

CS 161 Intro. To Artificial Intelligence

Week 9, Discussion 1C

Li Cheng Lan



Today's Topics

- **BN Learning and Inference**

- Learning parameters
- Learning BN structure
- Model-oriented vs. Query-oriented learning

- **Decision Tree and Random Forest**

- Entropy and Conditional Entropy
- Classifier
- BN Classifier

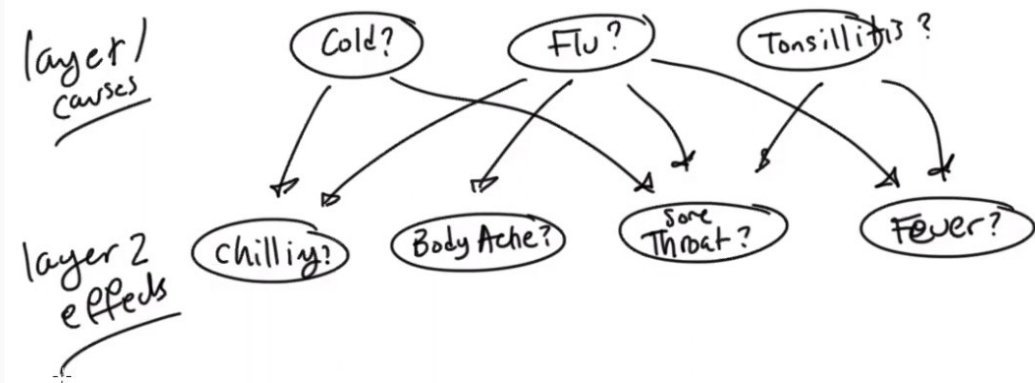
BN Learning - Example

Example

The flu is an acute disease characterized by fever, body aches and pains, and can be associated with chilling and a sore throat. The cold is a bodily disorder popularly associated with chilling and can cause a sore throat. Tonsillitis is inflammation of the tonsils which leads to a sore throat and can be associated with fever.

- Ways to learn parameters (probabilities) and construct CPTs (conditional probability table):

- Problem statement
- Subjective beliefs
- **Learning from data**



CPTs can also be estimated from medical records of previous patients

Case	Cold?	Flu?	Tonsillitis?	Chilling?	Bodyache?	Sorethroat?	Fever?
1	true	false	?	true	false	false	false
2	false	true	false	true	true	false	true
3	?	?	true	false	?	true	false
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

*CPT₁ → BN₁ → find prob. of all examples → score₁
→ score₂*

- Complete data (e.g. 2nd case): if **every** row is complete, dataset is complete → efficient
- Incomplete data (e.g. 1st case): if **any** row is incomplete, dataset is incomplete
 - Use algorithm such as expectation maximization (EM) to find maximum likelihood parameters
- **BN structure + CPTs = BN**
 - May come out multiple BNs → The BN with max score (computed by multiplying prob. assigned to all cases based on each BN) is better
 - It's called maximum likelihood principle

Learning Parameters

Goal: Estimate/Learn parameters in CPTs

- E.g. $\theta_{S|h}$, $\theta_{\neg S|h}$, $\theta_{S|\neg h}$, $\theta_{\neg S|\neg h}$ in the CPT of node S

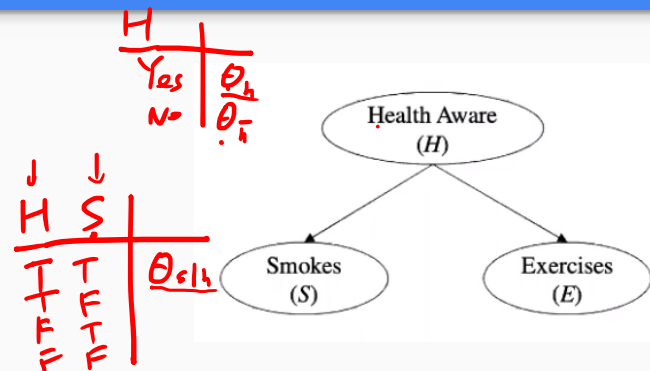
Case 1: using **complete** data

- Construct empirical distribution table from data
- Compute parameters using conditional probability based on empirical distribution

E.g. $\theta_{S|h} = \Pr(s|h) = \frac{\Pr(s \wedge h)}{\Pr(h)}$

Handwritten calculation: $\frac{2}{16} \rightarrow \frac{12}{16} = \frac{1}{6}$

dataset:



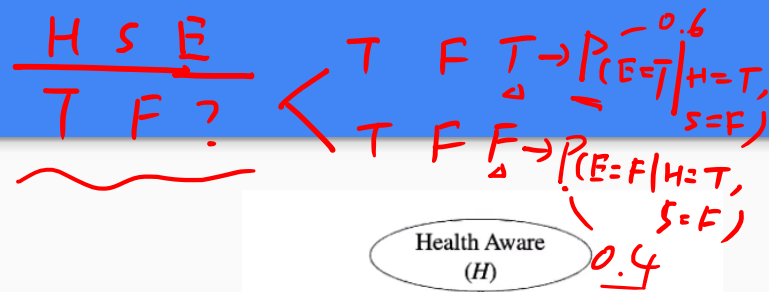
Case	H	S	E		
1	T	F	T		
2	T	F	T		
3	F	T	F		
4	F	F	T		
5	T	F	F		
6	T	F	T		
7	F	F	F		
8	T	F	T		
9	T	F	T		
10	F	F	T		
11	T	F	T		
12	T	T	T		
13	T	F	T		
14	T	T	T		
15	T	F	T		
16	T	F	T		

Empirical distribution				
H	S	E	Pr _D (.)	
T	T	T	2/16	}
T	T	F	0/16	
T	F	T	9/16	
T	F	F	1/16	
F	T	T	0/16	}
F	T	F	1/16	
F	F	T	2/16	
F	F	F	1/16	

Learning Parameters

Case 2: using **incomplete** data

- Apply **Expectation Maximization (EM)** algorithm
 - Use estimate parameters to construct CPTs, BN and compute distribution. If for initialization, randomly choose parameters to get $CPT_1 \rightarrow BN_1 \rightarrow Pr_1(\cdot)$.
 - Complete incomplete data cases/rows with all possible instances, and compute corresponding prob. using the current distribution
 - E.g. $Pr(E = T | H = T, S = F)$
 - Build empirical distribution based on prob. in step 2
 - Estimate parameters using empirical distribution to get updated CPTs, and **iteratively** perform steps 1-4 until the likelihood/score of parameters converge



Case	H	S	E
1	T	F	T
2	T	F	T
3	F	T	F
4	F	F	T
5	T	F	F
6	T	F	T
7	F	F	F
8	T	F	T
9	T	F	T
10	F	F	T
11	T	F	T
12	T	T	T
13	T	F	T
14	T	T	T
15	T	F	T
16	T	F	T

Empirical distribution

↓

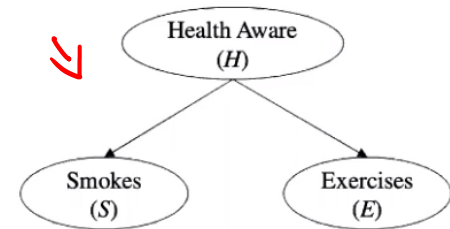
H	S	E	$Pr_{\mathcal{D}}(\cdot)$
T	T	T	2/16
T	T	F	0/16
T	F	T	9/16
T	F	F	1/16
F	T	T	0/16
F	T	F	1/16
F	F	T	2/16
F	F	F	1/16

Handwritten: 9.6/16, 1.4/16

Learning Parameters

Case 2: using **incomplete** data

- Apply **Expectation Maximization (EM)** algorithm
 - It's guaranteed that the likelihood/score of parameters won't get worse through iterations
 - Likelihood is computed by multiplying probabilities of getting all cases
 - Likelihood will either be improved by iterations or converges to a number
 - EM method is similar to local search algorithms in terms of initializing at random point and optimize through iterations
 - Where EM converges to depends on where it started
- may have multiple maximum likelihood parameters if data is incomplete



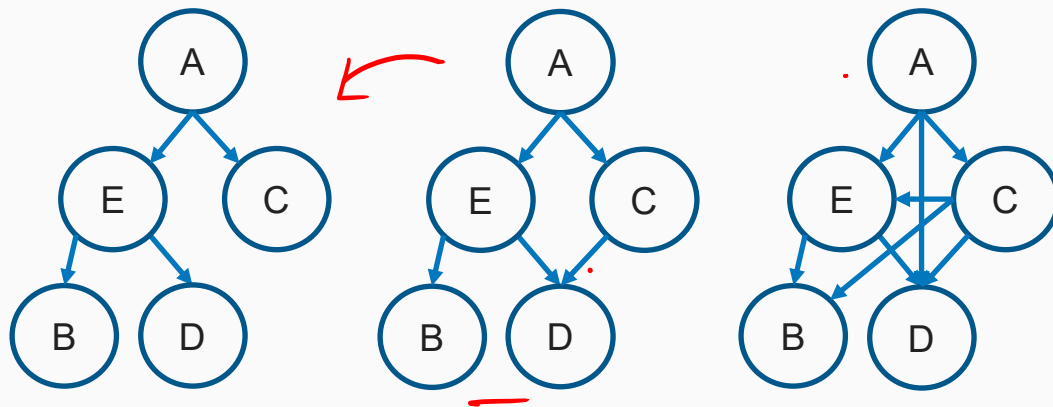
Case	H	S	E
1	T	F	T
2	T	F	T
3	F	T	F
4	F	F	T
5	T	F	F
6	T	F	T
7	F	F	F
8	T	F	T
9	T	F	T
10	F	F	T
11	T	F	T
12	T	T	T
13	T	F	T
14	T	T	T
15	T	F	T
16	T	F	T

Empirical distribution

↓

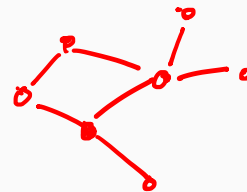
H	S	E	$\Pr_{\mathcal{D}}(.)$
T	T	T	2/16
T	T	F	0/16
T	F	T	9/16
T	F	F	1/16
F	T	T	0/16
F	T	F	1/16
F	F	T	2/16
F	F	F	1/16

Learning BN Structure

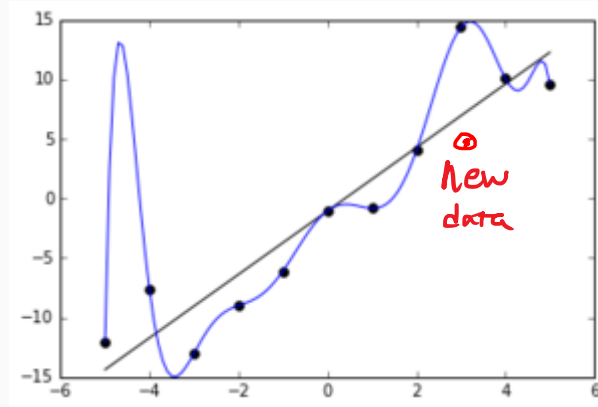
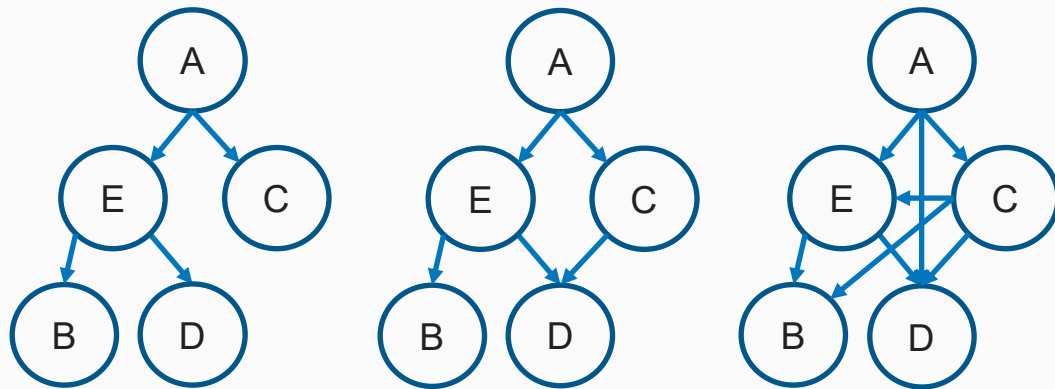


Goal: learn the BN structure given some scoring function

- Methods to find the structure with max. score:
 - Local search methods: add, remove, reverse an edge
 - Approximate method with no guarantee to find optimal structure, but efficient
 - Systematic search methods: A* search
 - Exact method to find optimal structure, but expensive



Learning BN Structure



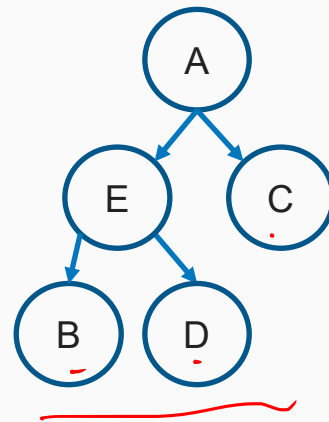
- Maximum likelihood is not enough due to over-fitting problem
 - Over-fitting: a statistical model is overfitting if it contains more parameters than can be justified by the data → very likely fail to fit additional data
 - A penalty term that penalizing the # of parameters in BN is needed to calculate scores
 - A good way is using MDL score, which uses both likelihood and penalty terms

Model-oriented vs. Query-oriented Learning

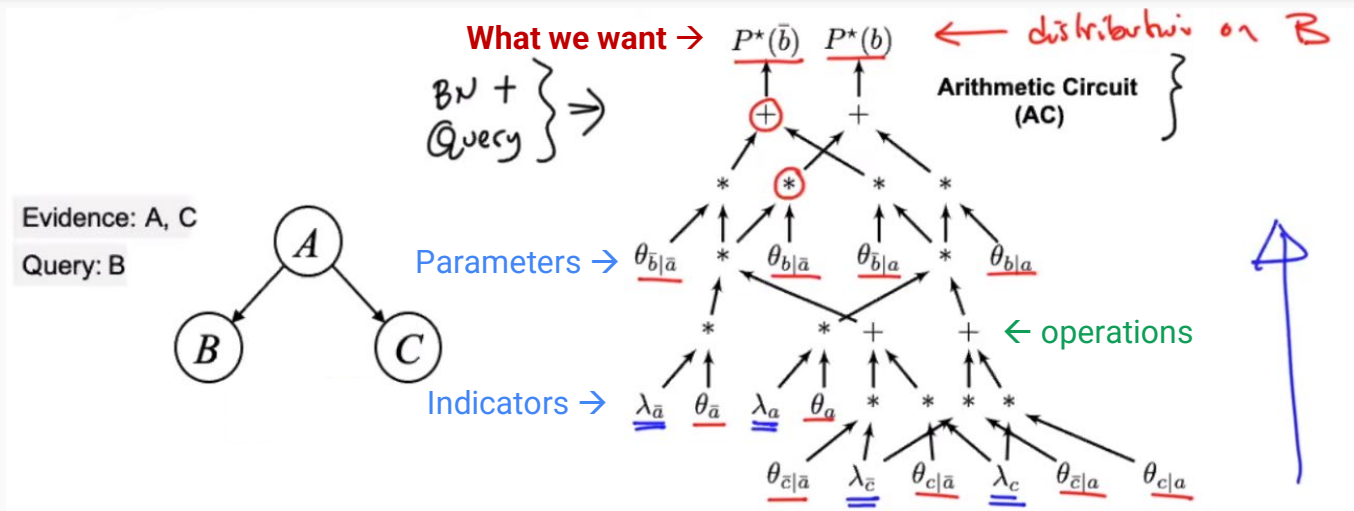
A.k.a. unsupervised vs. supervised learning, or un-labeled vs. labeled learning

Difference:

- Sometimes we want to learn the model to answer all kinds of queries
- Sometimes we only interest in one query, e.g. $\Pr(A \mid B, D, C)$, we need more focused learning \rightarrow optimize performance while answering the query
 - We need different optimization score .



Arithmetic Circuit



A	λ_a	$\lambda_{\bar{a}}$
T	1	0
F	0	1
?	1	1

\rightarrow evidence (input)

$A = T, C = F$

$\lambda_a = 1$
 $\lambda_{\bar{a}} = 0$
 $\lambda_c = 0$
 $\lambda_{\bar{c}} = 1$

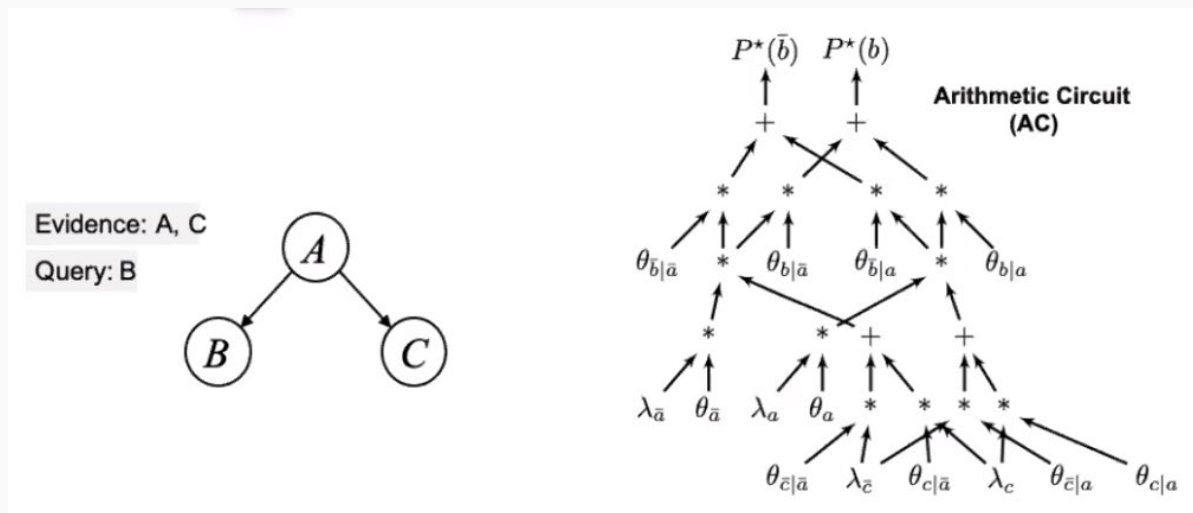
$A = F$

$\lambda_a = 0$
 $\lambda_{\bar{a}} = 1$
 $\lambda_c = 1$
 $\lambda_{\bar{c}} = 1$

Supervised Learning in BN

- Arithmetic circuit (AC) can be compiled from a BN in $\mathcal{O}(n \cdot d^w)$ time
 - N: # of variables, d: max # of values of each variable, w: treewidth

Q: If we don't know the parameters but has labeled data, what should we do?



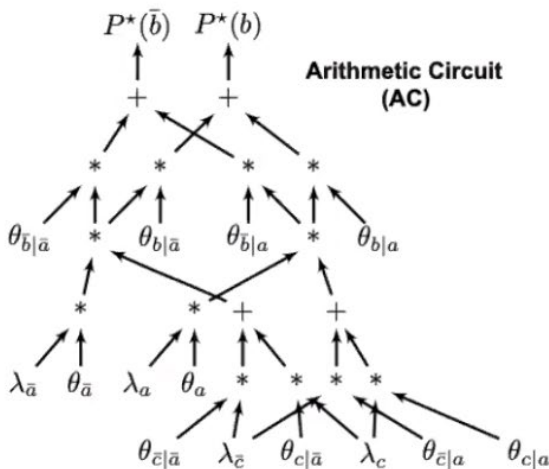
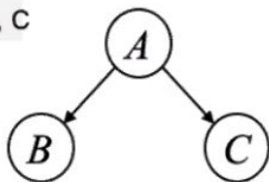
Labeled Data

input		output
A	C	B
T	T	F
T	F	F
⋮	⋮	⋮
T	T	F

Supervised Learning in BN

Evidence: A, C

Query: B



Labeled Data

	input		output
	A	C	B
case 1	T	T	F
case 2	T	F	F
⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮
	T	T	F

Use **loss function** to optimize parameters:

- Minimize the mean cross entropy between prediction in AC and one-hot distribution obtained from labeled data
 - Cross entropy**: measures how close two distributions are, it takes $P(x)$ and $Q(x)$ and returns a number

$\Rightarrow A=T, C=T \Rightarrow \underline{P_c(B)}$ (Prediction)
 $A=T, C=T \Rightarrow \begin{array}{c|c} B & \\ \hline T & 0 \\ F & 1 \end{array}$ (label)
 } one-hot distribution
 for each case

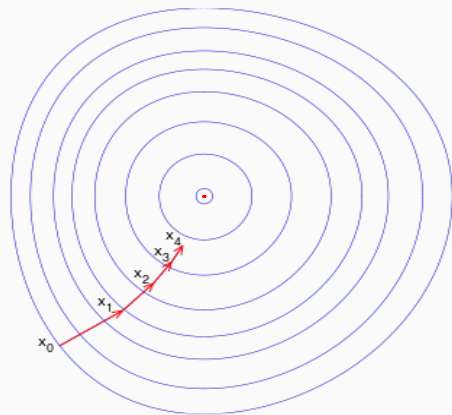
Supervised Learning in BN

Cross entropy (CE) calculation:

- Given two distributions: $P(X)$ – prediction, $Q(X)$ – label
- $CE = \sum_x Q(x) \log_2(P(x))$ → We try to minimize CE as a loss function using gradient descent

Gradient descent is often used to optimize loss function and find parameters together with neural network.

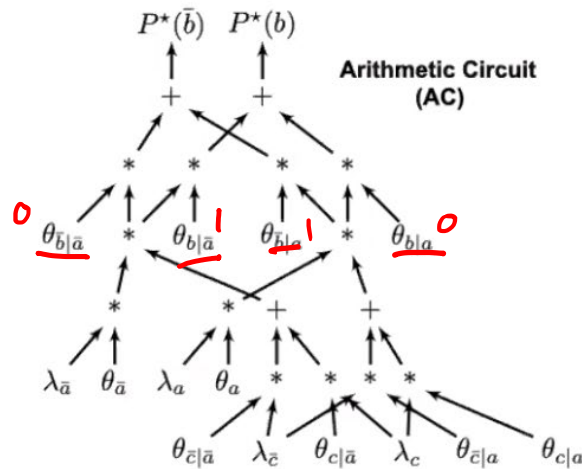
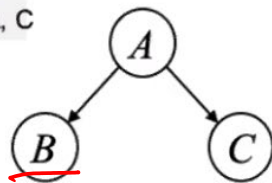
- It improves parameters by computing derivatives of loss function w.r.t those parameters, and move in the direction of the derivatives



Background Knowledge in BN

Evidence: A, C

Query: B



Assume we know “B = T iff A = F”, we then have background knowledge:

A	B	
T	T	0
T	F	1
F	T	1
F	F	0

background knowledge

→ we can replace parameters in AC with numbers (no need to learn!)

- Background knowledge allows the model to achieve certain accuracy for prediction with less data
- Background knowledge makes the model more robust

Today's Topics

- BN Learning and Inference
 - Learning parameters
 - Learning BN structure
 - Model-oriented vs. Query-oriented learning
- **Decision Tree and Random Forest**
 - Entropy and Conditional Entropy
 - Classifier
 - BN Classifier
- Hint of HW9

Entropy

Entropy can quantify uncertainty:

- Entropy is non-negative, higher entropy means more uncertainty
- $\text{ENT}(\mathbf{X}) = -\sum_x \text{Pr}(x) \log_2(\text{Pr}(x))$

E.g.

	Earthquake	Burglary	<u>Alarm</u>
True	$\text{Pr}(E=T) = 0.1$	0.2	0.2442
False	$\text{Pr}(E=F) = 0.9$	0.8	0.7558
E	<u>0.469</u>	0.722	<u>0.802</u>

$$\begin{aligned}\text{ENT}(E) &= -0.1 \times \log_2(0.1) - 0.9 \times \log_2(0.9) \\ &= 0.469\end{aligned}$$

Conditional Entropy

We are interested in $ENT(X)$ and we observe variable $Y = y$:

- $ENT(X | y) = -\sum_x Pr(x|y) \log_2(Pr(x|y))$

If we want to know the entropy of X when we observe variable Y , but don't know the value of Y yet:

- $ENT(X | Y) = \sum_y Pr(y) ENT(x|y)$
- $ENT(X | Y) \leq ENT(X) \rightarrow$ indicates that we never lose by observing a variable
- $ENT(X | Y) = ENT(X)$ if X and Y are independent

Conditional Entropy - Example

	Burglary	Burglary A=true	Burglary A=false
True	0.2	0.741	0.025
False	0.8	0.259	0.975
E	<u>0.722</u>	<u>0.825</u>	<u>0.169</u>

- Unlike $ENT(X | Y)$, $ENT(X | y)$ can be higher or lower than $ENT(X)$.

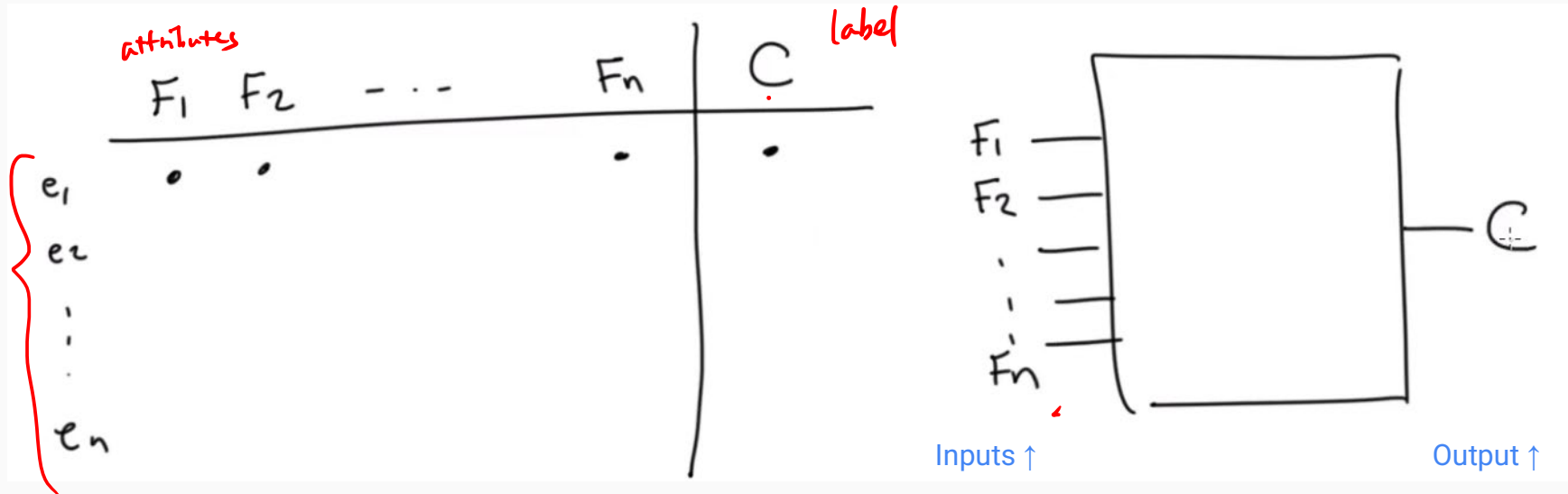
- $ENT(\text{Burglary} | \text{Alarm}) = ENT(B | a) Pr(a) + ENT(B | \neg a) Pr(\neg a)$
 $= 0.329 \leq 0.722$

$ENT(\text{Burglary})$

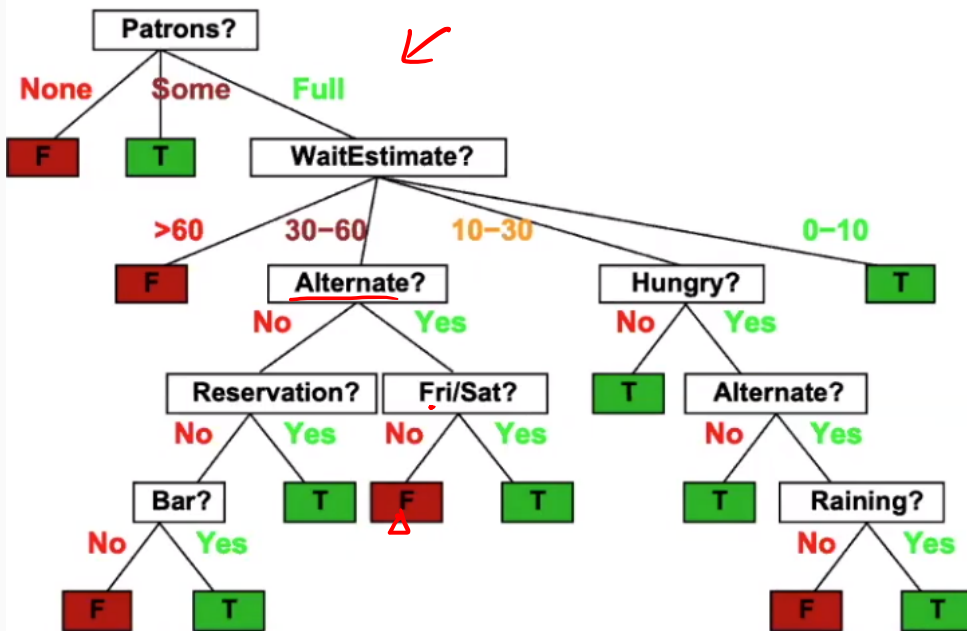
Classifiers

Supervised learning with labeled data

- Take features/attributes as input and output class labels
- A binary classifier has only two types of output



Decision Tree - Example



labeled data

Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait
X ₁	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X ₂	T	F	F	T	<u>Full</u>	\$	F	F	Thai	30-60	<u>F</u>
X ₃	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X ₄	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X ₅	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X ₆	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X ₇	F	T	F	F	None	\$	T	F	Burger	0-10	F
X ₈	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X ₉	F	T	T	F	Full	\$	T	F	Burger	>60	F
X ₁₀	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X ₁₁	F	F	F	F	None	\$	F	F	Thai	0-10	F
X ₁₂	T	T	T	T	Full	\$	F	F	Burger	30-60	T

label

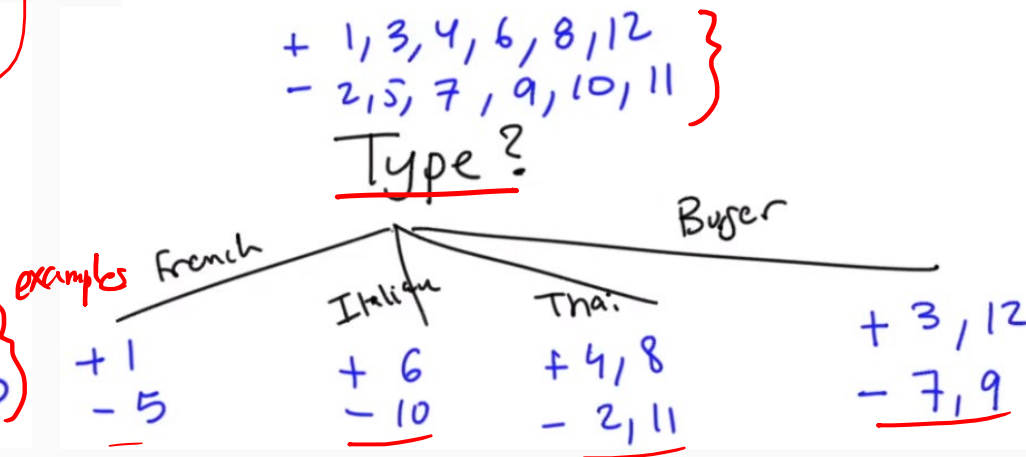
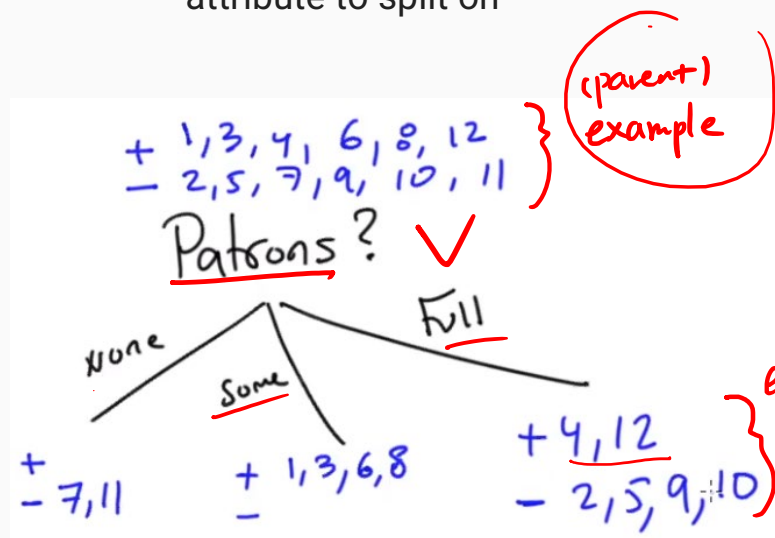
- This decision tree is interpretable
- Classification using decision tree is very simple (e.g. follow attributes to assign a missing label)

Build a Decision Tree

Key: decide which attribute to split on

- We want to split on the attribute with higher discriminating power
 - E.g. Patrons is better than Type as an attribute to split on

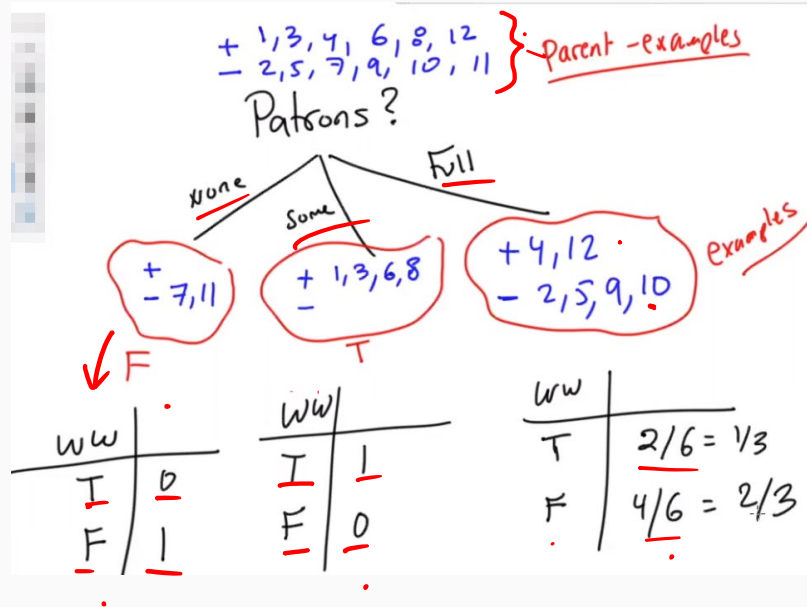
Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait
X ₁	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X ₂	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X ₃	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X ₄	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X ₅	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X ₆	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X ₇	F	T	F	F	None	\$	T	F	Burger	0-10	F
X ₈	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X ₉	F	T	T	F	Full	\$	T	F	Burger	>60	F
X ₁₀	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X ₁₁	F	F	F	F	None	\$	F	F	Thai	0-10	F
X ₁₂	T	T	T	T	Full	\$	F	F	Burger	30-60	T



Build a Decision Tree

Steps:

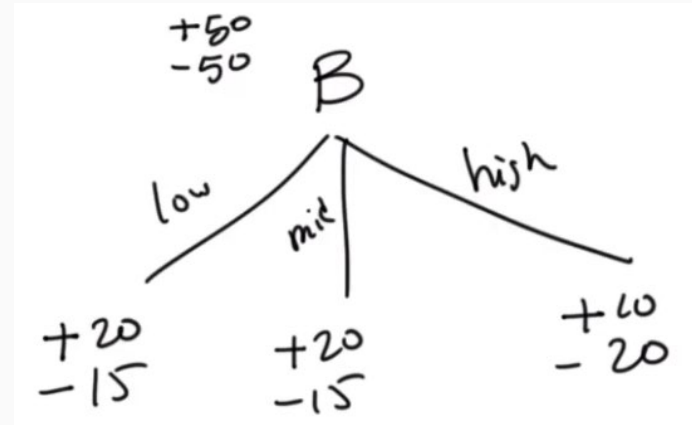
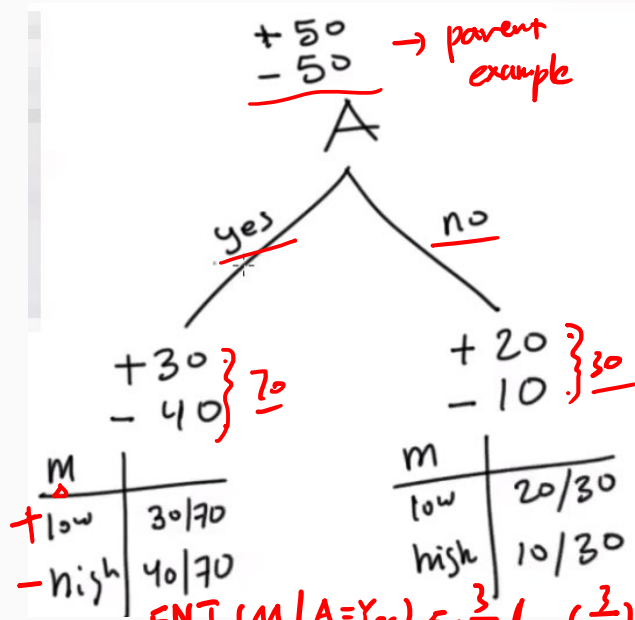
- Consider each of split case as a distribution
- Use **entropy** to quantify the distributions (how certain we are about the class label) after the split
- Compute **conditional entropy** of the output label given current attribute \rightarrow score of attribute
 - Recall $ENT(X | Y) = \sum_y Pr(y) ENT(x|y)$
 - E.g. $ENT(\text{WillWait} | \text{Patrons})$



Build a Decision Tree – Example

Let's say we want to output labels of $M = \{low, high\}$ *100 data examples*

We have two attributes to choose – A and B



Score of A:

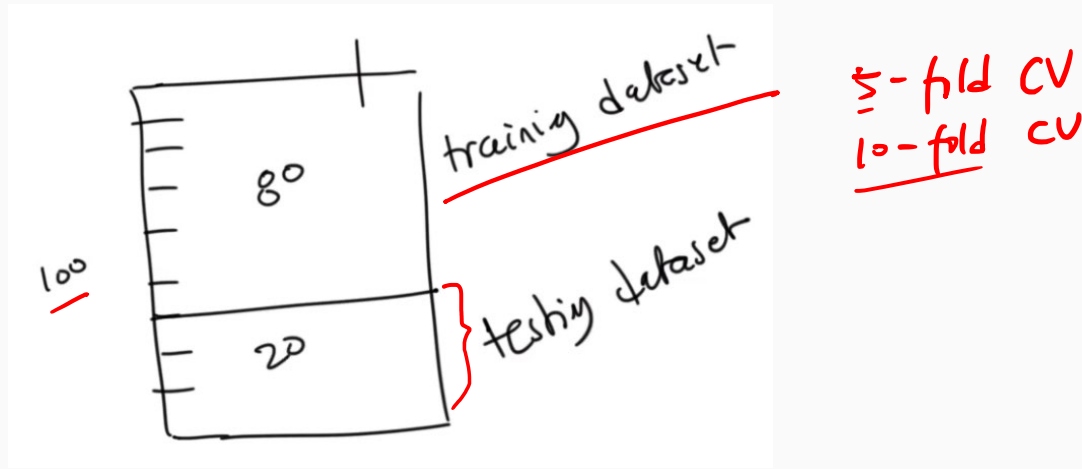
$$= ENT(M/A) = 0.7 \times 0.985 + 0.3 \times 0.918 \approx \underline{0.965}$$

$\Rightarrow ENT(M/A=Yes) = -\frac{3}{7} \log_2(\frac{3}{7}) - \frac{4}{7} \log_2(\frac{4}{7}) \approx 0.985$
 $\Rightarrow ENT(M/A=No) = -\frac{2}{3} \log_2(\frac{2}{3}) - \frac{1}{3} \log_2(\frac{1}{3}) \approx 0.918$

Cross Validation

Idea: split dataset to training and testing portions, build decision trees multiple times using different combinations of training and testing portions, take the average score

- Very useful to deal with over-fitting problem



Random Forests

Random Forests (RF) is an ensemble learning method that combine the results of multiple decision trees

- It does **majority voting** to do classification

Techniques of construct individual decision trees in RF:

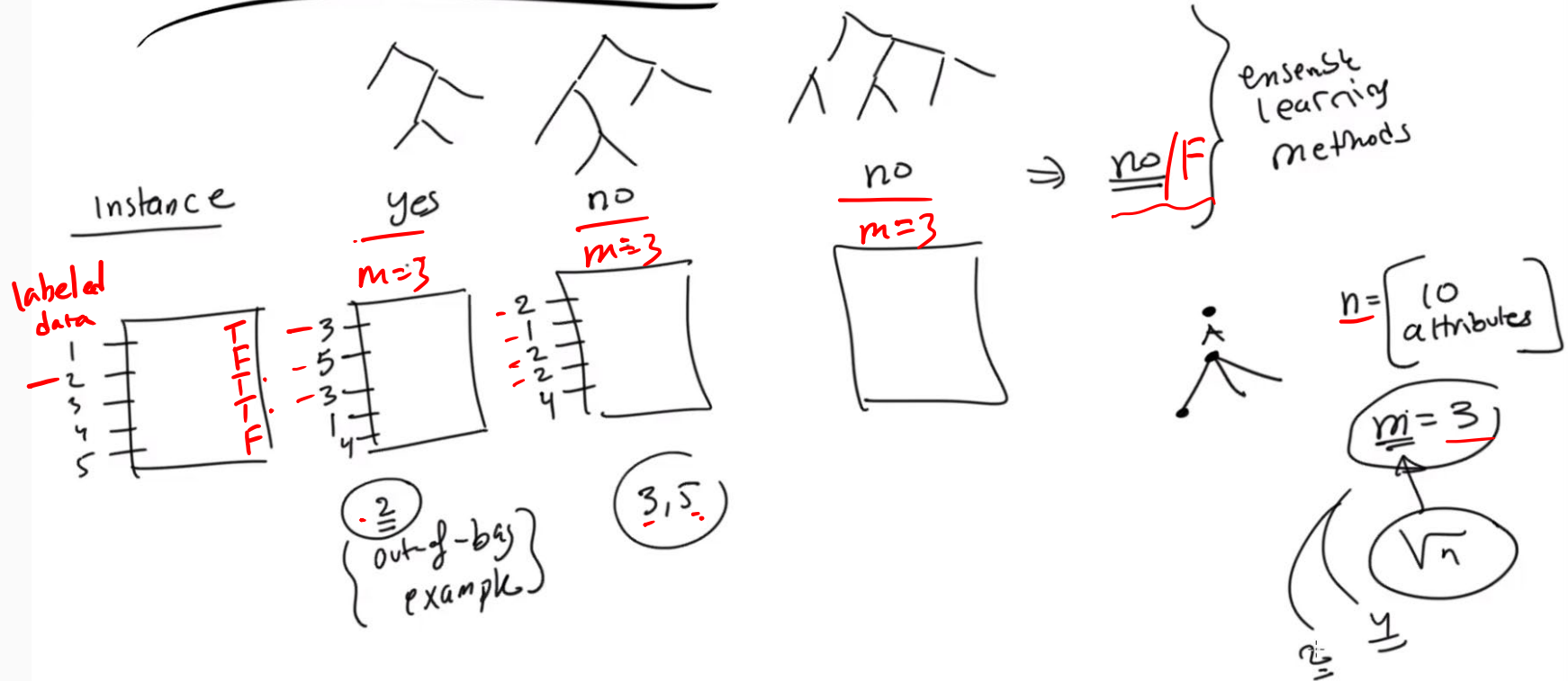
- **Bootstrapping dataset:** randomly sample (with replacement) from the dataset to construct each tree
 - Unselected examples are called out-of-bag examples
- Randomly choose a subset of attributes, then scoring and choose among them

Evaluate the RF:

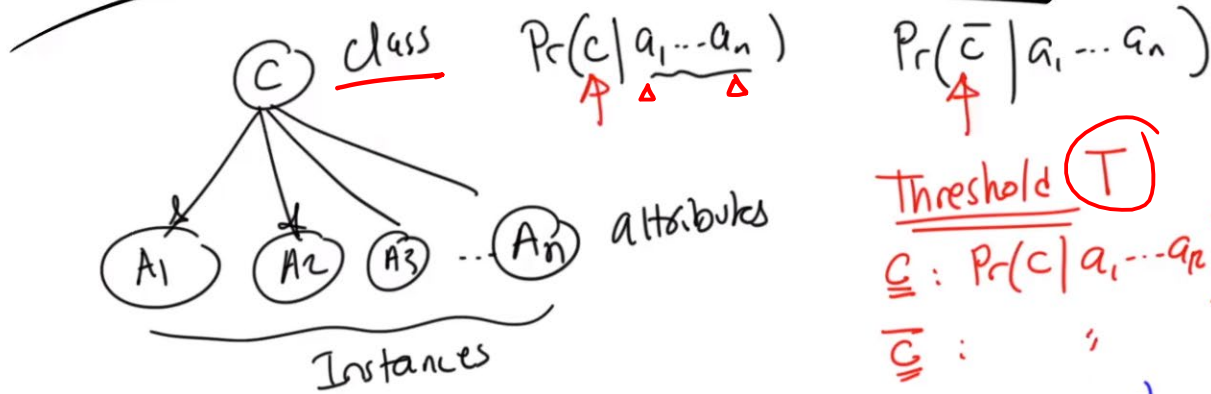
- Run the constructed RF on out-of-bag examples check how well it performs
- Change parameters (e.g. # of picked attributes) to improve performance

Random Forests – Example

Random Forests



Bayesian Network Classifier



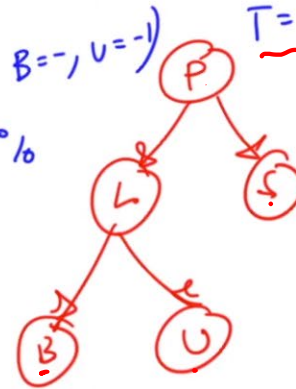
Threshold T

$$\underline{c} : \Pr(c | a_1, \dots, a_n) \geq T$$

$$\bar{c} : \Pr(\bar{c} | a_1, \dots, a_n) < T$$

Example: $\Pr(p | s=+, B=-, U=-) \geq 90\%$ $T = 90\%$

not pregnant + pregnant



Instance

S	B	U
+	-	-

Questions?

My slides take the following materials as references:

- Prof. Darwiche's lecture video

Thank you!