

Connect

beyond learning – reasoning; why

..... logic

..... and its limits

..... fundamental, uncertainty

..... reasoning under uncertainty

..... back to learning - from text

connecting the dots: motivation

“who is the leader of USA?”

facts ... [X is prime-minister of C] ... [X is president of C]

no such fact [X is leader of USA] ... now what?

X is president of C \Rightarrow X is leader of C – *rules (knowledge)*

✓ Obama is president of USA \Rightarrow Obama is leader of USA

example of *reasoning* ..

reasoning can be tricky:

Manmohan Singh is prime-minister of India

Pranab Mukherjee is president of India

“who is the leader of India”

... *much* more *knowledge* is needed

reasoning and web-intelligence

“book me an American flight to NY ASAP”

“this New Yorker who fought at the battle of Gettysburg was once considered the inventor of baseball”

Alexander Cartwright or Abner Doubleday – *Watson got it right*

“who is the Dhoni of USA?”

– *analogical reasoning* - X is to USA what Cricket is to India (?)

+ *abductive reasoning* – there *is no* US baseball team ... so ?

find *best possible answer*^

+ *reasoning under uncertainty* ... who is the “most” popular ?

Semantic Web:

- web of linked *data, inference* rules and engines, query
 - pre-requisite: extracting *facts* from text, as well as *rules*

logic: propositions

A, B – ‘propositions’ (either True or False)

A and B is True: $A=\text{True}$ and $B=\text{True}$ ($A \wedge B$)

A or B is True: either $A=\text{True}$ or $B=\text{True}$ ($A \vee B$)

if A then B (same as if $A=\text{True}$ then $B=\text{True}$)

is the same as saying $A=\text{False}$ or $B=\text{True}$

also written as:

$A \Rightarrow B$ is equivalent to **$\sim A \vee B$**

check: $A=T$, $\sim A=F$, so $(\sim A \vee B) = T$ only when $B=T$

Important:

if $A=F$, $\sim A=T$, so $(\sim A \vee B)$ is true regardless of B being T or F

logic: predicates

Obama is president of USA:

isPresidentOf (Obama, USA) - *predicates, variables*

X is president of C => X is leader of C

R: isPresidentOf (X, C) => isLeaderOf (X, C)

plus – the above is stating a rule for *all* X,C - *quantification*

“Obama is president of USA”: *fact*

F: isPresidentOf (Obama, USA)

using rule R and fact F,

isLeaderOf (Obama, USA) is *entailed*

(*unification: X bound to Obama; C bound to USA*)

Q: isLeaderOf (X, USA) – *query*

reasoning = answering queries or deriving new facts

using *unification + inference = resolution*

semantic web vision

facts and rules in RDF-S & OWL-..

web of *data* and *semantics*

web-scale inference

Google²; Wolfram-Alpha; Watson^{*}

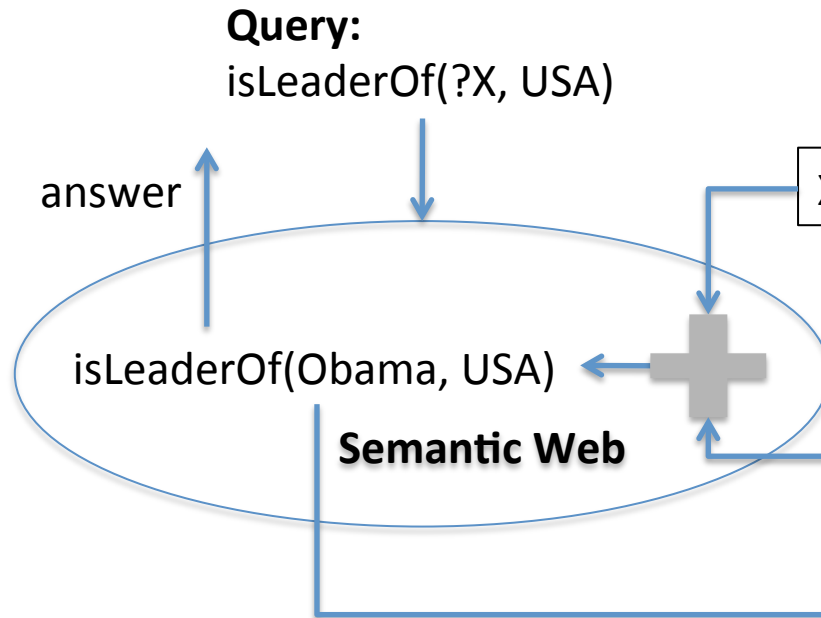
Manmohan Singh is prime-minister of India
Pranab Mukherjee is president of India
Vladimir Putin is president of Russia
Obama is president of USA
... is president of ...
... is premier of ...
a.com

inductive reasoning (rule learning)

X is president of C => X is leader of C c.com

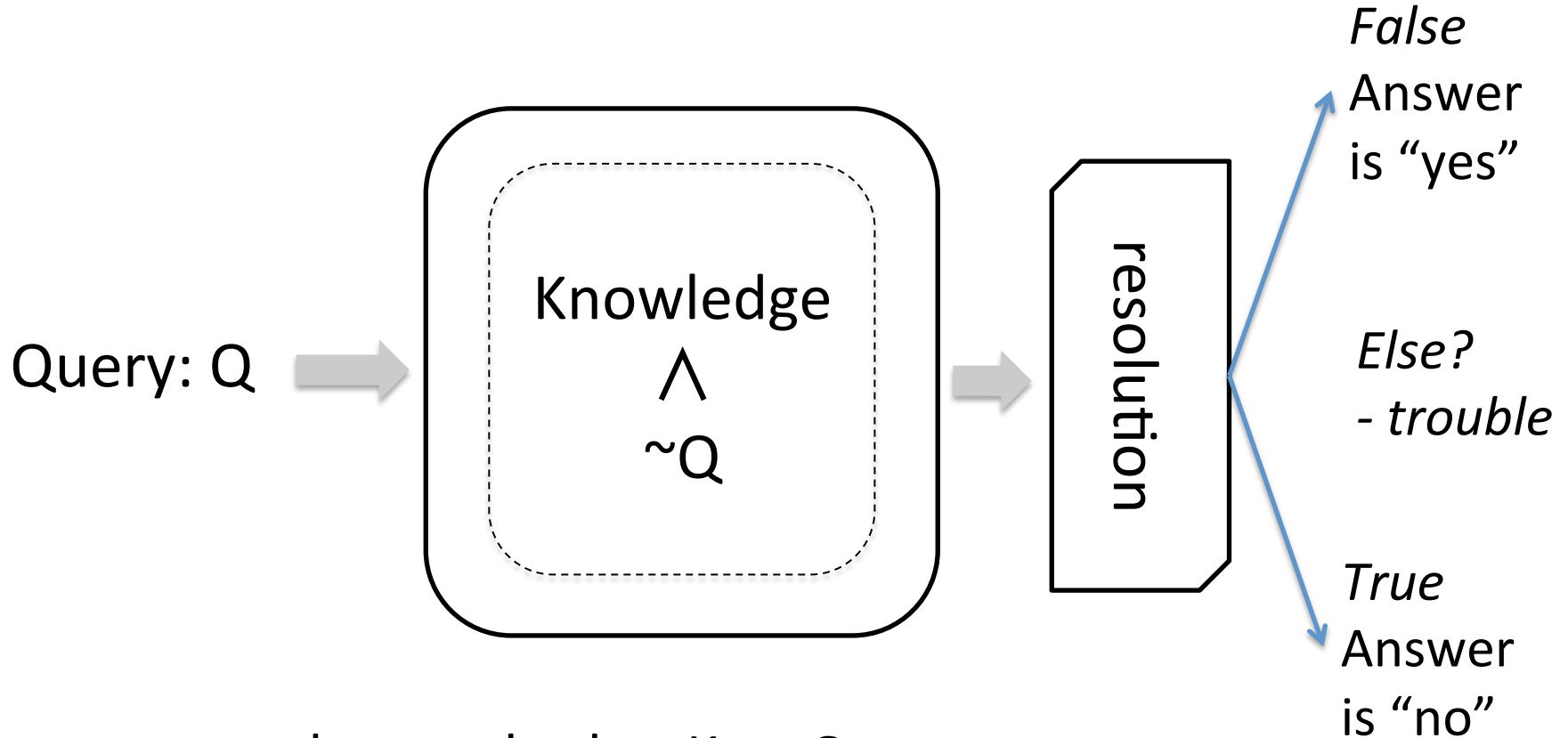
deductive reasoning
(logical inference)

isLeaderOf(Manmohan Singh, India)
isLeaderOf(Zuma, South Africa)
isLeaderOf(Putin, Russia)
.....
b.com



*don't use RDF, OWL or semantic-web
technology though they have similar intent, spirit ...

logical inference: resolution



we want to know whether $K \Rightarrow Q$

i.e. $\sim K \vee Q$ is True

i.e. $K \wedge \sim Q$ is False !

in other words K augmented with $\sim Q$ *entails* falsehood, for sure

logic: fundamental limits

resolution may never end; *never* (whatever algorithm!)

➤ undecidability

predicate logic undecidable (Godel, Turing, Church ...)

➤ intractability

propositional logic is decidable, but intractable (SAT and NP ..)

? whither automated reasoning, semantic-web..?

fortunately:

OWL-DL, OWL-lite (description logic: $\text{leader} \sqsubseteq \text{person} \dots$)

decidable; still intractable in worst case

Horn logic (rules, i.e., $\text{person} \wedge \text{bornIn}(C) \Rightarrow \text{citizen}(C) \dots$)

undecidable (except with caveats); but tractable

logic and uncertainty

predicates A , B , C

1. For all x , $A(x) \Rightarrow B(x)$.

2. For all x , $B(x) \Rightarrow C(x)$

1 and 2 *entail* For all x , $A(x) \Rightarrow C(x)$ fundamental

however, consider the *uncertain* statements:

1': For **most** x , $A(x) \Rightarrow B(x)$. “*most firemen are men*”

2'. For **most** x , $B(x) \Rightarrow C(x)$. “*most men have safe jobs*”

it *does not* follow that “For **most** x , $A(x) \Rightarrow C(x)$ ” !



logic and causality

- if the sprinkler was on then the grass is wet

$$S \Rightarrow W$$

- if the grass is wet then it had rained

$$W \Rightarrow R$$

therefore it follows, i.e. $S \Rightarrow R$ is *entailed*

which states “the sprinkler is on, so it had rained”

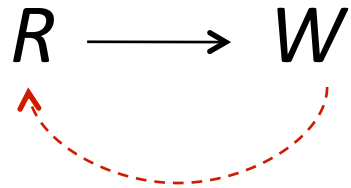
- *problem is that causality was treated differently in each statement \Rightarrow absurdity*

causality and classification

if S then W (W is an observable *feature* of S)

$$S \longrightarrow W$$

if R then W (W is an observable *feature* of R)



if W is observed then R happened abduction

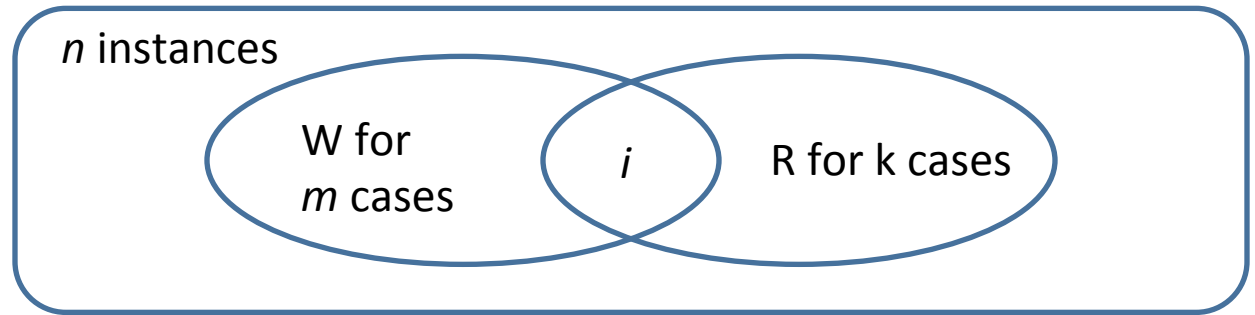
concluding which *class of event* observed S or R

abductive reasoning

= from effects to likely causes

probability tables and 'marginalization'

#	W	R
1	y	n
2	y	y
3	n	n



consider $p(R, W)$

to get $p(R)$ we can 'sum out' W: $p(R) = \sum_W p(R, W)$

this is called *marginalization* of W

notice that marginalization is equivalent to *aggregation* on column P:

$$\sum_W p(R, W) = {}_R G_{\text{SUM}(P)} T^{R, W}$$

or, in SQL:

```
SELECT R, SUM(P) from TR,W
GROUP BY R
```

R	P		R	W	P
y	k/n	= \sum_W	y	y	i/n
n	(n-k)/n		n	y	(m-i)/n
			y	n	(k-i)/n
			n	n	(n-m-k+i)/n

$$P(R, W) = T^{R, W}$$

probability tables and Bayes rule ...

R	W	P
y	y	i/n
n	y	$(m-i)/n$
y	n	$(k-i)/n$
n	n	$(n-m-k+i)/n$

$p(R, W)$

=

R	W	P
y	y	i/m
n	y	$(m-i)/m$
y	n	$k-i/(n-m)$
n	n	$(n-m-k+i)/(n-m)$

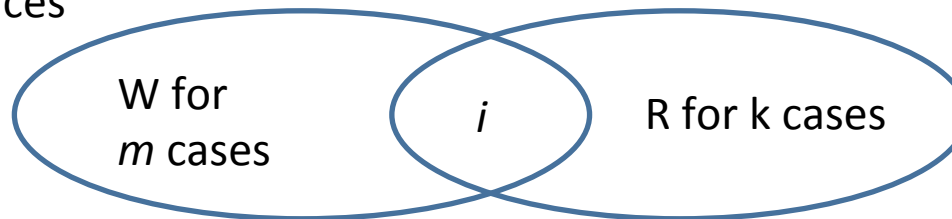
$p(R | W)$

*

W	P
y	m/n
n	$(n-m)/n$

$p(W)$

n instances



probability tables and Bayes rule ...

R	W	P
y	y	i/n
n	y	$(m-i)/n$
y	n	$(k-i)/n$
n	n	$(n-m-k+i)/n$

$p(R, W)$

$T_0^{R, W}$

=

R	W	P
y	y	i/m
n	y	$(m-i)/m$
y	n	$k-i/(n-m)$
n	n	$(n-m-k+i)/(n-m)$

$p(R | W)$

$T_1^{R, W}$

*

W	P
y	m/n
n	$(n-m)/n$

$p(W)$

T_2^W

notice that the product $p(R | W) p(W) = T_1^{R, W} \bowtie_B T_2^W$

i.e., the *join* of the two tables T_1 and T_2 on the common attribute W !
so, probability tables (also called *potentials*) can be multiplied in SQL!

SELECT R, SUM(P1*P2) from $T_1^{R, W}$, T_2^W WHERE W1=W2 GROUP BY R

probability tables and *evidence*

R	W	P
y	y	i/n
n	y	(m-i)/n
y	n	(k-i)/n
n	n	(n-m-k+i)/n

$$P(R, W)$$

$$= T^{R, W}$$

$$\mathbf{e}^{(B=y)} =$$

R	W	P
y	y	i/n
n	y	(m-i)/n

$$P(R, W) \mathbf{e}^{(W=y)}$$

SELECT R,W,P from $T^{R, W}$ WHERE $W=y$

$$=$$

R	W	P
y	y	i/m
n	y	(m-i)/m

$$P(R | W=y)$$

$$* p(W=y)$$

$$* m/n$$

if we restrict $p(R, W)$ to entries where *evidence* $W=y$ holds:

$$p(R, W) \mathbf{e}^{(W=y)} = p(R | W=y) * p(\mathbf{e}^{(W=y)})$$

applying evidence is equivalent to the *select* operator on $T^{R, W}$

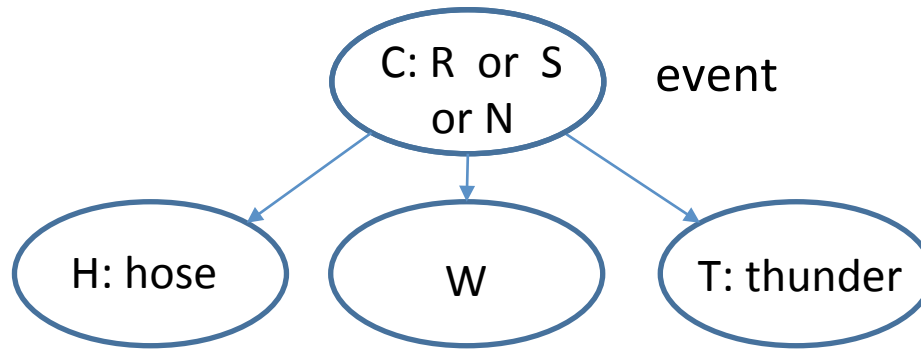
$$P(R, W) \mathbf{e}^{(W=y)} = \sigma_{W=y} T^{R, W}$$

so the *a posteriori* probability of R given evidence \mathbf{e} is just:

$$P(R | \mathbf{e}^{(W=y)}) = P(R, W) \mathbf{e}^{(W=y)} / p(\mathbf{e}^{(W=y)})$$

A	P
y	i/m
n	(m-i)/m

naïve Bayes classifier



assumption – independence of features $H, W, T \mid C \Rightarrow$

$$p(C \mid H, W, T) = \sigma \ p(H, W, T \mid C) = \sigma \ p(H \mid C) \ p(W \mid C) \ p(T \mid C)$$

and in general for n features:

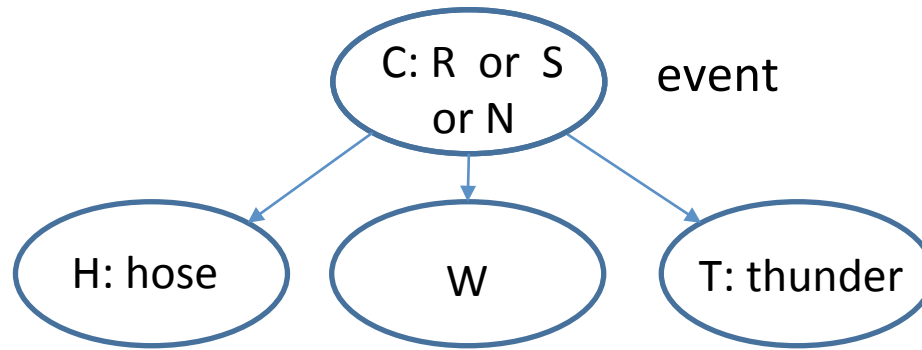
$$p(C \mid F_1 \dots F_n) = \sigma \ p(F_1 \dots F_n \mid C) = \sigma \ p(F_1 \mid C) \dots p(F_n \mid C)$$

- remember, these are tables (multiplied as before: SQL!)

now given observations $\mathbf{e}^{f_1, \dots, f_n}$ we get the likelihood rule

$$p(C \mid F_1 \dots F_n) \ \mathbf{e}^{f_1, \dots, f_n} = \sigma' \ p(f_1 \dots f_n \mid C) = \sigma' \ p(f_1 \mid C) \dots p(f_n \mid C)$$

naïve Bayes classifier and partial evidence



given observations $\mathbf{e}^{f_1, \dots, f_n}$ we get the likelihood rule

$$p(C | F_1 \dots F_n) \mathbf{e}^{f_1, \dots, f_n} = \sigma' p(f_1 \dots f_n | C) = \sigma' p(f_1 | C) \dots p(f_n | C)$$

again, ... even if some features are *not* measured, e.g. F_1 :

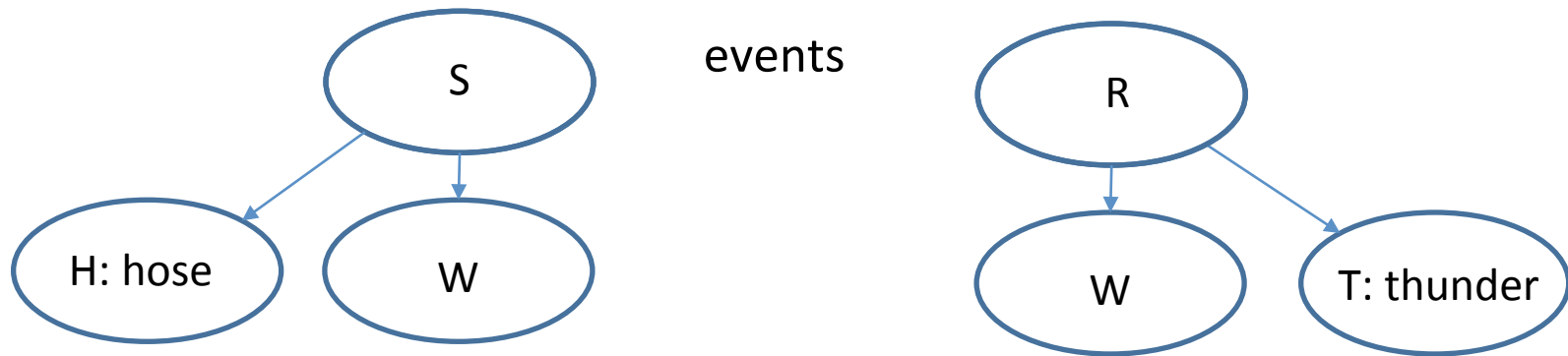
$$p(C | F_1 F_2 \dots F_n) \mathbf{e}^{f_2, \dots, f_n} = \sigma'' \sum_{F_1} p(F_1 | C) p(f_2 | C) \dots p(f_n | C)$$

in SQL:

```
SELECT C, SUM( $\prod_i P_i$ ) FROM T1..Tn WHERE F2=f2 ... Fn=fn {evidence}  
AND  
GROUP by C
```

(finally, normalize so that $\sum_C = 1$, i.e. σ'' can effectively be ignored)

multiple naïve Bayes classifiers



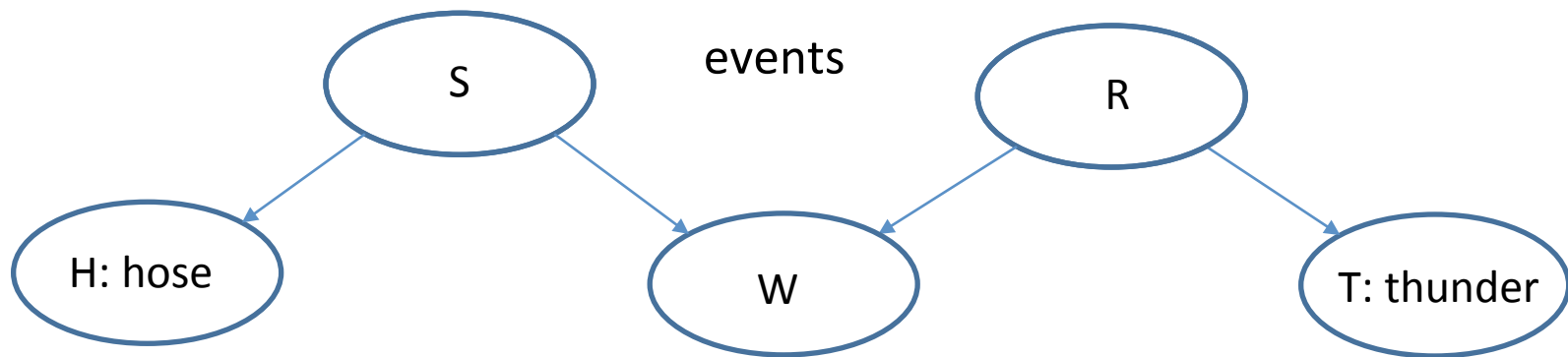
but ... R and S can happen *together*, so we need 2 classifiers

$$P(R|W,T) = \sigma_1 p(W|C) p(T|C)$$

$$P(S|H,W) = \sigma_2 p(H|C) p(W|C)$$

but ... W is the same observation ...

Bayesian network



$$P(R|H,W,T,S) = p(H,W,T,S|R) [p(R) / p(H,W,T,S)]$$

$$p(R,H,W,T,S) = p(H,W,T,S|R) p(R) = \sigma p(H,W,T,S|R)$$

assumption – independence of features $H, T, W | S, R \Rightarrow$

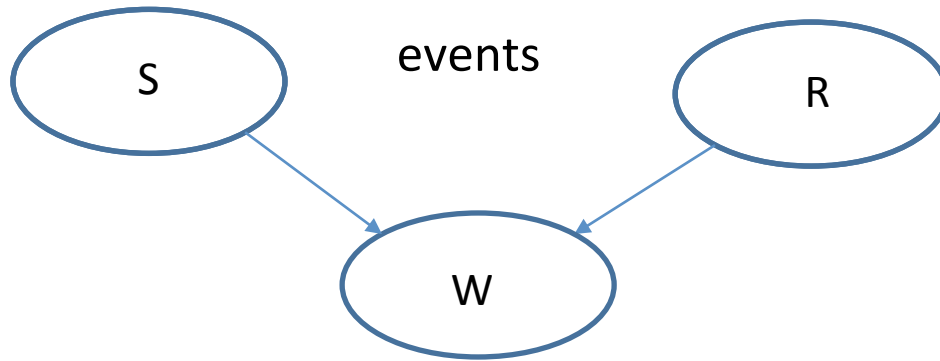
$$p(R,H,W,T,S) = \sigma p(H,W,T,S|R) = \sigma p(H|S,R) p(W|S,R) p(T|S,R)$$

But ... and this is tricky ... H, R and S, T also independent

$$p(R,H,W,T,S) = \sigma p(H|S) p(W|S,R) p(T|R) \quad \square$$

once we have the joint – “sum out everything but R” – SQL!

simple example



CPT
 $p(W|S,R)$
not joint!

W	S	R	P
y	y	y	.9
y	y	n	.7
y	n	y	.8
y	n	n	.1
n	n	n	.9
n	n	y	.2
n	y	n	.3
n	y	y	.1

$$P(W,R,S) = p(W|S,R) p(S) p(R) \quad \square$$

evidence₁: “grass is wet”, $W=y$

$$P(R|W) = \sum_S P(W,R,S) \mathbf{e}^{W=y} = \sum_S \sigma P(W|R,S) \mathbf{e}^{W=y} \quad \text{in SQL:}$$

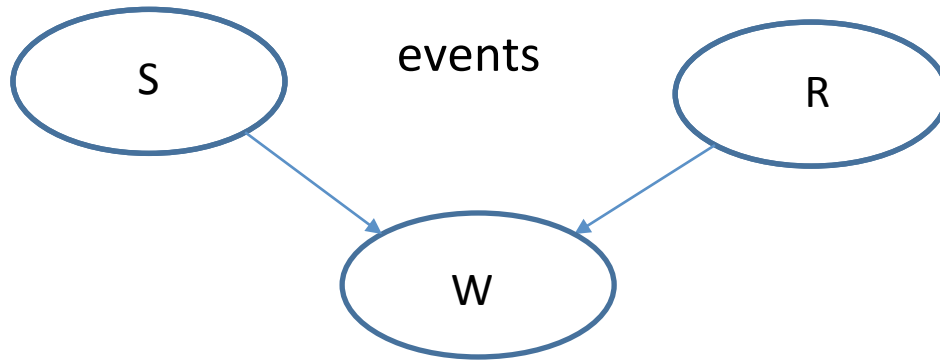
SELECT R, SUM(P) FROM T WHERE W=Y GROUP BY R

normalizing so that sum is 1:

$$p(R=y|W=y) = 1.7/(1.7+.8) = .68, \text{ i.e. } 68\%$$

W	R	P
y	y	1.7
y	n	.8

example continued: “explaining away” effect



W	S	R	P
y	y	y	.9
y	y	n	.7
y	n	y	.8
y	n	n	.1
n	n	n	.9
n	n	y	.2
n	y	n	.3
n	y	y	.1

evidence₁: “grass is wet”, $W=y$

AND evidence₂: “sprinkler on”, $S=y$

$P(R | W, S) = P(W, R, S) \mathbf{e}^{W=y, S=y} = p(R) P(W | R, S) \mathbf{e}^{W=y, S=y}$ in SQL:

SELECT R, SUM(P) FROM T WHERE W=Y, S=y GROUP BY R

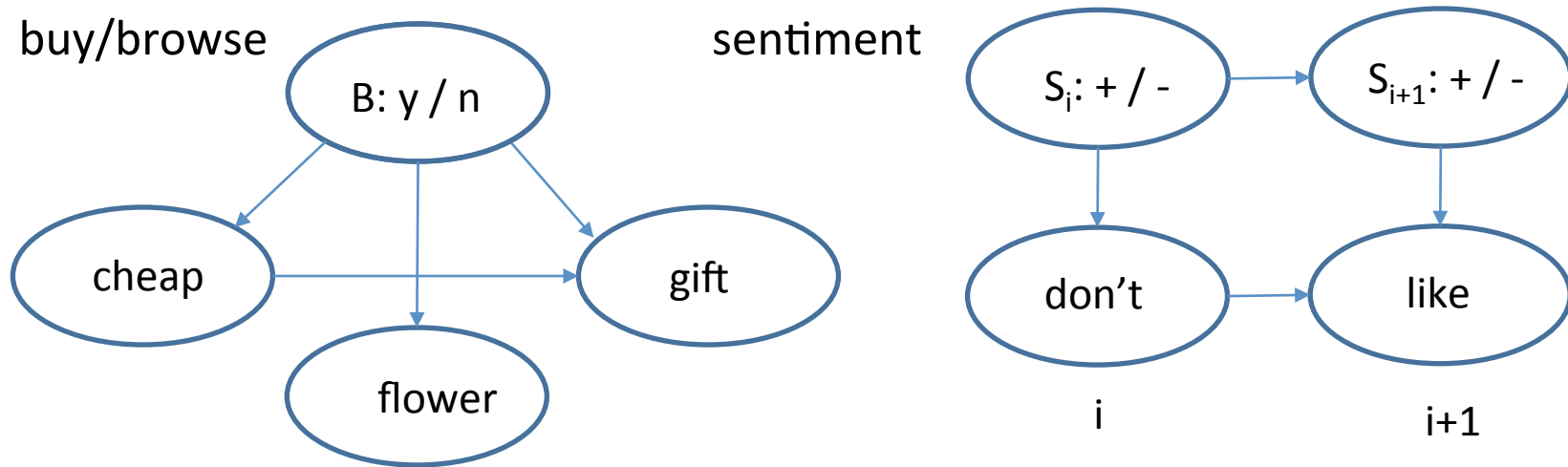
normalizing so that sum is 1:

$p(R=y | W=y, S=Y) = .9/1.6 = .56$, i.e. 56%

less than the earlier 68% - belief propagation

W	R	P
y	y	.9
y	n	.7

Bayes nets: beyond independent features



if 'cheap' and 'gift' are *not* independent, $P(G|C,B) \neq P(G|B)$

(or use $P(C|G,B)$, depending on the order in which we *expand* $P(G,C,B)$)

“I don't like the course” and “I like the course; don't complain!”

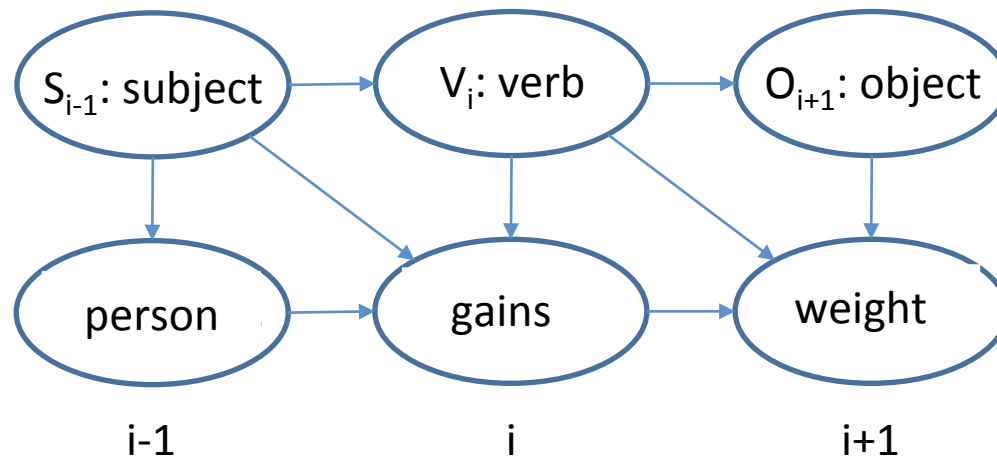
first, we might include “don't” in our list of features (also “not” ...)

still – might not be able to disambiguate: need *positional order*

$P(x_{i+1} | x_i, S)$ for each position i : hidden markov model (HMM)

we may also need to accommodate 'holes', e.g. $P(x_{i+k} | x_i, S)$

where do facts come from? learning from text



suppose we want to learn *facts* of the form <subject, verb, object> from text
single class variable is not enough; (i.e. we have many y_j in data $[Y, X]$)

further, positional order is important, so we can use a (different) HMM ..

e.g. we need to know $P(x_i | x_{i-1}, S_{i-1}, V_i)$

whether 'kill' following 'antibiotics' is a verb will depend on whether 'antibiotics' is a subject
more apparent for the case <person, gains, weight>, since 'gains' can be a verb or a noun

problem reduces to estimating *all* the a-posterior probabilities $P(S_{i-1}, V_i, O_{i+1})$
for every i , and also allowing 'holes' (i.e., $P(S_{i-k}, V_i, O_{i+p})$) and find the *best* facts
from a collection of text? many solutions; apart from HMMs - CRFs
after finding all facts from lots of text, we cull using support, confidence, etc.

open information extraction

Cyc (older, semi-automated): 2 billion facts

Yago – largest to date: 6 billion facts, linked i.e., a graph

e.g. <Albert Einstein, wasBornIn, Ulm>

Watson – uses facts culled from the web internally

REVERB – recent, lightweight: 15 million S,V,O triples

e.g. <potatoes, are also rich in, vitamin C>

1. part-of-speech tagging using NLP classifiers (trained on labeled corpora)
2. focus on verb-phrases; identify nearby noun-phrases
3. prefer proper nouns, especially if they occur often in other facts
4. extract more than one fact if possible:

“Mozart was born in Salzburg, but moved to Vienna in 1781” yields

<Mozart, moved to, Vienna>, in addition to <Mozart, was born in, Salzburg>

belief networks: learning, logic, big-data & AI

- network *structure* can be learned from data
- applications in [genomic] medicine
 - medical diagnosis
 - gene-expression networks
 - how do phenotype traits arise from genes
- logic and uncertainty
 - belief networks bridging the gap:
 - (Pearl Turing award; Markov logic n/w ...)
- big-data
 - inference can be done using SQL – map-reduce works!
- hidden-agenda:
 - deep belief networks
 - linked to connectionist models of brain

recap and preview

search is not enough for Q&A: reasoning

logic and semantic web

... but there are limits, fundamental + practical

reasoning under uncertainty Bayesian inference using SQL

... Bayesian networks and PGMs in general

Next few weeks:

next week (7) – 1 programming assignment

lecture videos only to explain but start preparing

week 8 (final lecture week) – “predict”

putting everything together! 4th prog assgn.

week 9

complete all assignments + final exam