

# Decision Theory

We all take actions to accomplish our goals, and it's our \*preferences\* that have set those goals to begin with.

There are a variety of stepping stones that we take on the journey towards our goals, some that are more instrumental to accomplishing our task than others.

**i** **Decision theory** is the art of designing rational agents that act based on their utility-based preferences and their perception of the world.

In the first-order logic we've seen previously, the results of our actions were definitive.

⚙ Recall that a **state** in first-order logic planning is composed of a set of set of conjoined fluents.

⚙ The **deterministic outcome** of a given action  $a$  from state  $s_0$  is  $s_1$ , denoted:

$$Result(s_0, a) = s_1$$

```
; In first order logic planning:
Action(Eat(Cake))
  Preconditions: Have(Cake)
  Effect: ¬Have(Cake) ∧ Eaten(Cake))

; ...where states:
s0 = Have(Cake)
s1 = [¬Have(Cake) ∧] Eaten(Cake)
; s1 need not have ¬Have(Cake) under the
; closed world assumption

; ...and so:
Result(s0, Eat(Cake)) = s1

; i.e., the result of being in state s0
; and eating your cake is that you *definitely*
; arrive in state s1
```

However, as opposed to some strong assumptions with planning that we made in past chapters, the consequences of our actions might not always be known.

In this case, where we're dealing with incomplete knowledge or stochastic outcomes, we have to take probabilities into account.

### Example

☑ Forney Industries has been developing a Self-Replicating Cake, such that while eating it, it actually has a chance to make a copy of itself! (OK, you can make the examples if you don't like this one)

This means that eating a cake does *\*not\** necessarily imply:  $\neg \text{Have}(\text{Cake})$

⚙ The **stochastic outcome** of action  $a$  from state  $s_0$  to state  $s_1$  under evidence  $e$  is given by:

$$\sum_{s' \in s_0} \text{Pr}(\text{Result}(s_0, a) = s_1 | a) * \text{Pr}(s_0 = s' | e)$$

In other words, sum over all the states of the probability of being in the initial state given the evidence  $\text{Pr}(s_0 = s' | e)$  times the probability of reaching state  $s_1$  with action  $a$ .

⚙ Often, the probability of being in the initial state,  $s_0$ , is implicit or unimportant for our reasoning, and so we'll sum out  $s_0$  and write the above as:

$$\text{Pr}(\text{Result}(a) = s_1 | a, e)$$

```
; So now, instead of:
Result(s0, a) = s1

; ...we have:
Pr(Result(a) = s1 | a, e)
```

So now we have some notion of a distribution over possible states that we can reach from taking an action under the current evidence, we need to know *\*which\** action we should take!

In other words, we need to find an action that has the highest probability of landing us in the most favorable state.

We do this by defining our agent's Utility function.

**i** A **Utility Function**  $U(s)$  assigns a single number to score the desirability of a state (the higher the number, the more desirable).

```
; Cake! Yay!
U( {Have(Cake)} ) = 100

; Oh why... why did I eat that
; whole cake? ;_
U( {Eaten(Cake)} ) = -50
```

**i** The **Expected Utility (EU)** of action  $a$  given evidence  $e$  is simply the average utility value of the possible outcomes weighted by the probability that the outcome occurs, or formally:

$$EU(a|e) = \sum_s Pr(Result(a) = s|a, e) * U(s)$$

In other words, the expected utility leverages the chance of being in a given state times the desirability of that state.

Once we know the expected utility of taking some action, if we have multiple actions to consider, we simply take the one with the highest expected utility!

Wasn't it The Rolling Stones who said, "You can't always get what you want, but if you try sometimes you might find... maximizing your expected utility to be sufficient?"

**i** The **maximum expected utility (MEU)** criterion simply stipulates that an agent will choose an action  $a$  (amongst all those that are possible  $A$ ) that has the highest expected utility amongst its choices, written:

$$ActionChoice = \operatorname{argmax}_{a \in A} EU(a|e)$$

In other words, for all action choices, choose the one that has the highest expected utility given the evidence.

### Example

☑ In the following examples, determine which action will be taken based on the maximum expected utility principle.

**Problem 1:** Andrew is deciding if he should use his day off to go to the beach, or stay inside and watch reruns of Golden Girls... Going to the beach would be a lot of fun and staying in would just be OK, but forecasts gave a 60% chance of rain for the day, which would spoil a trip to the beach.

Action	Weather	U(Action, Weather)
home	clear	2
home	raining	3
beach	clear	4
beach	raining	1

2 Compute the expected utilities of each action given by:

$$EU(a|e) = \sum_s Pr(Result(a) = s|a, e) * U(s)$$

...and determine where Andrew should spend his day off. (Click for solution)

### Example

✓ Recompute the previous example given that we observed `weather = clear` outside.

2 Click for solution.

## Preferences

So the next question you might be asking is: where do we get these utility values?

Do we need to have explicit numerical values mapped to states as indications of their desirability or can these numbers simply be relative and derived?

**i** **Preferences** are our agent's rankings of desirable states having considered their relative likelihoods of being reached.

We can use preference orderings to recover utility functions...

...but we also have to consider how desirable the state is *\*in concert with\** the probability of reaching it.

In other words, two actions A1 and A2 might both be capable of reaching a desirable state, but if one of those two actions is more *\*likely\** to reach that state, we will want to choose it over the other.

First, to formalize the possible outcomes of an action, we turn to the notion of a lottery:

**i** A **lottery** is a representation of the possible outcome states from taking some action with the probabilities of reaching each outcome, with the format:

$$L = [prob_1, outcome_1; prob_2, outcome_2; \dots; prob_n, outcome_n]$$

**i** An **outcome** is either a primitive state (i.e., a set of fluents) or, recursively, another lottery

For primitive states such as:

$$S1 = Eaten(Cake), S2 = \dots$$

We might have two lotteries:

$$L1 = [0.5, S1; 0.25, S2; 0.25, L2]$$

$$L2 = [0.3, S3; 0.6, S1; 0.1, S5]$$

...indicating that from Lottery L1, we have a 50% chance of transitioning into state S1, etc.

**i** A lottery is called **complex** if it includes an outcome that is itself another lottery.

The term "lottery" is intuitive because taking an action associated with a lottery is like buying a Lotto ticket, and hoping that you "win" the state you desired.

Now that we know how we can model outcomes of taking a particular action, let's talk about preferences.

For any two lotteries A and B, we can use the following notation to describe preferences:

- $A > B$  the agent prefers A over B
- $A \sim B$  the agent is indifferent between A and B
- $A \geq B$  the agent prefers A over B or is indifferent between them

The primary goal of utility theory is to determine how preferences between complex lotteries are related to preferences on the primitive states that compose them.

To do this, we can define some axioms on lotteries that, if violated by our intelligent systems, would lead to irrational behavior.

Axiom	Interpretation	Formalism
<b>Orderability</b>	An agent cannot avoid deciding between two actions (i.e., either one is preferred or they are equally preferable) and must assign one of the following relationships to lotteries A and B.	$(A > B), (B > A), \text{ or } (A \sim B)$
<b>Transitivity</b>	If an agent prefers lottery A to lottery B, and also prefers lottery B to lottery C, then it also prefers lottery A to lottery C.	$(A > B) \wedge (B > C) \Rightarrow (A > C)$
<b>Continuity</b>	If $A > B > C$ (i.e., some lottery B is between A and C in preference), then there is some probability p that we could find such that a certain outcome of B would be equally preferable to an outcome of A with probability p or of C with probability (1 - p)	$A > B > C \Rightarrow \exists p [p, A; 1 - p, C] \sim [1, B]$
<b>Substitutability</b>	If two lotteries are equally preferable, then you may substitute one in for the other in some other complex lottery.	$A \sim B \Rightarrow [p, A; 1 - p, C] \sim [p, B; 1 - p, C]$

Axiom	Interpretation	Formalism
<b>Monotonicity</b>	If we prefer outcome A to outcome B, then we must also prefer lotteries that have a higher probability to reach A than B.	$A > B \Rightarrow$ $(p > q \Leftrightarrow$ $[p, A; 1 - p, B] > [q, A; 1 - q, B])$
<b>Decomposability</b>	We can reduce any number of complex lotteries down to a simpler one simply by the laws of probability.	$L1 = [p, A; 1 - p, L2]$ $L2 = [q, B; 1 - q, C]$ $L1 \sim [p, A; (1 - p)q, B; (1 - p)(1 - q), C]$

❓ Substitutability tells us that two equally preferable lotteries can be substituted for one another in some other lottery. Does the following also hold?

$$(A > B) \Rightarrow [p, A; 1 - p, C] > [p, B; 1 - p, C]$$

Alright, so we have lotteries that abide by these axioms... what were we trying to do again?

Oh yeah... get some notion of what to use for utility functions...

Well, we have the following two consequences of preferential axioms that can allow us to solve for some (non-unique) utility function:

Consequence	Interpretation	Formalism
<b>Existence of Utility Function</b>	If an agent's preferences abide by the above axioms, then there exists a utility function such that: $U(A) > U(B)$ if and only if A is preferred to B, and $U(A) = U(B)$ if and only if the agent is indifferent between A and B.	$U(A) > U(B) \Leftrightarrow A > B$ $U(A) = U(B) \Leftrightarrow A \sim B$

Consequence	Interpretation	Formalism
<b>Expected Utility</b>	The utility of a lottery is the sum of the probability of each outcome times the utility of that outcome.	$EU([p_1, S_1; \dots; p_n, S_n]) = \sum_i p_i * U(S_i)$

From these two consequences, we see that, while it *does* matter what numbers we choose for our utility functions (dependent upon the scenario and how much more state  $s_0$  is desirable compared to state  $s_1$ ), there exists a utility function capable of respecting our preference ordering.

### Example

☑ Read the following preferences, observe the action lotteries, and then decide which action our intelligent system would choose based on the MEU criterion.

```
; Our agent has the following preferences,
; which abide by the 6 axioms above:
```

1.  $A > B$
2.  $B \sim C$
3.  $C \geq D$

```
; Utilities:
```

```
      A  B  C  D
U(S)  3  2  2  1
```

```
; Action 1 corresponds to lottery:
```

```
L1 = [0.2, A; 0.3, B; 0.5, L3]
```

```
; Action 2 corresponds to lottery:
```

```
L2 = [0.5, B; 0.3, C; 0.2, L3]
```

```
L3 = [0.6, D; 0.4, B]
```

```
; Which action ({1, 2}) should we take?
```

🔗 Click for solution. [Hint: Use decomposition and then compute the expected utility of each lottery]



# Multi-attribute Utility

Previously, we've dealt with states that have had an atomic utility value assigned to them, for example:

Whole state of being at home and it raining out gets assigned, statically, the utility value of 3

$$U(\text{home}, \text{rainy}) = 3$$

BUT, what if I wanted to assess \*components\* of states and treat them with different utility contributions?

$$U(\text{home}) = x, U(\text{rainy}) = y \Rightarrow U(\text{home}, \text{rainy}) = f(x, y)$$

## Example

✔ The following example from the book analyzes variables for planning airport construction.

Siting an airport (i.e., determining where to build one based on analyzed factors) requires a variety of considerations. If we're choosing between some number of land plots on which to build our airport, then we can analyze the putative values of some variables of interest based on our possible choices / actions. Let's say we were on the city planning committee of such and determined the following variables would comprise our decision criteria:

- **Cost:** the price of the land required to build upon, plus construction costs, plus the price of legal fees and processing, etc.
- **Noise:** the amount of noise generated by the airport.
- **Deaths:** possible deaths from construction hazards and other airport-related risks (propensities of nearby residents to stand on the runway? IDK...).

If we're considering a bunch of construction sites, it might not be feasible to sit down and assign a utility value to each one individually...

BUT, we might have some data on our variables of interest based on projections from studies and other sources of information, so perhaps we can construct our utility values from these facts.

That said, we often don't have deterministic information available (e.g., we don't know for certain that construction on site A will incur exactly 3 deaths), so we need to use estimation techniques.

Additionally, we don't always treat all of our variables of interest equally -- some might need to be weighted as more important than others.

❗ For some vector of variables  $X = \{X_1, X_2, \dots, X_n\}$  and their value functions  $(f_1, f_2, \dots, f_n)$ , where each  $f_i$  corresponds to a variable  $X_i$ , the utility of a state with variable instantiation:  $x = \{x_1, x_2, \dots, x_n\}$  is given by:

$$U(x_1, x_2, \dots, x_n) = F[f_1(x_1), f_2(x_2), \dots, f_n(x_n)]$$

...where  $F$  is a function that combines the value functions in a meaningful way (usually just the sum).

❗ When our metrics of interest exhibit **mutual preferential independence**, it means that a change in the value of one variable will not necessarily cause a change in the value of another variable.

⚙ If all of our metrics of interest exhibit mutual preferential independence, then our agent's choice boils down to a simple maximization of:

$$F[f_1(x_1), f_2(x_2), \dots, f_n(x_n)] = \sum_i f_i(x_i)$$

Here,  $F$  is a simple function like addition that would aggregate all of the individual variable-weighting **value functions**  $f_1, f_2, \dots, f_n$  for variable values  $x_1, x_2, \dots, x_n$ .

The idea is that we want to condense all of the variable information into a single utility value based on the importance weights of each variable.

### Example

✔ For our airport siting example, let's consider the following two sites and the value functions that provide the proper variable weighting. We'll assume our variables of interest exhibit mutual preferential independence. Determine, using the mutual preferential independence utility function above with  $F$  being simple summation, which is the superior site.

Variable	Value Function	Site 1 Value	Site 2 Value
Cost	$f_{\text{cost}}(x) = -x$	\$10,000,000	\$15,000,000
Noise	$f_{\text{noise}}(x) = -x * 1000$	200dB	180dB
Deaths	$f_{\text{deaths}}(x) = -x * 10^9$	2	1

🔍 Click for solution

📌 **Strict dominance** is the case where a state is superior on all metrics.

🔍 Site2 was the dominant choice in our example above; is it strictly dominant to Site1?

It turns out, however, that the utility cost of our estimated number of deaths far outweighs the monetary cost.

## Decision Networks

Often, however, our decisions have factors that are not necessarily the effect of anything we can control.

To model this uncertainty, we can combine Bayesian Networks, which model our knowledge of the way the world works, with the preference model we've been discussing to assess a utility value for possible actions under consideration.

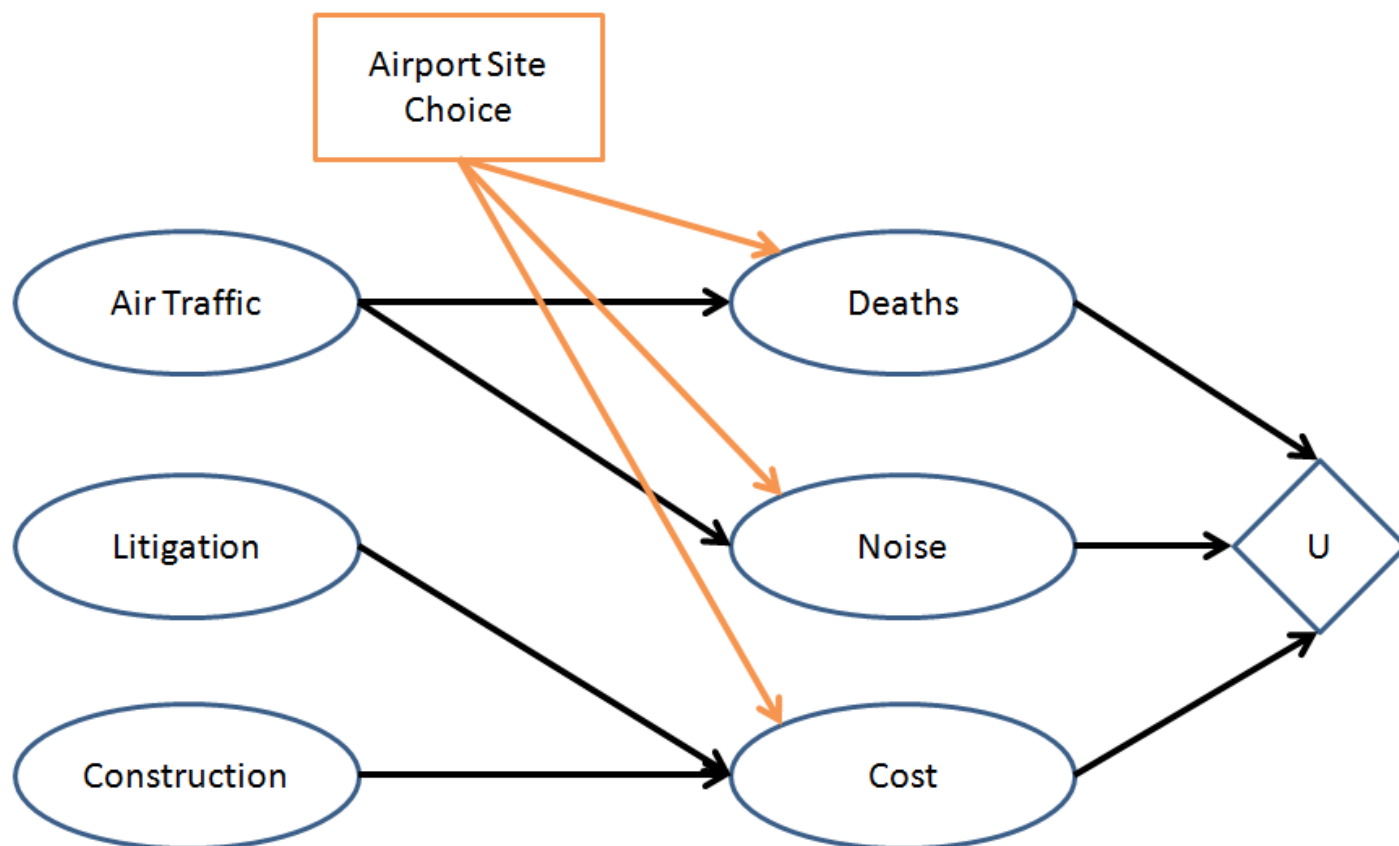
📌 A **decision network** unifies the stochastic modelling capacities of a Bayesian network with the action-choice and utility concepts of preference models by adding two additional nodes to the Bayesian network model.

📌 Decision network **chance nodes** (ovals) represent random variables, and illustrate uncertainty about the values of some of our variables (same way that Bayesian networks handled this: conditional probability tables).

📌 Decision network **decision nodes** (rectangles) represent points where the agent has a choice of actions.

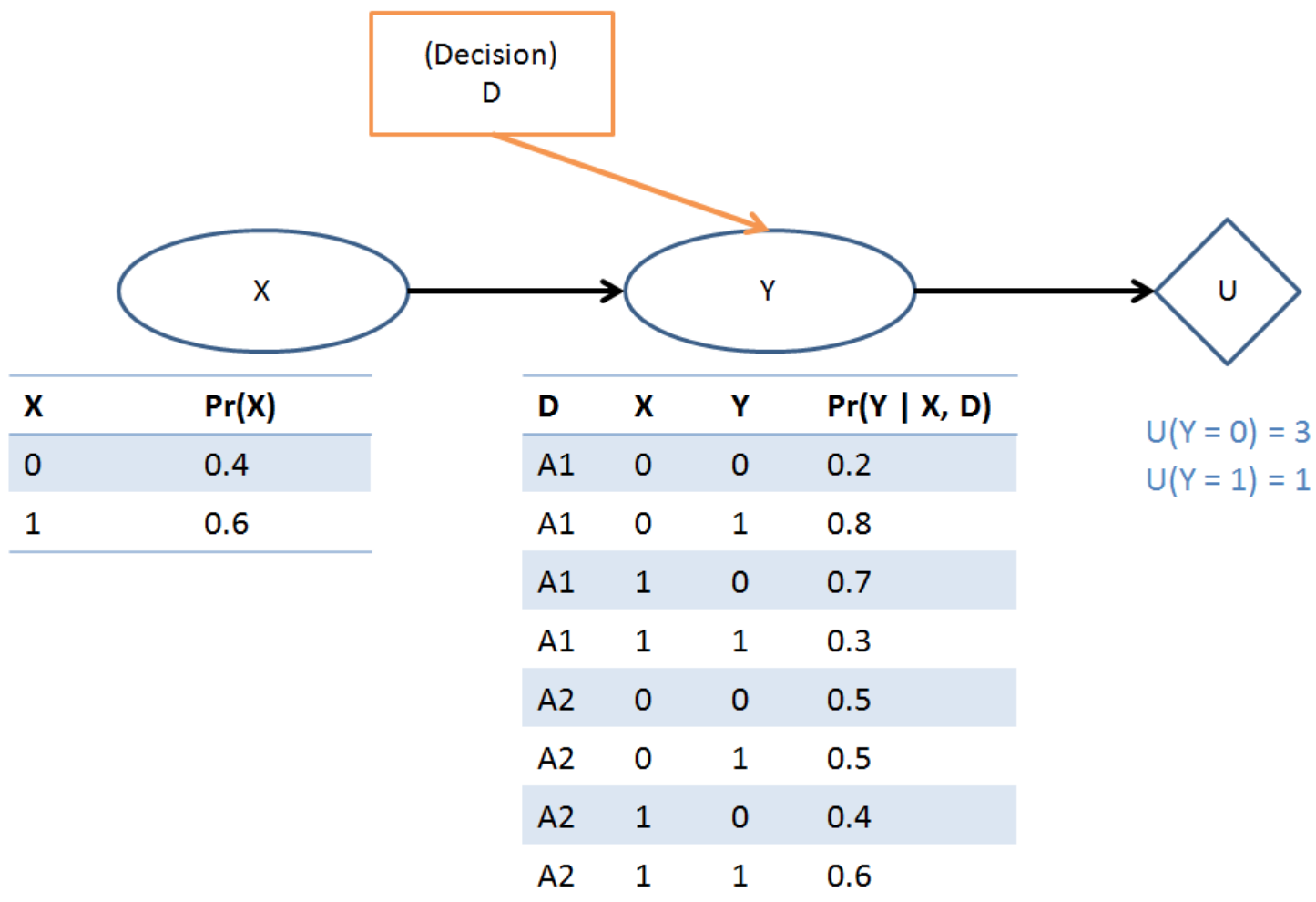
📌 Decision network **utility nodes** (diamonds) represent the agent's utility function for a given instantiation of decisions and inference.

Here is the book's example for our airport siting problem, with some additional chance nodes added in to represent our knowledge of the world:

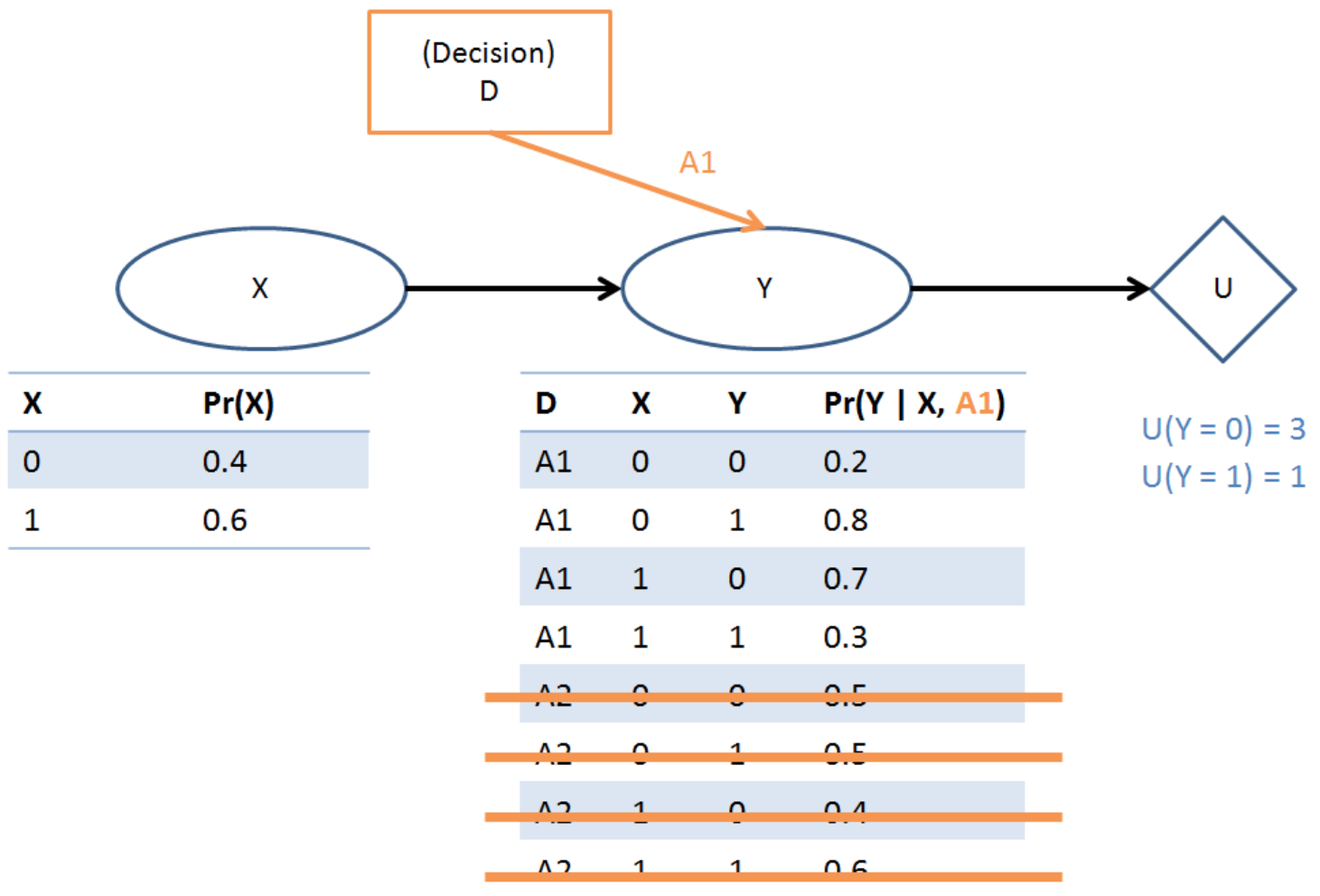


BUT, that example uses continuous variables, which we're not used to seeing... let's look at a quick example using discrete variables:

(here our decision node might have some nondeterministic action consequence)



So, we can ask what the utility of a given action is by setting the action, which causes the decision node to act like a "given" chance node:



What is the expected utility of A1 as listed above? (click for solution)

### Example

Compute the EU(A2) using the above example.

So here are some observations to make:

- **Purpose of decision networks:** judge which actions from our decision (rectangle) nodes produce the highest utility.
- **Data formats:**
  - Chance nodes are CPTs just like in Bayesian networks, except parents of chance nodes can be decision nodes as well (which are always given because we're evaluating between action choices)
  - Decision nodes for a given action simply become given chance nodes.

- Utility nodes are an evaluation of its direct causes (parents) based on supplied value functions. In our example above,  $U = F(f_{\text{death}}(\text{Deaths}), f_{\text{noise}}(\text{Noise}), f_{\text{cost}}(\text{Cost}))$
- **Inference:** based on the choice for our decision nodes, we can use an inference algorithm (like variable elimination from last week) that gives us the utility node's parents' posterior probabilities to determine not only the weighted utility of a given variable value, but also accounting for the chance that it will occur based on any evidence.

And that's a decision network in a nutshell!

Of course, all of this relies on us as humans knowing what data and variables to incorporate into our decision networks...

Our next topic will be the ability to learn how to construct a decision mechanism from raw data alone!

# Learning

The art of machine learning is little different from how we as humans acquire new information, and has particularly close analogy to how we teach children.

❶ **Machine learning** is the process by which programmers equip intelligent systems to modify their behavior based on training, observations, and other features available to the program during operation that might not have been available or convenient for the programmer.

So one of the first questions you might ask is: why give computers this adaptability at all? Shouldn't our programmers have all the tools necessary to empower their agents when they sit down to program them in the first place?

As it turns out, the answer is: not always; there are a variety of cases in which it is not possible for a programmer to explicitly define an agent's behavior.

❷ What are some examples of cases where we want our systems to learn optimal behavior rather than coding it top-down as programmers?

Learning algorithms generally suit one of a few different flavors:

Learning Class	Description	Example

<b>Unsupervised Learning</b>	Unsupervised learning is a type of pattern extrapolation engine whose typical task is <b>clustering</b> , the act of finding similarities between various input items (like images, text, etc.) and lumping them into some sort of group. In other words, it attempts to find hidden structures in unlabeled data.	Image Classification: given lots of images of faces and lots of images of computers, an unsupervised learning algorithm should be capable of lumping the faces together apart from the computers.
<b>Reinforcement Learning</b>	Just like operant conditioning (for humans), reinforcement learning provides an input for our program, on which it will make some sort of decision and be rewarded (if that was the right decision) or punished (if that was the wrong decision), figuratively speaking.	Animat Modeling: simulation study of animat populations where taking some action (like eating poison berries) has a harmful consequence that discourages that action in the future.
<b>Supervised Learning</b>	Just like having an instructor or parent correct your mistakes or reward your triumphs, supervised learning has some "oracle" that will tell you the correct answer on some problem during training so that you'll know what to change when you're wrong and what to keep when you're right. Additionally, labeled input data.	Object Recognition: given lots of images of objects with corresponding labels, e.g. a picture of a chair with label "chair", object recognition will be able to find other chairs in the future based on the label given *and* know that the particular object represents a chair (unlike unsupervised learning, which can only cluster)
<b>Semi-supervised Learning</b>	Almost exactly like supervised learning except that the labels or corrections given to us by our, now imperfect, oracle may not be entirely trustworthy. This creates extra noise that sets semi-supervised learning apart from its accurately supervised counterpart.	User Input Classification: determining age based on images of people and their self-reported age (on which they might have lied).

So, the first task of machine learning is to choose a model class that best fits the data that we have available.

What exactly constitutes our definition of "best" depends on the task at hand...

Two tools we'll discuss for some applications are decision trees and Bayesian networks.

## Classification

One of those most common machine learning tasks involves classifying different items evaluated on their variable settings and sorting them into the proper group.



**i** **Classification** attempts to learn a **model** from a collection of examples that will predict a classification given the observable attributes.

Just what model we should use is dependent on the question we're asking, i.e., the classification we're trying to make, and the data that we have available.

So, let's look at some example classification problems, and the models appropriate for them!

# Decision Trees

One type of classification is making a decision:

Given some input features about a situation, we have to evaluate whether to make one decision or another.

To model this task, we might consider a decision tree... but first...

Looking at a simplified version of the book's example:

## Example

### Dining Dilemma

Forney Industries has ~~recently contracted with the NSA~~ decided to branch out into the phone app business, and is debuting with its first app that learns your dining habits and can predict whether or not you will wait for a seat at a restaurant based on a number of attributes. The app will therefore generate a "yes" or "no" conclusion of whether or not you are going to wait based on the following inputs:

- **Patrons?** A measure of how many diners are currently at the restaurant, can be: {none, some, many}
- **Hungry?** Whether or not you are currently on the brink of starvation (ok, maybe that's too dramatic... you just haven't eaten since lunch), and can be: {yes, no}
- **Type?** The type of food the restaurant serves, which can be: {French, Italian, Thai, or Burgers} (you know... from the country of Burger?)
- **Fri / Sat?** Whether or not it is a peak-serving day of Friday or Saturday, can be: {yes, no}.

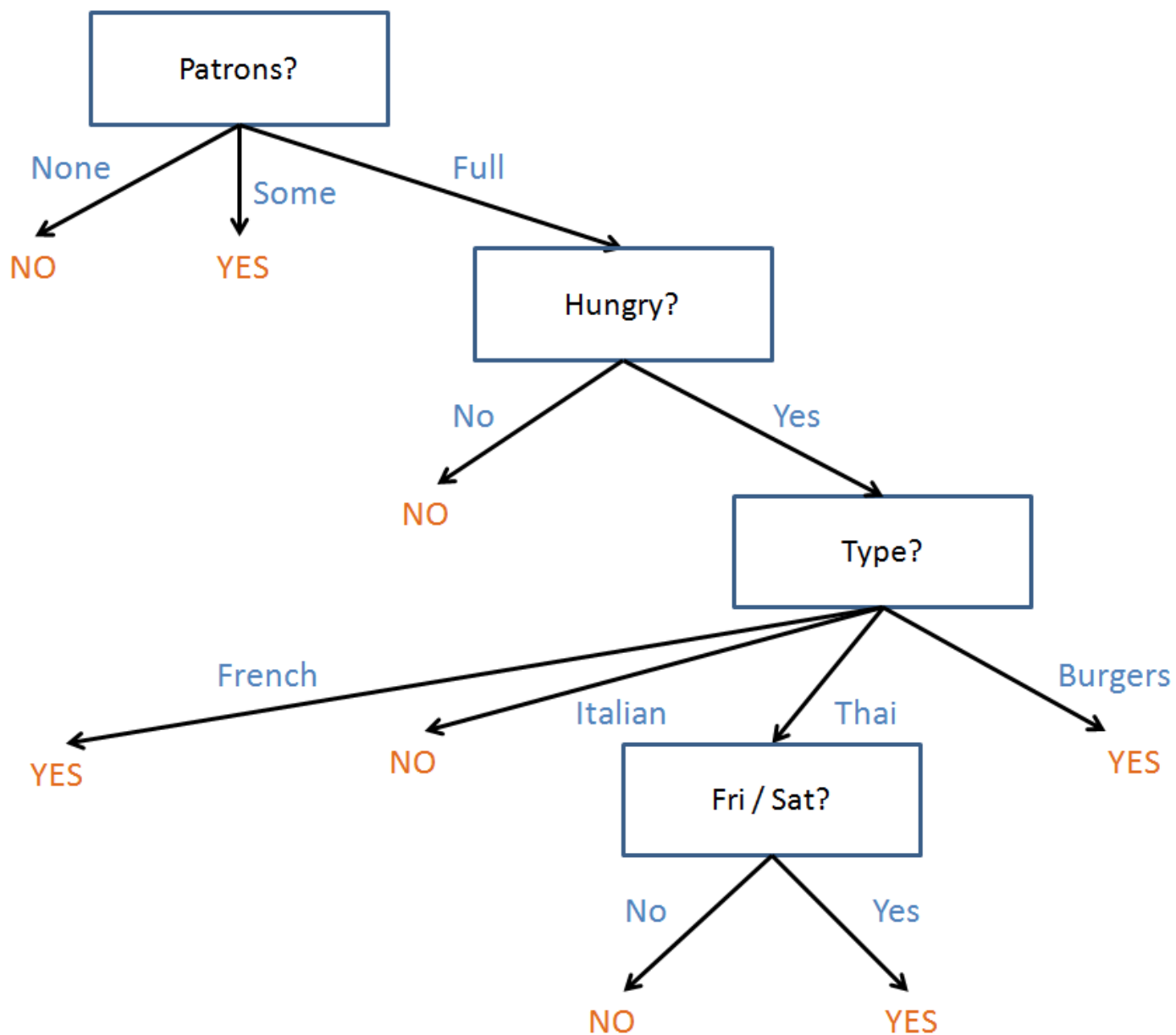
Our app has been learning your habits for awhile now and made the following determinations (expressed here just semantically):

- If there are no patrons you don't trust the restaurant quality and so will not wait, but if it's full, you'll only wait based on other attributes:
- If you're not really hungry, you're not going to wait, otherwise:
- You love French food and Burgers, so you'll wait at this point if it's either of those... but if it's Thai...
- You'll only wait if it's not a peak dining day of Friday or Saturday.

So if that's the plain-English description of our habits, how does an intelligent system predict whether or not we'll wait?

**i** A **decision tree (d-tree)** is a tree structure with nodes as attributes and leaves as deterministic classification outcomes.

The d-tree of our above observations might look like:



It's easy to see what a decision tree will classify for our decision; we simply start at the root and trace the edges until we end up at a leaf!

#### Example

☑ What will our above decision tree give us for the following samples?

Patrons?	Hungry?	Type?	Cost?	Fri / Sat?	Wait?
some	yes	burgers	\$\$	yes	???

Patrons?	Hungry?	Type?	Cost?	Fri / Sat?	Wait?
full	yes	italian	\$\$\$	no	???
full	yes	thai	\$	yes	???

You might notice: there was another attribute included in our sample data (Cost?) that was never used by our d-tree. That's fine, and sometimes preferable, as we'll soon find out...

So now that we see how a d-tree works... let's talk about how to learn them in the first place!

## Learning a D-Tree

To talk about learning d-trees, we imagine we have some large amount of data on which we can extract patterns to perform the learning aspect.

❗ A **training set** is a list of input / output pairings where, given the input characteristics of a particular sample, we tell our learning system what the expected outcome should be (under the assumption that it will be able to formulate a classification function from lots of examples).

So let's take a look at an arbitrary training set (Example credit to Evan Lloyd; it was too beautiful not to use)

### Example

☑ Imagine that we have arbitrary attributes A, B, C, and D. We must make a yes / no decision for some classification X. Observe the following training set over each feature and the expected decision for X.

A	B	C	D	X
grn	sml	1	0	<u>yes</u>
grn	sml	0	3	<u>yes</u>
red	med	0	5	<u>yes</u>
blu	med	0	5	<u>no</u>
grn	med	1	4	<u>no</u>
grn	lrg	1	1	<u>yes</u>
red	lrg	0	4	<u>yes</u>
blu	med	0	2	<u>no</u>
blu	lrg	1	4	<u>no</u>
blu	med	0	3	<u>no</u>
red	med	0	3	<u>yes</u>
grn	lrg	0	5	<u>yes</u>
grn	med	1	1	<u>no</u>
red	sml	1	2	<u>yes</u>
grn	lrg	1	3	<u>yes</u>

A	B	C	D	X
red	sml	0	1	<u>yes</u>
blu	sml	1	0	<u>no</u>
grn	lrg	1	2	<u>yes</u>
blu	med	1	4	<u>no</u>
grn	med	1	5	<u>no</u>
grn	med	0	2	<u>no</u>
red	sml	0	1	<u>yes</u>
grn	med	0	5	<u>no</u>
blu	med	1	1	<u>no</u>
red	sml	1	0	<u>yes</u>
blu	med	1	2	<u>no</u>
grn	lrg	1	1	<u>yes</u>
grn	lrg	1	1	<u>yes</u>
blu	lrg	1	0	<u>no</u>
grn	med	0	3	<u>no</u>

❓ What form of learning is this? Unsupervised, semi-supervised, or supervised?

Now, we want to find some decision tree that accurately gives the correct classification based on attributes A, B, C, and D.

But which one do we choose?

⚙️ A **single-node split** on a training set tries to maximize the classification accuracy based on a single attribute's values.

**i** In classification, a **hit** occurs whenever our split on attributes agrees with the training set, and a **miss** occurs when it disagrees.

When we split on a value of a variable, e.g.,  $A = \text{red}$ , we look at how many samples got classified to  $X = \text{yes}$  and  $X = \text{no}$  in the training set and try to make the classification that is most accurate (most hits / fewest misses) based only on  $A = \text{red}$ .

So... let's try splitting on  $A$  first. This would give us:

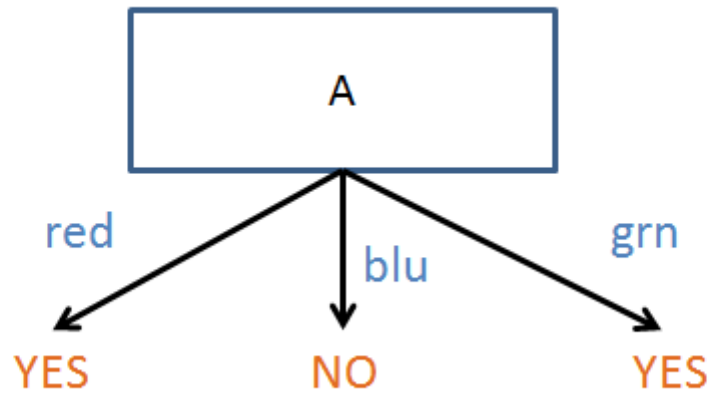
	$A = \text{red}$	$A = \text{blu}$	$A = \text{grn}$
$X = \text{yes}$	7	0	8
$X = \text{no}$	0	9	6

**?** Given this break down, what would be the most accurate classification to give the cases where  $A = \text{red}$ ?

**?** Given this break down, what would be the most accurate classification to give the cases where  $A = \text{blu}$ ?

**?** Given this break down, what would be the most accurate classification to give the cases where  $A = \text{grn}$ ?

This would give us a single node d-tree looking like:



So, using this single node split, we would achieve an 80% accuracy since 6 / 30 samples are misclassified by saying that  $X = \text{yes}$  whenever  $A = \text{grn}$

A	B	C	D	X
grn	sml	1	0	<u>yes</u>
grn	sml	0	3	<u>yes</u>
red	med	0	5	<u>yes</u>
blu	med	0	5	<u>no</u>
grn	med	1	4	<u>no</u>
grn	lrg	1	1	<u>yes</u>
red	lrg	0	4	<u>yes</u>
blu	med	0	2	<u>no</u>
blu	lrg	1	4	<u>no</u>
blu	med	0	3	<u>no</u>
red	med	0	3	<u>yes</u>
grn	lrg	0	5	<u>yes</u>
grn	med	1	1	<u>no</u>
red	sml	1	2	<u>yes</u>
grn	lrg	1	3	<u>yes</u>

A	B	C	D	X
red	sml	0	1	<u>yes</u>
blu	sml	1	0	<u>no</u>
grn	lrg	1	2	<u>yes</u>
blu	med	1	4	<u>no</u>
grn	med	1	5	<u>no</u>
grn	med	0	2	<u>no</u>
red	sml	0	1	<u>yes</u>
grn	med	0	5	<u>no</u>
blu	med	1	1	<u>no</u>
red	sml	1	0	<u>yes</u>
blu	med	1	2	<u>no</u>
grn	lrg	1	1	<u>yes</u>
grn	lrg	1	1	<u>yes</u>
blu	lrg	1	0	<u>no</u>
grn	med	0	3	<u>no</u>

Alright, so let's see if doing a single node split on another attribute can do better than 80%.

Shall we try to split on B?

	B = sml	B = med	B = lrg
X = yes	6	2	7
X = no	1	12	2



We actually do a bit better with the single var split on B if we choose to again use the plurality rule of classification, giving us only 5 / 30 misclassifications:

A	B	C	D	X
grn	sml	1	0	<u>yes</u>
grn	sml	0	3	<u>yes</u>
red	med	0	5	<u>yes</u>
blu	med	0	5	<u>no</u>
grn	med	1	4	<u>no</u>
grn	lrg	1	1	<u>yes</u>
red	lrg	0	4	<u>yes</u>
blu	med	0	2	<u>no</u>
blu	lrg	1	4	<u>no</u>
blu	med	0	3	<u>no</u>
red	med	0	3	<u>yes</u>
grn	lrg	0	5	<u>yes</u>
grn	med	1	1	<u>no</u>
red	sml	1	2	<u>yes</u>
grn	lrg	1	3	<u>yes</u>

A	B	C	D	X
red	sml	0	1	<u>yes</u>
blu	sml	1	0	<u>no</u>
grn	lrg	1	2	<u>yes</u>
blu	med	1	4	<u>no</u>
grn	med	1	5	<u>no</u>
grn	med	0	2	<u>no</u>
red	sml	0	1	<u>yes</u>
grn	med	0	5	<u>no</u>
blu	med	1	1	<u>no</u>
red	sml	1	0	<u>yes</u>
blu	med	1	2	<u>no</u>
grn	lrg	1	1	<u>yes</u>
grn	lrg	1	1	<u>yes</u>
blu	lrg	1	0	<u>no</u>
grn	med	0	3	<u>no</u>

But, as we've seen with our restaurant waiting example, we're not always interested in a single-node split on an attribute, but a multi-node one... why waste all that extra info?!

Now, of course, we have the question of which node to split on first... it's clear that order matters in deciding the most accurate classification!

**❶ Entropy** is a measure of uncertainty of a random variable and is the fundamental unit of information theory. Formally, for variable  $V$ , the entropy  $H(V)$  is given by:

$$H(V) = - \sum_{v \in V} Pr(v) * \log_2[Pr(v)]$$

In other words, the entropy is the sum, for all values of  $V$ , of the probability of seeing that value  $v$  times  $\log_2$  of that probability.

⚙ Intuitively, entropy weights the rarity of seeing a particular value of a variable by the same uncertainty its probability represents. So, the more, equally likely values a variable has, the higher its entropy's magnitude. The fewer, more certainly occurring values it has, the lower its entropy's magnitude.

Some things to note about entropy:

- It's measured in the bits it would require in order to represent all possible equally likely outcomes (manifest in the  $\log_2$  in the equation)
- A higher entropy means it's closer in value to a uniform distribution, like a coin flip (i.e., the higher the value, the more random)
- The lower the entropy, the more certain the outcome is.

❓ Find the entropy of a fair coin flip with  $V = \{heads, tails\}$ . Click for solution.

So, we see that a coin flip has 1 bit of entropy, since it only takes 1 bit (i.e., a 0 or 1 outcome uniformly distributed) to model the uncertainty.

Our goal in building decision trees, however, is to MINIMIZE uncertainty.

❗ In the classification problem, **Information gain** is the expected reduction in entropy from splitting our decision at a choice point on some attribute. Therefore, to maximize the information gain, we may equivalently minimize the entropy (uncertainty).

❗ We can model this in terms of **expected entropy**, which provides a metric for the entropy reducing contributions of each variable value. Computing the expected entropy reduction for splitting on a variable is a two step process described below.

Let  $V$  = the variable / feature being analyzed

$v$  = a value of  $V$

$p$  = the positive classifications ( $X = \text{YES}$ ) for  $V = v$

$n$  = the negative classifications ( $X = \text{NO}$ ) for  $V = v$

**Step One:** The entropy for  $V = v$  is given by  $H_v$  and is a function of the positive and negative classifications for that value:

$$H_v(p, n) = - \frac{p}{p+n} * \log_2 \frac{p}{p+n} - \frac{n}{p+n} * \log_2 \frac{n}{p+n}$$

**Step Two:** The expected entropy for splitting on  $V$  is defined as:

$$EH(V) = \sum_{v \in V} Pr(v) * H_v(p, n)$$

In other words, for all values  $v$  in variable  $V$ , sum up their entropies split on positive and negative classifications.

### Example

☑ Compute the expected entropy of splitting on variable  $A$  in our example from before (table replicated below for ease):

$n = 30$	$A = \text{red}$	$A = \text{blu}$	$A = \text{grn}$
$X = \text{yes (p)}$	7	0	8
$X = \text{no (n)}$	0	9	6

🔗 **Step Zero - Setup:** Compute the probabilities of seeing each variable value from the table and write the equation for the expected entropy of a split on  $A$ . (Click for solution)

🔗 **Step One - Entropies:** Compute the entropy of each variable value across the positive and negative trials in our training set. (Click for solution)

⚙️ **Step Two - Expected Entropy:** plug in relevant parts from step 1 into our goal from step 0 and solve.

Whew! What a pain! I hope I never have to do that again...

### Example

✔️ What is the expected entropy for splitting on attribute B in our example? Table replicated below for convenience.

n = 30	B = sml	B = med	B = lrg
X = yes (p)	6	2	7
X = no (n)	1	12	2

⚙️ Click for brief solution

Conclusion:

$$EH(B) = 0.6434 > EH(A) = 0.4598$$

$$\text{Accuracy}(B) = 86\% > \text{Accuracy}(A) = 80\%$$

Therefore, we see that splitting on A actually has less uncertainty associated with it than splitting on B, even though A is less accurate by itself.

(Recall that splitting on A had 6 misses compared to splitting on B which had 5).

So how does knowing the expected entropy of A help us? This brings us to the next question:

Since A is accurate for A = red and A = blu, what if we then split its inaccurate component (i.e., A = grn) on B? How would we decide to do that?

⚙️ The **decision tree algorithm** provides means of automating the attribute splitting to minimize expected entropy across multiple variable splits.

### "Naive" Decision Tree Algorithm

1. Pick an attribute A with the minimal expected entropy, and set that as the root
2. For each **value** of A, if there are attributes remaining on which to split, and classification is not perfect, then recurse on each subtree with A removed.
3. Otherwise, there were no more attributes on which to split, so simply maximize your accuracy by classifying with a plurality vote (i.e., satisfy the most positives or negatives even though you won't be perfect).

### Example

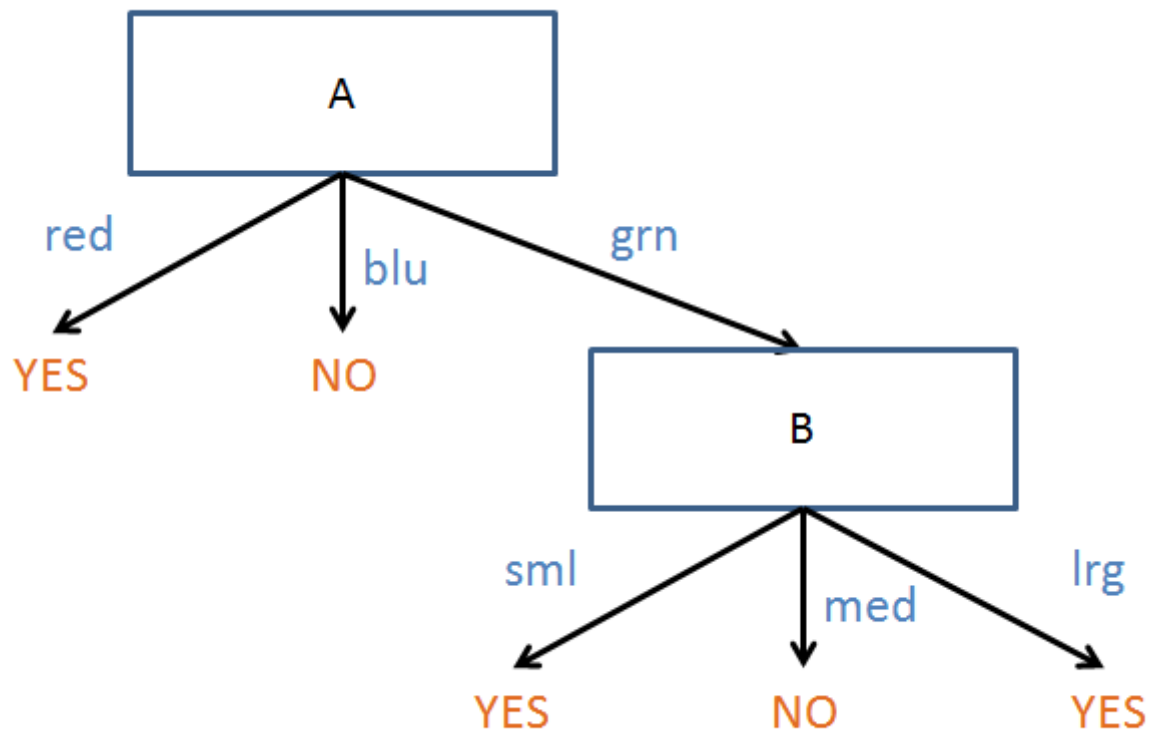
✔ Use the decision tree algorithm on our above example to perform a multi-variable split on the classification task.

1. Choose attribute A for the root since it has the minimal entropy.
2. We note that for A = red and A = blue, we have perfect accuracy (and therefore need not recurse on those subtrees), but we are not perfect for A = grn.
3. Recurse on all of the samples where A = grn and try to find a new attribute to split. In this case, splitting on B gives us a perfect split!

A	B	C	D	X
grn	sml	1	0	<u>yes</u>
grn	sml	0	3	<u>yes</u>
red	med	0	5	<u>yes</u>
blu	med	0	5	<u>no</u>
grn	med	1	4	<u>no</u>
grn	lrg	1	1	<u>yes</u>
red	lrg	0	4	<u>yes</u>
blu	med	0	2	<u>no</u>
blu	lrg	1	4	<u>no</u>
blu	med	0	3	<u>no</u>
red	med	0	3	<u>yes</u>
grn	lrg	0	5	<u>yes</u>
grn	med	1	1	<u>no</u>
red	sml	1	2	<u>yes</u>
grn	lrg	1	3	<u>yes</u>

A	B	C	D	X
red	sml	0	1	<u>yes</u>
blu	sml	1	0	<u>no</u>
grn	lrg	1	2	<u>yes</u>
blu	med	1	4	<u>no</u>
grn	med	1	5	<u>no</u>
grn	med	0	2	<u>no</u>
red	sml	0	1	<u>yes</u>
grn	med	0	5	<u>no</u>
blu	med	1	1	<u>no</u>
red	sml	1	0	<u>yes</u>
blu	med	1	2	<u>no</u>
grn	lrg	1	1	<u>yes</u>
grn	lrg	1	1	<u>yes</u>
blu	lrg	1	0	<u>no</u>
grn	med	0	3	<u>no</u>

On the 14 remaining cases where A = grn, splitting on B actually gives us a perfect fit! Let's see the decision tree that results:



And that's how you do decision trees!

---

## PAC

Sometimes, answers are not so clean-cut, and a descent down a decision tree is insufficient for classification.

So, we'll return to our old friend the Bayesian network in order to help with classification; let's look at our motivating example:

### Example

☑ GreeterBot 4000

Forney Industries is developing its most ambitious robotic project yet: GreeterBot 4000. GreeterBot 4000 needs to be able to walk around department stores greeting and answering questions of the customers. You've been tasked with programming GreeterBot 4000's ~~KillAllHumans~~ IdentifyHuman protocol, which must determine if some object in its environment is a human or not.

The main issue is that there are a variety of mannequins throughout the store that, although resembling humans, should not be greeted.

That said, GreeterBot 4000 is already equipped with input sensors that can assess the following attributes:

- **Movement (M):** we can assess whether an object is moving at, say, 3 discrete velocities: {not, slow, fast}
- **Height (H):** we can assess the height of an object at, say, 3 discrete heights relative to humans: {sml, med, lrg}
- **Speaking (S):** we have audio sensors capable of determining if some object is speaking a language: {0, 1}
- **Form (F):** we have some scoring algorithm that returns a score from 0 - 5 based on whether an object has human parts like a head, torso, arms, etc.

The goal, then, is to determine for each object we encounter assessed with the above 4 attributes, whether or not that object is a human.

(ignoring, of course, that this is just a re-labelling of the previous example with different semantics...)

**⚠ Complications:** our classification might not be so simple as splitting on a single attribute, some confounding cases might be:

- A loudspeaker might announce the latest deals and be perceived as "speaking" even though it's not a human.
- A mannequin might have a good human form, and the right height, even though it's not human.
- A true human might be any height, have good human form (if our visual sensors aren't occluded), and simply not be speaking at a given moment.

So, needless to say, there are a variety of different cases that need to be handled, but different attributes confer some amount of human-ness.

Because of the amount of noise implicit in this problem, however, we may not be able to get away with clean classification divisions like with decision trees.

**📍 Probably Approximately Correct (PAC)** classification strategies are those that, given a sufficiently large training set, can make a decision based on the most likely outcome of any witnessed evidence.



Frankly, I think the term "Probably Approximately Correct" sounds like a rationalizing scientist who isn't quite sure of his invention...

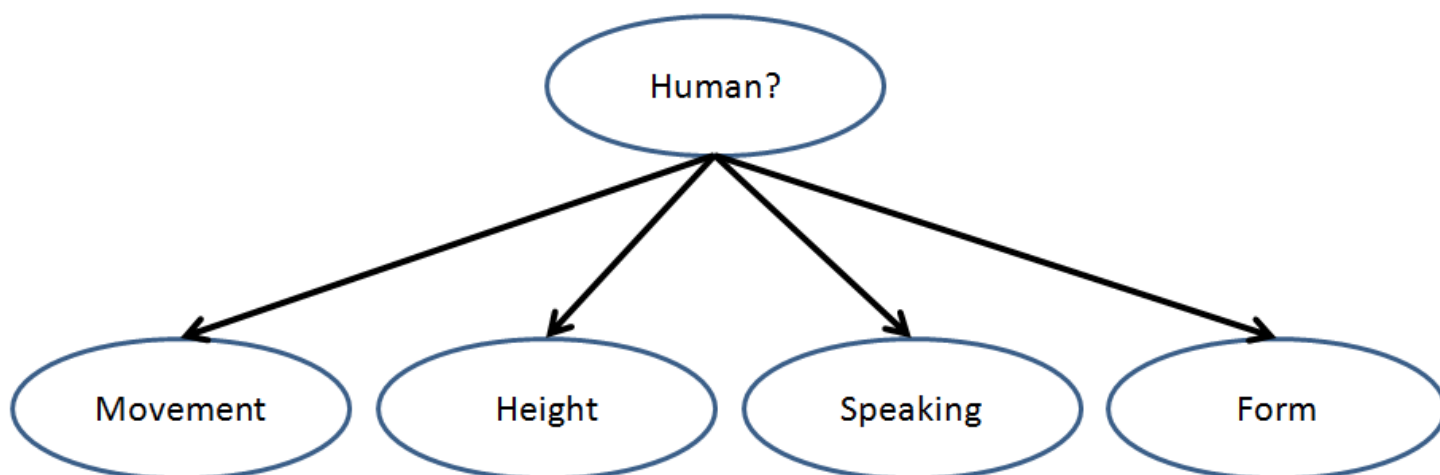


PAC strategies are good when we don't need a completely accurate answer where all cases are considered, but we have a good gist for when to say "Yes" and when to say "No"

❗ The **Naive Bayes Classifier** is a PAC classification structure where we consider some class  $C$  to be the cause of witnessing some number of identifiers, or effects,  $F_i$ . It is called "Naive" because given the class, we assume that all of the indicators are conditionally independent.

The Naive Bayes Classifier is a Bayesian network where we have a single cause (the class) being the reason for certain indications of that class.

For our example, the Naive Bayes structure might look like this:



This classifier says, "If you are a human, then you'll probably exhibit some movement, a certain height distribution, the propensity for speech, and some amount of human form."

Furthermore, if we know you're human, then the indicators are all conditionally independent of one another because knowing that a human can speak tells us nothing more about a human's ability to move.

Being a Bayesian network, these indicating propensities are all implicit within the network's conditional probability tables.

⚙ Naive Bayes classification then simply asks: having observed some object with indicators  $F_1, F_2, \dots, F_n$ , which classification (i.e., value of the class  $C = c$ ) is most likely?

Let's look at the probability statements of interest:

We are interested in determining which classification in some set of class variable ( $C$ ) values is most likely, having witnessed some object with indicators  $F = f$

So, we are interested in the table:

$$Pr(C|f_1, f_2, \dots, f_n)$$

...to see which value of  $C$  is most probable, we would then find the value for  $C$  for which:

$$\operatorname{argmax}_{c \in C} Pr(c|f_1, f_2, \dots, f_n)$$

So, for our example, let's say we encountered an object exhibiting the following:

Movement	Height	Speaking	Form
slow	med	1	4

Here we have some object that is moving slowly, of moderate height, that's speaking, and closely resembles the human form.

Chances are good that it's a human!

Let's derive what we'll need to make that determination.

Note that  $f_1, f_2, \dots, f_n$  is our evidence so we can simply use a familiar theorem to find our answer:

$$Pr(C|f_1, f_2, \dots, f_n) = \frac{Pr(f_1, f_2, \dots, f_n|C) * Pr(C)}{Pr(f_1, f_2, \dots, f_n)}$$

Notice that  $Pr(C)$  is the prior on  $C$ , and since our evidence is conditionally independent given  $C$ , we can write the chain-rule factorization to achieve:

$$= \frac{Pr(f_1|C) * Pr(f_2|C) * \dots * Pr(f_n|C) * Pr(C)}{Pr(f_1, f_2, \dots, f_n)}$$

So, in our example, we would have each indicator  $f_i$  be a different attribute about our witnessed object.

So our example's quantity of interest is (for  $Hu?$  = Human?,  $M$  = movement,  $H$  = height,  $S$  = speech,  $F$  = form):

$$\frac{Pr(Hu?|M = slow, H = med, S = 1, F = 4) = Pr(M = slow|Hu?) * Pr(H = med|Hu?) * Pr(S = 1|Hu?) * Pr(F = 4|Hu?) * Pr(Hu?)}{Pr(M = slow, H = med, S = 1, F = 4)}$$

Almost there! Now we just "need" one more element to complete our classification:

Now all we need to find out is the  $Pr(f_1, f_2, \dots, f_n)$ . We have a strategy to do this: Case analysis!

$$\begin{aligned} Pr(f_1, f_2, \dots, f_n) &= \sum_{c \in C} Pr(f_1, f_2, \dots, f_n|c) * Pr(c) \\ &= \sum_{c \in C} Pr(f_1|c) * Pr(f_2|c) * \dots * Pr(f_n|c) * Pr(c) \end{aligned}$$

And so, substituting back into our original equation:

$$\frac{Pr(f_1|c) * Pr(f_2|c) * \dots * Pr(f_n|c) * Pr(C)}{Pr(f_1, f_2, \dots, f_n)} = \frac{Pr(f_1|C) * Pr(f_2|C) * \dots * Pr(f_n|C) * Pr(C)}{\sum_{c \in C} Pr(f_1|c) * Pr(f_2|c) * \dots * Pr(f_n|c) * Pr(c)}$$

Notice: This gives us a \*table\* since we have the class variable C in the numerator; we really want to determine \*which\* value of C is most likely.

We also notice that we don't need the denominator to make this determination since it's simply the normalizing constant!

So, we simply compare the different values of C and classify based on whichever is most probable.

This outcome is defined by:

The most likely class for all classes in C will be represented by:

$$\operatorname{argmax}_{c' \in C} \frac{Pr(f_1|c') * Pr(f_2|c') * \dots * Pr(f_n|c') * Pr(c')}{\sum_{c \in C} Pr(f_1|c) * Pr(f_2|c) * \dots * Pr(f_n|c) * Pr(c)}$$

Noting that the division is simply the normalizing constant:

$$\propto \operatorname{argmax}_{c' \in C} Pr(f_1|c') * Pr(f_2|c') * \dots * Pr(f_n|c') * Pr(c')$$

...which looks intimidating, but just says, "Choose the value of class variable C that has the highest probability given the witnessed indicator values for the object"

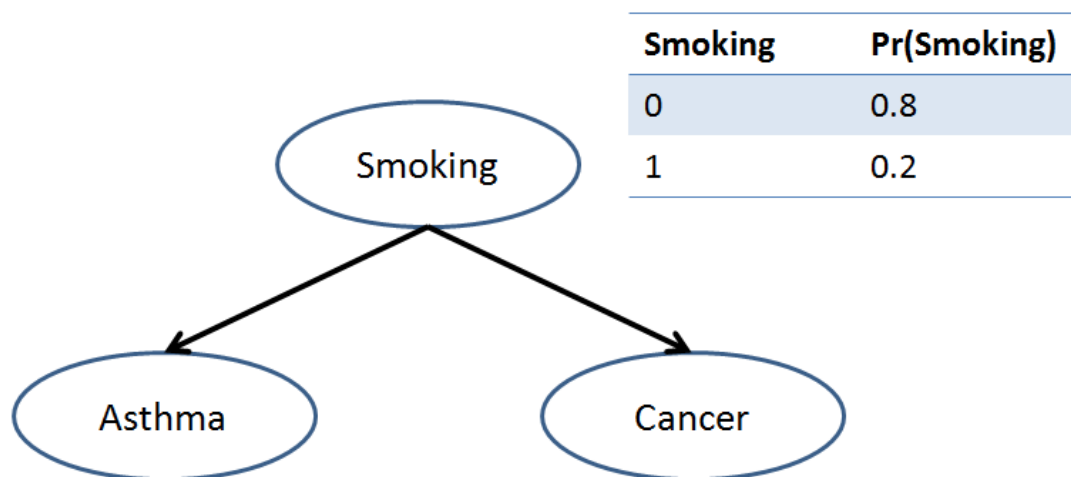
In our example, we would be choosing between: C = Human? = {yes, no}

### Example

☑ Recall our example Bayesian network from last week (which is actually a Naive Bayes structure with Smoking as the class variable). Suppose that we observe:

$$\text{Asthma} = 1, \text{Cancer} = 0$$

Is it more likely that this patient is a smoker or non-smoker?



Smoking	Asthma	Pr(Asthma   Smoking)
0	0	0.9
0	1	0.1
1	0	0.4
1	1	0.6

Smoking	Cancer	Pr(Cancer   Smoking)
0	0	0.8
0	1	0.2
1	0	0.1
1	1	0.9

[Click for solution.](#)

## Final

Your final is fast approaching! Engage panic!

Just a quick note that I've prepared a massive practice final that you can find in my course site's materials section for your perusal.

If you want any type of problem to be particularly represented, let me know and I'll be sure to include an instance!