i=0}^{\infty} \gamma^i R(s_i) \leq \sum_{i=0}^{\infty} \gamma^i R_{\max} = R_{\max}/(1 - \gamma)$$

2. *Proper policy*: This requires there to be a terminate state, so infinite series are never encountered.

3. *Average reward*: Averaging the rewards gives a way to handles the infinities.

Note, discounted rewards are the best.

Let the random variable  $S_t$  denote the state the agent reaches at time step  $t$  when using policy  $\pi$ . the expected utility then becomes:

$$U^\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t R(S_t)\right] \quad (17.2)$$

Next let  $\pi_s^*$  denote the optimal policy of state  $s$ , and it is chosen by:

$$\pi_s^* = \operatorname{argmax}_{\pi} U^\pi(s) \quad (17.3)$$

And the best action can be calculated by:

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} \sum_{s'} P(s' | s, a) U(s') \quad (17.4)$$

To calculate the optimal utility one of two algorithms can be used; *value iteration* or *policy iteration*. The value iteration algorithm is based on the *Bellman equation*, these also give unique solutions:

$$U(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s' \in NEIGH(s)} P(s' | s, a) U(s') \quad (17.5)$$

These equations are non linear and must therefore be solved using non-regular methods, one is iterative where each state value is updated for each iteration by the *Bellman update*:

$$U_{i+1} \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s' \in NEIGH(s)} P(s' | s, a) U_i(s') \quad (17.6)$$

By showing that the Bellman update is a *contraction* we guarantee that the value iteration converges. There are two properties of a contraction function; it has only one fixed point and for all inputs to the function, the output is closer to the fixed point. By viewing the Bellman update as an operator  $B$  and  $U_i, U'_i$  as vectors of utilities at iteration  $i$ ,  $U$  as the true utilities and the max norm of a vector  $\|U\| = \max_s |U(s)|$  it can be shown that value iteration converges:

$$\begin{aligned} \|BU_i - BU'_i\| &\leq \gamma \|U_i - U'_i\| && B \text{ is an contraction of factor } \gamma \\ \|BU_i - U\| &\leq \gamma \|U_i - U\| && \text{using that } BU = U \text{ it shows that for each iteration} \\ &&& \text{the error is reduced by at least } \gamma \end{aligned} \quad (17.7)$$

From this it is also possible to calculate number of iterations needed to reach an error of size  $\epsilon$ ;  $N = \lceil \log(2R_{\max}/\epsilon(1-\gamma))/\log(1/\gamma) \rceil$

Through contraction it is also shown that if the update is small then the error is also small;

$$\|U_{i+1} - U_i\| < \epsilon(1-\gamma)/\gamma \implies \|U_{i+1} - U\| < \epsilon$$

Further the *policy loss* is connected to error by the following, where  $U^{\pi_i}$  is the policy if  $\pi_i$  was executed:  $\|U_i - U\| < \epsilon \implies \|U^{\pi_i} - U\| < 2\epsilon\gamma/(1-\gamma)$

Next we look at the *policy iteration* that contains two steps:

1. *Policy evaluation*: given a policy  $\pi_i$ , calculate the utility of each state if  $\pi_i$  was executed,  $U_i = U^{\pi_i}$ .
2. *Policy improvement*: Calculate a new MEU policy  $\pi_{i+1}$ , using one-step look-ahead based on  $U_i$ .

To improve on value iteration, the max operation is removed to make the algorithm linear by using  $\pi_i$  instead of finding the best action, changing the Bellman equation to:  $U_i(s) = R(s) + \gamma \sum_{s'} P(s' | s, \pi_i(s)) U_i(s')$ , and the Bellman update to:  $U_{i+1}(s) \leftarrow R(s) + \gamma \sum_{s'} P(s' | s, \pi_i(s)) U_i(s')$ . This algorithm is called the *modified policy iteration*.

Next we look at *partially observable MDPs* (POMDP), which uses the same elements as a MDP, but extends it with a *senor model*,  $P(e | s)$ . Given the old belief state  $b(s)$ , an action  $a$  and perceived evidence  $e$  the next belief state becomes the following:  $b'(s') = \alpha P(e | s') \sum_s P(s' | s, a) b(s)$  or  $b' = FORWARD(b, a, e)$

## 18 Learning from examples

*Learning* is when an agent improves its performance by making observations of the environment. The improvements and techniques used for improvements are dependent on the following factors:

- Which components are improved; these can be for i.e. A mapping from the current state to actions, infer relevant properties from the precept sequence...
- What prior knowledge the agent has
- What representation is used, in todays machine learning this usually is a *factored representation*(feature vectors).
- What feedback is available. This is the difference between *unsupervised learning*(no feedback), *reinforcement learning*(gets rewards or punishments) and *supervised learning*(mapping input to output, so it knows the answer from the beginning). *Semi-supervised learning* is having a partially labelled training set.

Finding a function/rule from input-output pairs is called *inductive learning*. *Analytic* or *deductive learning* is to deduct new rules that entails from the previous known rules.

A *training set* is a set of input-output pairs, and a system tries to give a mapping(function) from the input to the output, this mapping is called the *hypotheses*,  $h$ . The hypotheses is tested using a *testing set*. A hypotheses *generalises* well if it scores high on the testing set. When the output,  $y$ , is an element of a finite set its a *classification*, while if  $y$  is a number its called a *regression*. The set of possible hypotheses is called the hypotheses space,  $\mathbb{H}$ . If  $h$  agrees with all the data points it's called *consistent*. The problem of choosing between several consistent hypotheses is called the *Ockham's razor*. The *true function*,  $f$ , is the correct mapping from input to output and a problem is *realisable* if the  $f \in \mathbb{H}$ . There is a trade-off between the expressiveness of  $\mathbb{H}$  and the complexity of finding a good  $h$ .

A *decision tree* is a representation of a multiattribute function that returns a single output. The output variable is called the *goal predicate* and its values are the leaf nodes of the tree. The tree learning algorithm is greedy, as it chooses the most important attribute first, and solves the problem recursively by removing that attribute from the problem and the examples it can classify with just that attribute. For this recursion to work there are 4 cases to handle:

1. If the remaining examples are all of the same classification, we are done.
2. If there are several classifications of the examples, choose the best attribute to split the examples on.
3. If there is no examples left we have found an unknown combination of attributes the only thing that can be done is to return the *plurality* classification of all the examples(the most common output value).
4. If there are no attributes left, but examples of different classifications the only thing that can be done is to return the *plurality* classification of the remaining examples.

*Entropy* is the measure of uncertainty of a random variable, and is defined as:  $H(V) = -\sum_k P(v_k) \log_2 \frac{1}{P(v_k)} = -\sum_k P(v_k) \log_2 P(v_k)$

This is generalised to Boolean variables with a probability of  $q$  of being true:  $B(q) = -(q \log_2 q + (1 - q) \log_2 (1 - q))$

If a training set has  $p$  positive examples and  $n$  negative, the entropy of the goal attribute is:  $H(Goal) = B(\frac{p}{p+n})$ . A attribute  $A$  with  $d$  distinct values divides the training set  $E$  into  $d$  subsets. Each subset  $E_k$  has  $p_k$  and  $n_k$  pos/neg examples so that we need  $B(p_k/(p_k + n_k))$  extra bits of information to answer a query. A randomly chosen example has the  $k$ th value for the attribute with a probability of  $(p_k + n_k)/(p + n)$  so the remaining entropy after testing on attribute  $A$  is:

$Remainder(A) = \sum_{k=0}^d \frac{p_k+n_k}{p+n} B(\frac{p_k}{p_k+n_k})$ . The *information gained* from testing on  $A$  is the reduction in entropy:  $Gain(A) = B(\frac{p}{n+k}) - Remainder(A)$

*Overfitting* is when the learner deduce rules that are just accidentally in the data set, like if the learner only has two examples of a die being thrown and it becomes a 6, and the thrower also has a white t-shirt, the learner might say that each time a die is thrown by someone in a white shirts it will be a 6. To avoid overfitting in decision tree learning *decision tree pruning* is used, iterate over the nodes above the leaf nodes and see if they are irrelevant (there is only noise in the leaf nodes) it is turned into a leaf node. This discovery is done by a *significance test* where we assume the *null hypotheses* that there are no underlying pattern in the data and we calculate a distribution that then either verifies or rejects the hypotheses, in decision tree algorithms this is usually done using  $\chi^2$ -pruning. A combination of  $\chi^2$ -pruning and information gain is called *early stopping*.

To extend the applicability of decision trees some of the following issues must be addressed:

- *Missing data*: Not having values for every attribute in every example
- *Multivalued attributes*: If an attribute gives a classification for several examples it will have great information gain, but it might not be useful (splitting on time for example). This is solved by using *gain ratio*.
- *Continuous and integer-valued input attributes*: Find the best place for a *split point* (where to split on height to get the biggest information gain).
- *Continuous-valued output attributes*: If the output should be continuous the tree becomes a *regression tree* and the leaf node is usually a function.

When choosing the best hypotheses the *stationary assumption* is made, each example is independent of the previous one and each example has the same probability distribution of its attributes. These examples are called *independent and identically distributed (I.I.D)*. The *error rate* is defined as the proportion of mistakes it makes, i.e.  $h(x) = y', y \neq y'$  for an  $(x, y)$  example. Splitting examples into training and test set is called *holdout cross-validation*, but this reserves some of the examples for testing, and reduces the number of examples for training. A solution is *k-fold cross-validation*; Split the data set into  $k$  sets, perform  $k$  rounds of training with  $1/k$  of the subset as testing data. If  $n = k$  this is called *leave-one-out cross-validation (LOOCV)*. Another problem is *peeking* at the test data, usually done by tweaking the settings of the system based on error rates of the test set, this is avoided by using a *validation set* in place of a test set.

There are two ways to improve the choosing of the best hypotheses; improving *model selection*, defining a good hypotheses space; or *optimisation* finding the best hypotheses. An algorithm for this is a *wrapper* that tests learners based on a variable, *size*, which corresponds to the complexity of the hypotheses space, and returns the best learner after looking at every size of complexity in the hypotheses space.

Next we want to maximise the expected utility of the learner, not just minimise the error rate, this is usually done by minimising *loss function*,  $L(x, y, y') = U(f(x) = y) - U(h(x) = y')$ . The *generalisation loss* of  $h$  with respect to  $L$  is using the set of all examples  $\epsilon$ :  $GenLoss_L(h) = \sum_{(x,y) \in \epsilon} L(y, h(x))P(x, y)$  and the best hypotheses is defined as:  $h^* = \operatorname{argmin}_h GenLoss_L(h)$  but since  $P(x, y)$  is unknown this loss is estimated with *empirical loss* on a set of examples,  $E$ :  $EmpLoss_{L,E}(h) = \frac{1}{N} \sum_{(x,y) \in E} L(y, h(x))$  and the estimated best hypotheses  $h' = \operatorname{argmin}_h EmpLoss_{L,E}(h)$ .  $h'$  and  $h^*$  differ for 4 reasons:

- *Unreliability*,  $f \notin \mathbb{H}$



- *Variance*, different sets of examples generates different learners
- *Noise*,  $f$  is non deterministic,  $f(x) = y \wedge f(x) = z$  (also called *noisy*)
- *Computational complexity*, if  $\mathbb{H}$  is big or complex it might be too hard to search it all.

*Regularisation* is the process of penalising complex hypotheses. Another way to “decomplexify”  $\mathbb{H}$  is by *feature selection*, where some features are ignored.

## Neural networks

A NN is a network of nodes connected by links. A link between the nodes  $i$  and  $j$  propagates the *activation*  $a_i$  from  $i$  to  $j$ , it also has a *weight*  $w_{ij}$  which determines the strength of the activation from  $i$  to  $j$ . To activate, node  $j$  sums its inputs:  $in_j = \sum w_{ij}a_i$  and then applies its *activation function*,  $g$ :  $a_j = g(in_j)$ . A network where the links only go one way is called *feed-forward network* where each node has its inputs from the upstream nodes and sends the activation to the downstream nodes. A *recurrent network* sends its activation’s back through the network. Using several *hidden layers* a NN can do *nonlinear regression*. Learning is done by *back-propagation*, which changes the weights in the network based on a learning rate,  $\alpha$ , and that weights “fault” for loss:

$$w_{jk} \leftarrow w_{jk} + \alpha a_j \delta_k \quad (18.11)$$

where  $\delta_k = Err_k g'(in_k)$  and  $Err_k$  is the  $k$ th component of the error vector  $\vec{y} - \vec{h}_w$ . NNs can be overfitted and as cross-validation can be used to avoid peeking. *Optimal brain damage* can also be used in NN to remove unused links.

Models where the number of input variables are bounded are called *parametric models*, the opposite are called *nonparametric models*. Examples of these are *instance-based learning* or *memory-based learning* where the hypotheses uses previous examples to infer a solution to the next problem. One example is the  $k$  – *nearestneighbors*, where the “closest” looking examples are used to classify/find a solution to a new example. The distance here is typically measured in *Minkowski distance*, or  $L^p$ -norm:  $L^p(\vec{x}_j, \vec{x}_q) = (\sum_i |x_{ji} - x_{qi}|^p)^{1/p}$ . To avoid one dimension of the attribute vector to dominate, the components are *normalised*.

## 21 Reinforcement learning

The agent gets reward for good or bad behaviour. Goal of reinforcement learning is to learn an optimal policy for the environment. Until now we have three agent designs:

- *Utility-based agent*: Learns a utility function, using it to select actions that maximise the expected outcome.
- *Q-learning*: Learns an action-utility(Q) function, given the expected utility of doing an action in the given state.
- *Reflex agent*: Learns a policy that maps directly from states to actions.

*Passive learning* is when the policy is fixed and the utility function is unknown. *Active learning* is when the agent also must find what to do (the transition model is unknown), where the issue is *exploration*. The simplest utility estimation is the *direct utility estimation*, this method reduces

reinforcement learning to inductive learning, but it ignores the fact that the utilities of each state is dependent on other states;

$$U^\pi(s) = R(s) + \gamma \sum_{s'} P(s' | s, \pi(s)) U^\pi(s') \quad (21.2)$$

*Adaptive dynamic programming*(ADP) uses the constraints among the utilities of states by learning the transition model and solving the Markov decision process. There are two possible improvements to this approach:

1. *Bayesian reinforcement learning*: Assume hypotheses'  $h$  and the prior probability  $P(h)$  and obtain the posterior probability  $P(h | e)$  through observation of data, when the agent is done learning the optimal policy is obtained from the highest expected utility:  $\pi^* = \operatorname{argmax}_\pi \sum_h P(h | e) u_h^\pi$
2. *Robust control theory*: Similar, but instead defines the optimal robust policy to be the best hypotheses in the worst case:  $\pi^* = \operatorname{argmax}_\pi \min_h u_h^\pi$

*Temporal-difference learning* uses difference in utilities between successive states:

$$U^\pi(s) \leftarrow U^\pi(s) + \alpha(R(s) + \gamma U^\pi(s') - U^\pi(s)) \quad (21.3)$$

Here  $\alpha$  is the *learning rate*, can be a constant or a function. Notice that the update rule for TD does not need the transition model.

Passive learning agents works with fixed policy's, active learning agents discovers the best actions. This is done using Bellman equations:

$$U(s) = R(s) + \gamma \max_a \sum_{s'} P(s' | s, a) U(s') \quad (21.4)$$

ADP is an *greedy agent* and they usually don't find the optimal policy as they get stuck on local optima. Agents must choose in the tradeoff between *exploration* and *exploitation*. This is introduced to the agent through the use of an *exploitation function*,  $f(u, n)$ , where  $u$  represents the wanted values(the greed) and  $n$  represents the how often an action have been tried. Implemented in Eq. 21.4, where  $N(s, a)$  gives the number of times the action in that state have been tried:

$$U(s) = R(s) + \gamma \max_a f\left(\sum_{s'} P(s' | s, a) U(s'), N(s, a)\right) \quad (21.5)$$

Next we create an *active TD learning agent*. This is done using *Q-learning*, which learns an action-utility function using a stored value,  $Q(s, a)$ , of a given action in a state that translates to the utility function by  $U(s) = \max_a Q(s, a)$ . As it uses  $Q(s, a)$  it no longer needs the transition model and the  $Q$  values are constrained using the equation:

$$Q(s, a) = R(s) + \gamma \sum_{s'} P(s' | s, a) \max_{a'} Q(s', a') \quad (21.6)$$

The update equation for TD Q-learning becomes:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a)) \quad (21.7)$$

The TD Q-learning is called a *off-policy* algorithm as it does not care about the actual policy, *SARSA* is a related *on-policy*, it is updated using:  $Q(s, a) \leftarrow Q(s, a) + \alpha(R(s) + \gamma Q(s', a') - Q(s, a))$

## 22 Philosophical Foundations

*Weak AI* is the assertion that machines can act as if they were intelligent, while *strong AI* is the assertion that machines are actually thinking (not simulating). Alan Turing defined the *Turing test* to check if a machine could act as a human: An interrogator has 5 minutes to talk to a man or machine and must after that try and decide if he talked to a man or machine, if the machine tricks the interrogator 30% of the time it passes and can act as a human. Turing, and others, also gave more objections to the possibility of intelligent machines:

1. *Disability*: A machine can never do X (be kind, resourceful, friendly...)
2. *Gödel*:  $G(F)$  is a sentence of  $F$ , but cannot be proved within  $F$ , if  $F$  is consistent, then  $G(F)$  is true.
3. *Informality*: Behaviour is not rational and informal, you can't program spontaneous chatter. The *qualification problem* is that it is impossible to program to handle everything.

Other philosophical jargon:

- *Consciousness*: A machine must be aware of its mental state and actions
- *Phenomenology*: A machine must feel its emotions
- *Intentionally*: Is a machine's beliefs, states, desires and other representations actually "real"
- *Dualism*: What is the connection between mind and body, is it just as inferring over sensors and responding with actuators.
- *Functionalism*: A mental state is any intermediate causal condition between input and output. Brain replacement experiment: Say we can change out any neuron in the brain with an exact artificial copy. If we swap out every neuron one by one, are you the same person? Are you even a person anymore?
- *Biological naturalism*: Our brain is just electronic signals in the brain, what's the difference from artificial neural networks:
- *Technological singularity*: If humans create a machine smarter than them, they would be able to create something smarter than them and so on until humans are left in the dust.
- *Transhumanism*: A movement that looks forward to a future where humans are merged or replaced with machines.

Ethical problems:

- People might lose their jobs to automation
- People might lose their sense of uniqueness
- People might get too much/little leisure time
- AI might be used to reach undesirable goals (war, terminator, mass surveillance)
- AI might result in the loss of accountability
- AI might be the end of mankind

## 23 AI: The Present and Future

Which components of an agent can be improved?

- Interaction with the environment through sensors and actuators
- Keeping track of the state of the world
- Projecting, evaluating, and selecting future courses of action
- Utility as an expression of preferences
- Learning

Are we going in the right direction? There are 4 specifications of intelligent agents:

- *Perfect rationality*: Always chooses the best action, maximises utility
- *Calculated rationality*: Perfect rationality in action, by the time the perfect rational action has been found, time has passed it is no longer needed.
- *Bounded rationality*: Returns a good enough approximation of the optimal action. More realistic.
- *Bounded optimally*: Good enough, given its computational bounds. A variation is *Asymptotic bounded optimality*(ABO), where the bounds are set to help with comparison between systems.

## 2 CBR paper

### 1 Background

CBR is a problem solving paradigm that uses specific knowledge of previous experiences on new problems(cases). It is a incremental, sustained learning method, as it can add each solved problem to its case base, adapting to new problems and solutions. Because it learns and the machine learning community is the main contributors to the research forward, CBR is seen as a sub-field of machine learning. The paper states lots of empirical evidence that humans uses CBR in their reasoning. Examples of how humans use CBR can be:

- Physician: reuses diagnosis and treatments on patient that has the same type of symptoms as a previous patient.
- Financial consultant: Decides to give a loan to a new customer based on a similar previous customer.

Note that CBR does not just use the same solution on the same type of problem, it uses a solution to a problem that it remembers working.

CBR requires methods to; extract relevant knowledge, integrate each case into the knowledge base/structure and index the cases for later use. One issue with CBR, if cognitive plausibility is a guiding principle, is that it should also incorporate other types of knowledge as well, not only from previous cases.

### 2 History

A lot of people have worked with CBR-systems through the years, noting the CREEK system and integration framework from NTNU Trondheim developed by Agnar Aamodt and sintef, because NTNU and sintef.

### 3 Fundamentals of case-based reasoning methods

There are 5 main differences in CBR methods; 1. how they id the current case, 2. how they find past similar cases, 3. how they use those cases to propose a solution to the current case, 4. how they evaluate the proposed solution and 5. how they learn from that solution. Further more there are 5 main types of CBR:

- Exemplar-based reasoning classifies a new case based on the most similar previous case. Exemplar is just another word for case. It does not change the solution of the previous case in any way.

- Instance-based reasoning is a highly syntactic specialisation of exemplar-based reasoning. It uses simple representations(usually a feature-vector) and focuses on automated learning(no human interaction while learning). This is a non-generalisation approach and is basically very strict exemplar-based reasoning. Learning OR/AND examples.
- Memory-based reasoning collects cases into a huge memory and reasons through searching and accessing this memory. Therefore has a big focus on organising and accessing this memory. The utilization of parallel processing techniques is a characteristic of these methods.
- Case-based reasoning covers the most typical methods. Cases are usually assumed to have a certain degree of richness and complexity in information. These methods are capable of modifying/adapting solutions and utilises general background knowledge. They also borrows theories from cognitive psychology.
- Analogy-based reasoning solves new cases based on previous cases from another problem domain. It is a sub-field concerning identification and utilization of cross-domain analogies. The mapping problem is a central problem in this field: How to map the solution from an analogue(called the source or base) to the present problem(called the target).

All CBR system implements the following four core function, called the CBR cycle:

- Retrieve the most similar case/cases
- Reuse the information/knowledge of that case
- Revise the proposed solution
- Retain parts of the experience likely to be useful in the future

Together with the CBR cycle, a CBR system can be analysed in a task oriented view where the system is viewed as an agent that has goals and the means to achieve them. A system can be described from three different perspectives; Tasks, methods and domain knowledge models.

There are 5 problem areas that a CBR methods must have coherent solutions for to be called a CBR:

- Knowledge representation
- Retrieval method
- Reuse method
- Revise method
- Retain method

## 4 Representation of cases/Knowledge representation

CBR uses past cases to solve new ones, so representing past cases(knowledge base) in a way that makes searching it both effective and time efficient, is key. It is also important to store new cases and solutions. From this the main representation problem in CBR becomes what to store in a case, how the cases content should be structured, how all the cases should be organised/indexed for effective retrieval and reuse, and how general domain knowledge can be integrated into this structure. The paper reviews two models; The dynamic memory model of Schank and Kolodner, and the category-exemplar model of Porter and Bareiss.

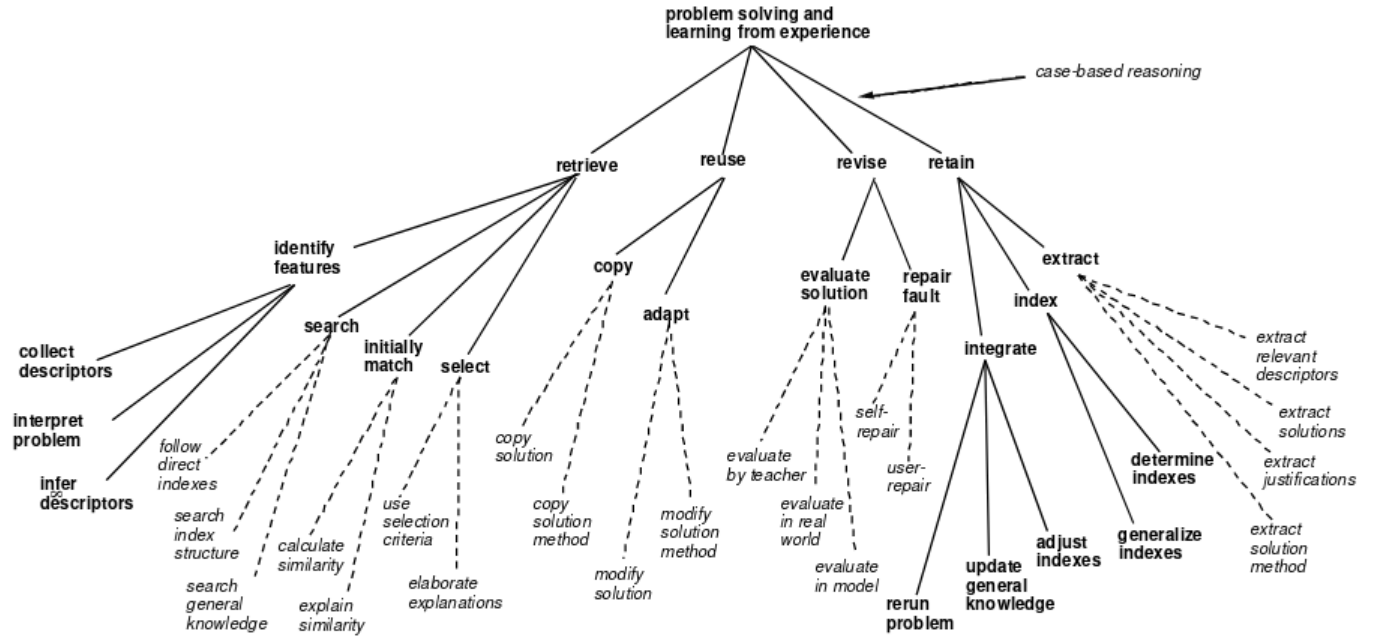


Figure 2. A task-method decomposition of CBR

## 4.1 The Dynamic Memory Model

The dynamic memory model is a hierarchical structure of episodic memory organisation packets. It organises specific cases with similar properties under a general structure (a generalized episode - GE), this episode contains three different types of objects: Norms, cases, and indices. The norms are the common features of all the cases of a GE and the indices are the distinct features for each case. And an index works like index do work, an index name and index value and points to either a GE or a specific case. Searching in the tree is done by “pushing down” the case through the network structure. The model is dynamic in the sense that when a new case is to be stored, it is again “pushed down” through the network and if it ends up with the same index as another a new GE is created. It is claimed that this model is suitable for learning both generalized and case specific knowledge and that it is a model of human reasoning and learning.

## 4.2 Category & Exemplar Model

The idea behind this model is that real world/natural concepts should be defined extensionally and that different features are assigned different importance’ to a case’s membership to a category. The case memory is embedded in a network structure of categories, semantic relations, cases and index pointers. Similar to the dynamic model a index pointer may point to either a category or a case. The indices are one of three types: feature links (pointing from features to categories), case links (pointing from categories to cases/categories) and difference links (pointing between similar cases). The cases of a category is stored according to how “typical” of that category they are (their prototypicality). Case matching (searching) is done by entering the new case’s features and returning the most prototypical cases of the matching category. Similar to searching, a new case is stored by finding a matching case and creating the appropriate feature

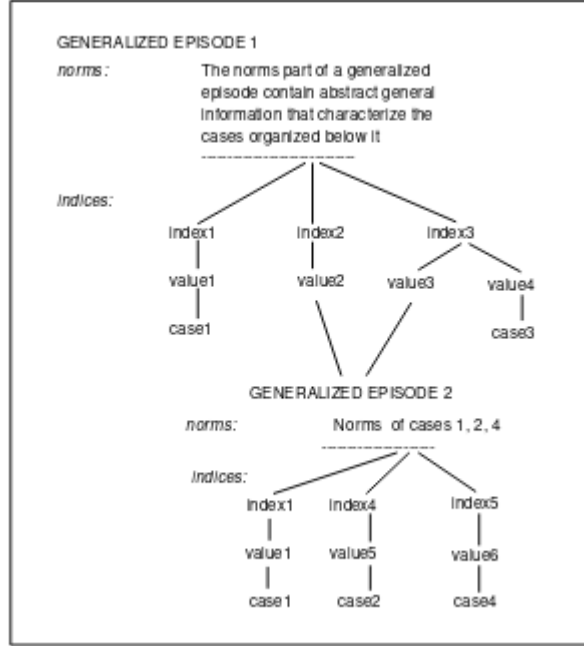


Figure 2.1: Structure of the dynamic model. It constructs a tree where each node is a GE with norms that are further separated into indexes. Each index-value pair points from one GE to either a new GE or a case.

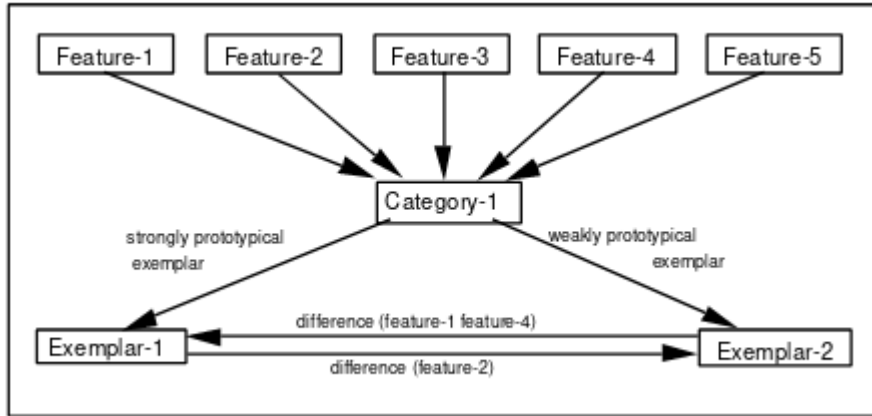


Figure 2.2: Structure of the Category & Exemplar model.

indices. If the case is sufficiently similar to an already existing case, it might not be retained at all, or merged with the matching case.

## 5 Case retrieval

The Retrive method takes a problem description and returns a best matching previous case. It is split into 4 subtasks;

- Identify Features: Identify a new case's features/problem descriptors for use to match similar previous cases.
- Initial Match: take the features found in the previous task and return a set of similar cases.
- Search: A more elaborate search to reduces the number of cases found initially.
- Select: Chooses the best solution based on the cases found in previous task.



## 6 Case Reuse

The reuse of the retrieved case solution in the context of the new case focuses on two aspects: (a) the differences among the past and the current case, and (b) what part of retrieved case can be transferred to the new case. The possible subtasks of Reuse are Copy and Adapt;

1. Copy: Just reuse the solution of a previous case, not considering (a) and (b)
2. Adapt: There are two ways to reuse past cases with adaption
  - Transformational reuse: Reuse and adapt the past case solution
  - Derivational reuse: Reuse and adapt the past *method* that constructed the past solution

## 7 Case Revision

When a solution from the reuse phase is not correct, an opportunity to learn from the failure arises. This is separated into two tasks;

- Evaluate solution: See the results of applying the solution to the environment.
- Repair fault: Detect the errors of the current solution and retrieve/generate explanations for the errors. Use these to modify(repair) the solution and retain usefull knowledge if possible. Then repeat the revision untill the solution is valid.

## 8 Case Retainment - Learning

This phase is triggered by the outcome of the case revision phase and involves selecting which information to retain, in what form it is retained, how to index the case, and how to integrate the case into the memory structure.

- Extraction: Incorporate the new case into the memory or generalize the chosen previous case by retaining information from the current case. This information being gather from the Reuse or Revision phase, as the solution method, not just the concrete solution.
- Index: The choice of indexing method influences how the cases are structured(its search space). This is a knowledge acquisition problem and should be handled in the modeling step<sup>1</sup>. For examples see Sections 4.1 and 4.2

---

<sup>1</sup>As in the system modelling step

## 3 Deep learning

A deep learning system is a neural network. Neurons<sup>1</sup> are set up in several layers, where a number of inputs are sent in to the first layer, and the last layer gives an output.

Its the best, deep learning(DL) systems have beaten records in image recognition, speech recognition, predicting the activity of potential drug molecules, analysing particle accelerator data, reconstructing brain circuits, and predicting the effects of mutations in non-coding DNA on gene expression and disease, and more.

### 1 Supervised learning

Here the system is trained using a data set of inputs and a class. It takes each input and sends it through the system, and uses an *error function*, also called *objective function* or *loss function*, to calculate the error(or distance) between the desired output and the actual output. This error is then used to modify the weights of the network, this in effect nudges the network to be more accurate when given a “similar” problem. This weight adjustment is done by calculating each weights gradient, which is a number representing that weights impact on the error/objective function. When classifying, usually the weighted sum of the outputs are used to find which category the input belongs to using thresholds.

From early on linear(also called shallow) classifiers(networks with few neurons and layer, i.e. the perceptron) could only classify simple linearly seperable problem spaces<sup>2</sup>. By using deeper(several layer) networks this problem can be solved.

### 2 Backpropagation

This is the process where the gradient(derivate) of the error function is calculated with respect to each weight. Given an example from the training set of inputs and target, backpropagation is done by first sending the inputs through the network, this is the *feed forward* step, this gives an output from the network that is compared with the target of the example. Then, working backwards from the output layer the gradient is calculated for each weight in the layer and the weight is updated. More formally the steps are:

1. Take inputs(the feature vector)  $\vec{x} = \langle x_1, x_2, ..x_d \rangle$  and the target  $y^*$ .

---

<sup>1</sup>A neuron is basicly a function(called the activation function) that takes inputs from the previous nodes in the network, toghether with weights and possible biases and the functions output(aka the neurons output) is sent to the next layer of neurons.

<sup>2</sup>A linearly seperable problem is a problem where, if the input(feature vector) and the target(correct classification of that input) is plotted on a graph(this plot is the search space/problem space) the different classes can be seperated by lines.

2. Feed the inputs through the network, using the weighted sum of the outputs of the previous layer as inputs to each activation function of the next layer until the prediction(output of the network) is calculated,  $y$ .
3. Calculate the error  $E = y - y^*$
4. backpropagate the error:
  - (a) TODO: Add backprop algorithm

A system using the algorithm presented above is called a feedforward neural network and are used to map fixed-sized inputs, to fixed-sized outputs, like categorizing an image. The nodes/neurons not in the input or output layers are called hidden nodes/neurons and are responsible for distorting the input in a way that makes the categories linearly seperable by the last layer.

### 3 Convolutional neural networks

CNN are an extention to regular neural networks that work well with multi dimentional inputs(images as 2d or 3d arrays) and can be constructed to be multidimentional(input, hidden and output layers have both a height and a width). In a convolutional layer the input of each node in the layer isn't connected to all the nodes of the previous layer, instead it takes inputs from a group of the nodes of the previous layer. There are also pooling layers, that work the same way, but instead of using an activation fuction like ReLU, it uses statistical functions, like output the max of the input values, or the average.

By using groups of inputs, the network can recognise shapes in the input, independent of location. For example, if the CNN can recognise a dog, it can recognise it anywhere in the picture, not just in the place where the dog was placed in the training set. Because of this, and the efficient use of GPUs, since 2012 CNN has taken over for older, classical computer-vision systems.

### 4 Distributed representations and language processing

### 5 Recurrent neural networks