Course plan



Machine learning



Variable-based models

Special cases:

- Constraint satisfaction problems
- Markov networks
- Bayesian networks

```
Key idea: variables-
```

- Solutions to problems \Rightarrow assignments to variables (modeling).
- Decisions about variable ordering, etc. chosen by inference.

Higher-level modeling language than state-based models

Factor graph







Factors

Definition: scope and arity-

Scope of a factor f_j : set of variables it depends on. **Arity** of f_j is the number of variables in the scope. **Unary** factors (arity 1); **Binary** factors (arity 2). **Constraints** are factors that return 0 or 1.





Example: map coloring-

Scope of $f_1(X) = [WA \neq NT]$ is $\{WA, NT\}$ f_1 is a binary constraint

Assignment weights



An assignment is **consistent** if Weight(x) > 0.

Objective: find the maximum weight assignment

 $\arg\max_x \mathsf{Weight}(x)$

A CSP is satisfiable if $\max_x \text{Weight}(x) > 0$.







Variables, factors: specify locally

Weight
$$({X_1 : B, X_2 : B, X_3 : R}) = 1 \cdot 1 \cdot 2 \cdot 2 = 4$$

Assignments, weights: optimize globally

Example: object tracking





(O) Noisy sensors report positions: 0, 2, 2.(T) Objects can't teleport.

What trajectory did the object take?



Example: object tracking CSP

Factor graph:



[demo]

- Variables $X_i \in \{0, 1, 2\}$: position of object at time i
- Observation factors $o_i(x_i)$: noisy information compatible with position
- Transition factors $t_i(x_i, x_{i+1})$: object positions can't change too much





• Decide on variables and domains

• Translate each desideratum into a set of factors

- Try to keep CSP small (variables, factors, domains, arities)
- When implementing each factor, think in terms of checking a solution rather than computing the solution

Backtracking search

Algorithm: backtracking search-

$\mathsf{Backtrack}(x, w, \mathsf{Domains})$:

- If x is complete assignment: update best and return
- Choose unassigned **VARIABLE** X_i
- Order **VALUES** Domain_i of chosen X_i
- For each value v in that order:
 - $\delta \leftarrow \prod_{f_j \in D(x, X_i)} f_j(x \cup \{X_i : v\})$
 - If $\delta = 0$: continue
 - Domains' \leftarrow Domains via **LOOKAHEAD**
 - If any Domains' is empty: continue
 - $\mathsf{Backtrack}(x \cup \{X_i : v\}, w\delta, \mathsf{Domains'})$

Partial assignment weights

Idea: compute weight of partial assignment as we go



Dependent factors

• Partial assignment (e.g., $x = \{WA : R, NT : G\}$)



Definition: dependent factors-

Let $D(x, X_i)$ be set of factors depending on X_i and x but not on unassigned variables.

 $D(\{\mathsf{WA}:\mathsf{R},\mathsf{NT}:\mathsf{G}\},\mathsf{SA})=\{[\mathsf{WA}\neq\mathsf{SA}],[\mathsf{NT}\neq\mathsf{SA}]\}$

Lookahead: forward checking

Key idea: forward checking (one-step lookahead)-

- After assigning a variable X_i, eliminate inconsistent values from the domains of X_i's neighbors.
- If any domain becomes empty, return.



Choosing an unassigned variable



Which variable to assign next?

Key idea: most constrained variable

Choose variable that has the smallest domain.

This example: SA (has only one value)

Ordering values of a selected variable

What values to try for Q?



2+2+2=6 consistent values



1+1+2=4 consistent values

Key idea: least constrained value-

Order values of selected X_i by decreasing number of consistent values of neighboring variables.

When to fail?



Most constrained variable (MCV):

- Must assign every variable
- If going to fail, fail early \Rightarrow more pruning

Least constrained value (LCV):

- Need to choose **some** value
- Choose value that is most likely to lead to solution

When do these heuristics help?

Most constrained variable: useful when some factors are constraints (can prune assignments with weight 0)

$$[x_1 = x_2] \qquad [x_2 \neq x_3] + 2$$

• Least constrained value: useful when **all** factors are constraints (all assignment weights are 1 or 0)

$$[x_1 = x_2] \qquad \qquad [x_2 \neq x_3]$$

• Forward checking: needed to prune domains to make heuristics useful!

Arc consistency



A variable X_i is **arc consistent** with respect to X_j if for each $x_i \in \text{Domain}_i$, there exists $x_j \in \text{Domain}_j$ such that $f(\{X_i : x_i, X_j : x_j\}) \neq 0$ for all factors f whose scope contains X_i and X_j .

Algorithm: enforce arc consistency-

EnforceArcConsistency (X_i, X_j) : Remove values from Domain_i to make X_i arc consistent with respect to X_j .

AC-3

Forward checking: when assign $X_j : x_j$, set $Domain_j = \{x_j\}$ and enforce arc consistency on all neighbors X_i with respect to X_j

AC-3: repeatedly enforce arc consistency on all variables



 X_i

Limitations of AC-3

• AC-3 isn't always effective:



- No consistent assignments, but AC-3 doesn't detect a problem!
- Intuition: if we look locally at the graph, nothing blatantly wrong...

Backtracking search



Greedy search



Beam search



Beam size K = 4

Beam search

Idea: keep $\leq K$ candidate list C of partial assignments

- 🗖 Algorithm: beam search————
Initialize $C \leftarrow [\{\}]$
For each $i = 1, \ldots, n$:
Extend:
$C' \leftarrow \{x \cup \{X_i : v\} : x \in C, v \in Domain_i\}$
Prune:
$C \leftarrow K$ elements of C' with highest weights

Not guaranteed to find maximum weight assignment!

[demo: beamSearch({K:3})]







- Beam size K controls tradeoff between efficiency and accuracy
 - K = 1 is greedy search (O(nb) time)
 - $K = \infty$ is BFS ($O(b^n)$ time)

Backtracking search : DFS :: beam search : pruned BFS

Search strategies

Backtracking/beam search: extend partial assignments



Local search: modify complete assignments





Example: object tracking



[demo]

One small step



Old assignment: (0, 0, 1); how to improve?

$$\begin{array}{ll} (x_1, {\boldsymbol v}, x_3) & \text{weight} \\ (0, 0, 1) & 2 \cdot 2 \cdot 0 \cdot 1 \cdot 1 = 0 \\ (0, 1, 1) & 2 \cdot 1 \cdot 1 \cdot 2 \cdot 1 = 4 \\ (0, 2, 1) & 2 \cdot 0 \cdot 2 \cdot 1 \cdot 1 = 0 \end{array}$$

New assignment: (0, 1, 1)

Exploiting locality



Weight of new assignment $(x_1, \boldsymbol{v}, x_3)$:

 $o_1(x_1)t_1(x_1,v)o_2(v)t_2(v,x_3)o_3(x_3)$

Key idea: locality When evaluating possible re-assignments to X_i , only need to consider the factors that depend on X_i .

Iterated conditional modes (ICM)

Algorithm: iterated conditional modes (ICM)₇

Initialize x to a random complete assignment Loop through i = 1, ..., n until convergence: Compute weight of $x_v = x \cup \{X_i : v\}$ for each v $x \leftarrow x_v$ with highest weight



[demo: iteratedConditionalModes()]

Convergence properties

- Weight(x) increases or stays the same each iteration
- Converges in a finite number of iterations
- Can get stuck in local optima
- Not guaranteed to find optimal assignment!









Algorithm	Strategy	Optimality	Time complexity
Backtracking search	extend partial assignments	exact	exponential
Beam search	extend partial assignments	approximate	linear
Local search (ICM)	modify complete assignments	approximate	linear

Course plan



Machine learning

Definition

Definition: Markov network-

A Markov network is a factor graph which defines a joint distribution over random variables $X = (X_1, \ldots, X_n)$:

$$\mathbb{P}(X = x) = \frac{\mathsf{Weight}(x)}{Z}$$

where $Z = \sum_{x'} \text{Weight}(x')$ is the normalization constant.

x_1	x_2	x_3	Weight(x)	$\mathbb{P}(X=x)$
0	1	1	4	0.15
0	1	2	4	0.15
1	1	1	4	0.15
1	1	2	4	0.15
1	2	1	2	0.08
1	2	2	8	0.31

$$Z = 4 + 4 + 4 + 4 + 2 + 8 = 26$$

Represents uncertainty!

Marginal probabilities

Example question: where was the object at time step 2 (X_2) ?

Definition: Marginal probability The marginal probability of $X_i = v$ is given by: $\mathbb{P}(X_i = v) = \sum_{x:x_i=v} \mathbb{P}(X = x)$

Object tracking example:

x_1	x_2	x_3	Weight(x)	$\mathbb{P}(X=x)$
0	1	1	4	0.15
0	1	2	4	0.15
1	1	1	4	0.15
1	1	2	4	0.15
1	2	1	2	0.08
1	2	2	8	0.31

 $\mathbb{P}(X_2 = 1) = 0.15 + 0.15 + 0.15 + 0.15 = 0.62$ $\mathbb{P}(X_2 = 2) = 0.08 + 0.31 = 0.38$ Note: different than max weight assignment!




Markov networks = factor graphs + probability

- Normalize weights to get probablity distribution
- Can compute marginal probabilities to focus on variables

CSPs	Markov networks
variables	random variables
weights	probabilities
maximum weight assignment	marginal probabilities

Gibbs sampling



[demo]

Search versus sampling

Iterated Conditional Modes	Gibbs sampling
maximum weight assignment	marginal probabilities
choose best value	sample a value
converges to local optimum	marginals converge to correct answer*

*under technical conditions (sufficient condition: all weights positive), but could take exponential time









- Objective: compute marginal probabilities $\mathbb{P}(X_i = x_i)$
- Gibbs sampling: sample one variable at a time, count visitations
- More generally: Markov chain Monte Carlo (MCMC) powerful toolkit of randomized procedures

Course plan



Machine learning

Markov networks versus Bayesian networks

Both define a joint probability distribution over assignments





Markov networks	Bayesian networks
arbitrary factors	local conditional probabilities
set of preferences	generative process



Review: probability

Random variables: sunshine $S \in \{0, 1\}$, rain $R \in \{0, 1\}$

Joint distribution (probabilistic database):

	$s \ r \ \mathbb{I}$	$\mathbb{P}(S=s,R=r)$
	00	0.20
(S, R) =	01	0.08
	10	0.70
	11	0.02

Marginal distribution:

 \mathbb{P}

 $\begin{pmatrix} \text{aggregate rows} \end{pmatrix} \\
\mathbb{P}(S) = \begin{vmatrix} s & \mathbb{P}(S = s) \\ 0 & 0.28 \\ 1 & 0.72 \end{vmatrix}$

Conditional distribution:

(select rows, normalize) $\mathbb{P}(S \mid R = 1) = \begin{bmatrix} s \ \mathbb{P}(S = s \mid R = 1) \\ 0 & 0.8 \\ 1 & 0.2 \end{bmatrix}$

Bayesian network (definition)



Definition: Bayesian network-

Let $X = (X_1, \ldots, X_n)$ be random variables.

A **Bayesian network** is a directed acyclic graph (DAG) that specifies a joint distribution over X as a product of local conditional distributions, one for each node:

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) \stackrel{\mathsf{def}}{=} \prod_{i=1}^n p(x_i \mid x_{\mathsf{Parents}(i)})$$

Probabilistic inference (definition)





Example: if coughing and itchy eyes, have a cold?

$$\mathbb{P}(C \mid H = 1, I = 1)$$



Bayesian network (alarm)



$$p(b) = \epsilon \cdot [b = 1] + (1 - \epsilon) \cdot [b = 0]$$

$$p(e) = \epsilon \cdot [e = 1] + (1 - \epsilon) \cdot [e = 0]$$

$$p(a \mid b, e) = [a = (b \lor e)]$$

$$\mathbb{P}(B = b, E = e, A = a) \stackrel{\text{def}}{=} p(b)p(e)p(a \mid b, e)$$

CS221



Explaining away





Suppose two causes positively influence an effect. Conditioned on the effect, further conditioning on one cause reduces the probability of the other cause.

$$\mathbb{P}(B = 1 \mid A = 1, E = 1) < \mathbb{P}(B = 1 \mid A = 1)$$

Note: happens even if causes are independent!







- Random variables capture state of world
- Directed edges between variables represent dependencies
- Local conditional distributions \Rightarrow joint distribution
- Probabilistic inference: ask questions about world
- Captures reasoning patterns (e.g., explaining away)

Probabilistic programs



- Y Key idea: probabilistic program A randomized program that sets the random variables.

Reduction to Markov networks



Reminder: single factor that connects all parents!

Conditioning on evidence



Markov network:

$$\mathbb{P}(C = c, A = a \mid H = 1, I = 1) = \frac{1}{Z}p(c)p(a)p(h = 1 \mid c, a)p(i = 1 \mid a)$$

Bayesian network with evidence = Markov network with $Z = \mathbb{P}(H = 1, I = 1)$

Solution: run any inference algorithm for Markov networks (e.g., Gibbs sampling)! [demo]

Leveraging additional structure: unobserved leaves



Markov network:

$$\mathbb{P}(C = c, A = a, I = i \mid H = 1) = \frac{1}{Z}p(c)p(a)p(h = 1 \mid c, a)p(i \mid a),$$

where $Z = \mathbb{P}(H = 1)$

Question: $\mathbb{P}(C = 1 \mid H = 1)$

Can we reduce the Markov network before running inference?

Leveraging additional structure: unobserved leaves



Markov network:

$$\begin{split} \mathbb{P}(C = c, A = a \mid H = 1) &= \sum_{i} \mathbb{P}(C = c, A = a, I = i \mid H = 1) \\ &= \sum_{i} \frac{1}{Z} p(c) p(a) p(h = 1 \mid c, a) p(i \mid a) \\ &= \frac{1}{Z} p(c) p(a) p(h = 1 \mid c, a) \sum_{i} p(i \mid a) \\ &= \frac{1}{Z} p(c) p(a) p(h = 1 \mid c, a) \end{split}$$

Throw away any unobserved leaves before running inference!

Leveraging additional structure: independence



Markov network:

$$\mathbb{P}(\mathbf{C} = \mathbf{c} \mid I = 1) = \sum_{a,h} \mathbb{P}(\mathbf{C} = \mathbf{c}, A = a, H = h \mid I = 1)$$
$$= \sum_{a,h} \frac{1}{Z} p(\mathbf{c}) p(a) p(h \mid \mathbf{c}, a) p(i = 1 \mid a)$$
$$= \sum_{a} \frac{1}{Z} p(\mathbf{c}) p(a) p(i = 1 \mid a)$$
$$= p(\mathbf{c}) \sum_{a} \frac{1}{Z} p(a) p(i = 1 \mid a)$$
$$= p(\mathbf{c})$$

Throw away any disconnected components before running inference!







- Condition on evidence (e.g., I = 1)
- Throw away unobserved leaves (e.g., H)
- Throw away disconnected components (e.g., A and I)
- Define Markov network out of remaining factors
- Run your favorite inference algorithm (e.g., manual, Gibbs sampling)

Inference questions



Question (**filtering**):

$$\mathbb{P}(H_2 \mid E_1 = 0, E_2 = 2)$$

Question (**smoothing**):

 $\mathbb{P}(H_2 \mid E_1 = 0, E_2 = 2, E_3 = 2)$

Note: filtering is a special case of smoothing if marginalize unobserved leaves

Lattice representation



- Edge start \Rightarrow $H_1 = h_1$ has weight $p(h_1)p(e_1 \mid h_1)$
- Edge $H_{i-1} = h_{i-1} \Rightarrow H_i = h_i$ has weight $p(h_i \mid h_{i-1})p(e_i \mid h_i)$
- Each path from start to end is an assignment with weight equal to the product of edge weights

Key: $\mathbb{P}(H_i = h_i \mid E = e)$ is the weighted fraction of paths through $|H_i = h_i|$

Forward and backward messages



Putting everything together



Running time: $O(n|\text{Domain}|^2)$

[demo]

Review: inference in Hidden Markov models



Filtering questions:

$$\mathbb{P}(H_1 \mid E_1 = 0)$$

$$\mathbb{P}(H_2 \mid E_1 = 0, E_2 = 2)$$

$$\mathbb{P}(H_3 \mid E_1 = 0, E_2 = 2, E_3 = 2)$$

Problem: many possible location values for H_i



Forward-backward is too slow $(O(n|Domain|^2))...$

Why sampling?





not representative

K with highest weight K sampled from distribution



more representative

Sampling is especially important when there is high uncertainty!

Particle filtering

Algorithm: particle filtering-Initialize $C \leftarrow [\{\}]$ For each $i = 1, \ldots, n$: Propose: $C' \leftarrow \{h \cup \{H_i : h_i\} : h \in C, h_i \sim p(h_i \mid h_{i-1})\}$ Weight: Compute weights $w(h) = p(e_i \mid h_i)$ for $h \in C'$ Resample: $C \leftarrow K$ particles drawn independently from $\frac{w(h)}{\sum_{h' \in C} w(h')}$

[demo: particleFiltering({K:100})]

Step 1: propose

Old particles: $\approx \mathbb{P}(H_1, H_2 \mid E_1 = 0, E_2 = 2)$

 $\{H_1: 0, H_2: 1\} \\ \{H_1: 1, H_2: 2\}$



New particles: $\approx \mathbb{P}(H_1, H_2, H_3 | E_1 = 0, E_2 = 2)$

 $\{H_1: 0, H_2: 1, H_3: 1\}$ $\{H_1: 1, H_2: 2, H_3: 2\}$

Step 2: weight

Old particles: $\approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 1)$

 $\{H_1: 0, H_2: 1: H_3: 1\}$ $\{H_1: 1, H_2: 2: H_3: 2\}$

For each old particle (h_1, h_2, h_3) , weight it by $p(e_3 = 2 \mid h_3)$.



New particles: $\approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 1, E_3 = 2)$

{ $H_1: 0, H_2: 1: H_3: 1$ } (1/4) { $H_1: 1, H_2: 2: H_3: 2$ } (1/2)

Step 3: resample

```
Old particles: \approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 2, E_3 = 2)
```

```
 \{H_1: 0, H_2: 1: H_3: 1\} (1/4) \Rightarrow 1/3 
 \{H_1: 1, H_2: 2: H_3: 2\} (1/2) \Rightarrow 2/3 

Key idea: resampling

Normalize weights and draw K samples to redistribute particles to more promising areas.
```

```
New particles: \approx \mathbb{P}(H_1, H_2, H_3 \mid E_1 = 0, E_2 = 2, E_3 = 2)
```

 $\{H_1: 1, H_2: 2: H_3: 2\}$ $\{H_1: 1, H_2: 2: H_3: 2\}$









$$\mathbb{P}(H_3 \mid E_1 = 0, E_2 = 2, E_3 = 2)$$

• Use particles to represent an approximate distribution

Propose (transitions)	Weight (emissions)	Resample
		Resumple

- Can scale to large number of locations (unlike forward-backward)
- Maintains better particle diversity (compared to beam search)

Where do parameters come from?



c p(c)	a	p(a)
1 ?	1	?
0 ?	0	?

c	a	h	$p(h \mid c, a)$
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

a	i	p(i	a)
0	0	?	
0	1	?	
1	0	?	
1	1	?	

Learning task

Training data-

 \mathcal{D}_{train} (an example is an assignment to X)

Parameters-

 θ (local conditional probabilities)

Parameter sharing

Key idea: parameter sharing-

The local conditional distributions of different variables can share the same parameters.



Impact: more reliable estimates, less expressive model

General case

Bayesian network: variables X_1, \ldots, X_n

Parameters: collection of distributions $\theta = \{p_d : d \in D\}$ (e.g., $D = \{\text{start}, \text{trans}, \text{emit}\}$)

Each variable X_i is generated from distribution p_{d_i} :

$$\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n p_{d_i}(x_i \mid x_{\mathsf{Parents}(i)})$$

Parameter sharing: d_i could be same for multiple i

General case: learning algorithm

Input: training examples $\mathcal{D}_{\mathsf{train}}$ of full assignments

```
Output: parameters \theta = \{p_d : d \in D\}
```

```
Algorithm: count and normalize

Count:

For each x \in D_{train}:

For each variable x_i:

Increment count_{d_i}(x_{Parents(i)}, x_i)

Normalize:

For each d and local assignment x_{Parents(i)}:

Set p_d(x_i \mid x_{Parents(i)}) \propto \text{count}_d(x_{Parents(i)}, x_i)
```

Maximum likelihood

Maximum likelihood objective:

```
\max_{\theta} \prod_{x \in \mathcal{D}_{\mathsf{train}}} \mathbb{P}(X = x; \theta)
```



Closed form — no iterative optimization!
Review: maximum likelihood

 $\mathcal{D}_{\mathsf{train}} = \{(\mathsf{d},4), (\mathsf{d},4), (\mathsf{d},5), (\mathsf{c},1), (\mathsf{c},5)\}$

				g	r	$count_R(g,r)$	$p_R(r \mid g)$
	g	$count_G(g)$	$p_G(g)$	d	4	2	2/3
:	d	3	3/5	d	5	1	1/3
	С	2	2/5	с	1	1	1/2
				с	5	1	1/2

Do we really believe that $p_R(r = 2 \mid g = c) = 0$?

Overfitting!

Laplace smoothing

Key idea: maximum likelihood with Laplace smoothing-

For each distribution d and partial assignment $(x_{Parents(i)}, x_i)$:

Add λ to count_d($x_{\text{Parents}(i)}, x_i$).

Further increment counts $\{\text{count}_d\}$ based on $\mathcal{D}_{\text{train}}$.

Hallucinate λ occurrences of each local assignment

Interplay between smoothing and data

Larger $\lambda \Rightarrow$ more smoothing \Rightarrow probabilities closer to uniform

g	$count_G(g)$	$p_G(g)$	g	$count_G(g)$	$p_G(g)$
d	1/2+1	3/4	d	1 +1	2/3
С	1/2	1/4	С	1	1/3

Data wins out in the end (suppose only see g = d):

g	$count_G(g)$	$p_G(g)$	g	$count_G(g)$	$p_G(g)$
d	1 +1	2/3	d	1 +998	0.999
С	1	1/3	С	1	0.001

Motivation



Genre $G \in \{ drama, comedy \}$ Jim's rating $R_1 \in \{1, 2, 3, 4, 5\}$ Martha's rating $R_2 \in \{1, 2, 3, 4, 5\}$

If observe all the variables: maximum likelihood = count and normalize

$$\mathcal{D}_{\mathsf{train}} = \{(\mathsf{d}, 4, 5), (\mathsf{d}, 4, 4), (\mathsf{d}, 5, 3), (\mathsf{c}, 1, 2), (\mathsf{c}, 5, 4)\}$$

What if we **don't observe** some of the variables?

$$\mathcal{D}_{\mathsf{train}} = \{(?, 4, 5), (?, 4, 4), (?, 5, 3), (?, 1, 2), (?, 5, 4)\}$$

Expectation Maximization (EM)

Intuition: generalization of the K-means algorithm

```
cluster centroids = parameters \theta cluster assignments = hidden variables H
```

```
Variables: H is hidden, E = e is observed
```



Maximum likelihood (count and normalize) on weighted examples to get θ



Summary



Maximum marginal likelihood: max $\prod \mathbb{P}(E = e; \theta)$

 $\max_{\theta} \prod_{e \in \mathcal{D}_{\text{train}}} \mathbb{P}(E = e; \theta)$

EM algorithm:

← probabilistic inference (E-step)

hidden variables q(h)



parameters θ

count and normalize (M-step) \Rightarrow

Applications: decipherment, phylogenetic reconstruction, crowdsourcing

Course plan



Machine learning

Modeling paradigms

State-based models: search problems, MDPs, games

Applications: route finding, game playing, etc. *Think in terms of* **states, actions, and costs**

Variable-based models: CSPs, Bayesian networks

Applications: scheduling, tracking, medical diagnosis, etc. *Think in terms of* variables and factors

Logic-based models: propositional logic, first-order logic

Applications: theorem proving, verification, reasoning Think in terms of logical formulas and inference rules



Natural language

Example:

- A **dime** is better than a **nickel**.
- A nickel is better than a penny.
- Therefore, a **dime** is better than a **penny**.

Example:

- A penny is better than nothing.
- Nothing is better than world peace.
- Therefore, a **penny** is better than **world peace**???

Natural language is slippery...

Two goals of a logic language

• **Represent** knowledge about the world



• **Reason** with that knowledge



Ingredients of a logic

Syntax: defines a set of valid **formulas** (Formulas)

Example: Rain \land Wet

Semantics: for each formula, specify a set of **models** (assignments / configurations of the world)



Inference rules: given f, what new formulas g can be added that are guaranteed to follow $(\frac{f}{g})$?

Example: from Rain \land Wet, derive Rain

Syntax versus semantics

Syntax: what are valid expressions in the language?

Semantics: what do these expressions mean?

Different syntax, same semantics (5):

 $2 + 3 \Leftrightarrow 3 + 2$

Same syntax, different semantics (1 versus 1.5):

3 / 2 (Python 2.7) $\Leftrightarrow 3 / 2$ (Python 3)

Propositional logic



Logics

- Propositional logic with only Horn clauses
- Propositional logic
- Modal logic
- First-order logic with only Horn clauses
- First-order logic
- Second-order logic

```
Key idea: tradeoff
Balance expressivity and computational efficiency.
```

. . .

Syntax of propositional logic

Propositional symbols (atomic formulas): A, B, C

Logical connectives: $\neg, \land, \lor, \rightarrow, \leftrightarrow$

Build up formulas recursively—if f and g are formulas, so are the following:

- Negation: $\neg f$
- Conjunction: $f \wedge g$
- Disjunction: $f \lor g$
- Implication: $f \to g$
- Biconditional: $f \leftrightarrow g$

Syntax of propositional logic





Model



A **model** w in propositional logic is an **assignment** of truth values to propositional symbols.

Example:

- 3 propositional symbols: A, B, C
- $2^3 = 8$ possible models w:

 $\{A:0,B:0,C:0\} \\ \{A:0,B:0,C:1\} \\ \{A:0,B:1,C:0\} \\ \{A:0,B:1,C:1\} \\ \{A:1,B:0,C:0\} \\ \{A:1,B:0,C:1\} \\ \{A:1,B:1,C:0\} \\ \{A:1,B:1,C:1\} \\ \{A:1,B:1,C:1\} \}$

Interpretation function



Definition: interpretation function-

Let f be a formula.

Let w be a model.

An interpretation function $\mathcal{I}(f, w)$ returns:

- true (1) (say that w satisfies f)
- false (0) (say that w does not satisfy f)



Interpretation function: example



Formula represents a set of models

So far: each formula f and model w has an interpretation $\mathcal{I}(f, w) \in \{0, 1\}$



Let $\mathcal{M}(f)$ be the set of **models** w for which $\mathcal{I}(f, w) = 1$.



Models: example

Formula:

 $f = \mathsf{Rain} \lor \mathsf{Wet}$

Models:



Knowledge base

Definition: Knowledge base-

A **knowledge base** KB is a set of formulas representing their conjunction / intersection:

$$\mathcal{M}(\mathsf{KB}) = \bigcap_{f \in \mathsf{KB}} \mathcal{M}(f).$$

Intuition: KB specifies constraints on the world. $\mathcal{M}(\text{KB})$ is the set of all worlds satisfying those constraints.

Let $KB = \{Rain \lor Snow, Traffic\}.$



Adding to the knowledge base

Adding more formulas to the knowledge base:



Shrinks the set of models:



How much does $\mathcal{M}(\textbf{KB})$ shrink?

[whiteboard]

Entailment



Intuition: f added no information/constraints (it was already known).



Example: Rain \land Snow \models Snow

Contradiction



Intuition: f contradicts what we know (captured in KB).

Definition: contradiction KB contradicts f iff $\mathcal{M}(KB) \cap \mathcal{M}(f) = \emptyset$.

Example: Rain \land Snow contradicts \neg Snow

Contradiction and entailment





Entailment:



Proposition: contradiction and entailment₇

KB contradicts f iff KB entails $\neg f$.

Contingency



Intuition: f adds non-trivial information to KB

 $\emptyset \subsetneq \mathcal{M}(\mathsf{KB}) \cap \mathcal{M}(f) \subsetneq \mathcal{M}(\mathsf{KB})$

Example: Rain and Snow

Satisfiability

Definition: satisfiability-

A knowledge base KB is **satisfiable** if $\mathcal{M}(KB) \neq \emptyset$.

Reduce Ask[f] and Tell[f] to satisfiability:



Model checking

Checking satisfiability (SAT) in propositional logic is special case of solving CSPs!

Mapping:

propositional symbol	\Rightarrow	variable
formula	\Rightarrow	constraint
model	\Leftarrow	assignment

Model checking

Example: model checking- $\mathsf{KB} = \{A \lor B, B \leftrightarrow \neg C\}$ Propositional symbols (CSP variables): $\{A, B, C\}$ CSP: $A \lor B$ $B \leftrightarrow \neg C$ В Α Consistent assignment (satisfying model): ${A:1, B:0, C:1}$

Propositional logic



Inference rules

Example of making an inference:

It is raining. (Rain) If it is raining, then it is wet. (Rain \rightarrow Wet) Therefore, it is wet. (Wet)



Inference framework



Wey idea: inference rules Rules operate directly on **syntax**, not on **semantics**.

Inference algorithm

Algorithm: forward inference-

Input: set of inference rules Rules. Repeat until no changes to KB: Choose set of formulas $f_1, \ldots, f_k \in KB$. If matching rule $\frac{f_1, \ldots, f_k}{g}$ exists: Add g to KB.

Definition: derivation

KB derives/proves f (KB \vdash f) iff f eventually gets added to KB.

Desiderata for inference rules

Semantics

Interpretation defines **entailed/true** formulas: $KB \models f$:



Syntax:

Inference rules **derive** formulas: $KB \vdash f$

```
How does \{f : \mathsf{KB} \models f\} relate to \{f : \mathsf{KB} \vdash f\}?
```

Truth



$$\{f: \mathsf{KB} \models f\}$$












Soundness and completeness

The truth, the whole truth, and nothing but the truth.

- **Soundness**: nothing but the truth
- **Completeness**: whole truth

Completeness: example

Recall completeness: inference rules derive all entailed formulas (f such that $KB \models f$)

```
Example: Modus ponens is incomplete
Setup:
     \mathsf{KB} = \{\mathsf{Rain}, \mathsf{Rain} \lor \mathsf{Snow} \to \mathsf{Wet}\}
     f = Wet
     \mathsf{Rules} = \left\{ \frac{f, \quad f \to g}{q} \right\} \text{ (Modus ponens)}
Semantically: KB \models f (f is entailed).
Syntactically: KB \not\vdash f (can't derive f).
                    Incomplete!
```

Fixing completeness

Option 1: Restrict the allowed set of formulas



Definite clauses

• Definition: Definite clause A **definite clause** has the following form: $(p_1 \land \dots \land p_k) \rightarrow q$ where p_1, \dots, p_k, q are propositional symbols.

Intuition: if p_1, \ldots, p_k hold, then q holds. Example: (Rain \land Snow) \rightarrow Traffic Example: Traffic Non-example: \neg Traffic Non-example: (Rain \land Snow) \rightarrow (Traffic \lor Peaceful)

Horn clauses



A Horn clause is either:

- a definite clause $(p_1 \land \cdots \land p_k \rightarrow q)$
- a goal clause $(p_1 \land \cdots \land p_k \rightarrow \mathsf{false})$

Example (definite): $(Rain \land Snow) \rightarrow Traffic$

Example (goal): Traffic \land Accident \rightarrow false

equivalent: $\neg(\mathsf{Traffic} \land \mathsf{Accident})$

Modus ponens

Inference rule:



Example:

Completeness of modus ponens

Theorem: Modus ponens on Horn clauses-

Modus ponens is **complete** with respect to Horn clauses:

- Suppose KB contains only Horn clauses and p is an entailed propositional symbol.
- Then applying modus ponens will derive *p*.

Upshot:

$\mathsf{KB} \models p$ (entailment) is the same as $\mathsf{KB} \vdash p$ (derivation)!

Resolution [Robinson, 1965]

General clauses have any number of literals:

 $\neg A \lor B \lor \neg C \lor D \lor \neg E \lor F$ **Example: resolution inference rule** $\frac{\text{Rain} \lor \text{Snow}, \quad \neg \text{Snow} \lor \text{Traffic}}{\text{Rain} \lor \text{Traffic}}$

Definition: resolution inference rule
$$\frac{f_1 \lor \cdots \lor f_n \lor p, \quad \neg p \lor g_1 \lor \cdots \lor g_m}{f_1 \lor \cdots \lor f_n \lor g_1 \lor \cdots \lor g_m}$$

Conjunctive normal form

So far: resolution only works on clauses...but that's enough!

Conjunctive normal form (CNF) Conjunctive normal form (CNF)

A CNF formula is a conjunction of clauses.

Example: $(A \lor B \lor \neg C) \land (\neg B \lor D)$

Equivalent: knowledge base where each formula is a clause



Every formula f in propositional logic can be converted into an equivalent CNF formula f':

$$\mathcal{M}(f) = \mathcal{M}(f')$$

Conversion to CNF: general

Conversion rules:

- Eliminate \leftrightarrow : $\frac{f \leftrightarrow g}{(f \rightarrow g) \land (g \rightarrow f)}$
- Eliminate \rightarrow : $\frac{f \rightarrow g}{\neg f \lor g}$
- Move \neg inwards: $\frac{\neg (f \land g)}{\neg f \lor \neg g}$
- Move \neg inwards: $\frac{\neg (f \lor g)}{\neg f \land \neg g}$
- Eliminate double negation: $\frac{\neg \neg f}{f}$
- Distribute \lor over \land : $\frac{f \lor (g \land h)}{(f \lor g) \land (f \lor h)}$

Resolution algorithm

Recall: relationship between entailment and contradiction (basically "proof by contradiction")

 $\mathsf{KB} \models f \qquad \longleftarrow \qquad \mathsf{KB} \cup \{\neg f\} \text{ is unsatisfiable}$



Algorithm: resolution-based inference-

- Add $\neg f$ into KB.
- Convert all formulas into **CNF**.
- Repeatedly apply resolution rule.
- Return entailment iff derive false.

Resolution: example

$$\mathsf{KB}' = \{A \to (B \lor C), A, \neg B, \neg C\}$$

Convert to CNF:

 $\mathsf{KB}' = \{\neg A \lor B \lor C, A, \neg B, \neg C\}$

Repeatedly apply **resolution** rule:



Limitations of propositional logic

All students know arithmetic.

```
AliceIsStudent \rightarrow AliceKnowsArithmetic
```

 $\mathsf{BobIsStudent} \to \mathsf{BobKnowsArithmetic}$

Propositional logic is very clunky. What's missing?

- Objects and predicates: propositions (e.g., AliceKnowsArithmetic) have more internal structure (alice, Knows, arithmetic)
- Quantifiers and variables: *all* is a quantifier which applies to each person, don't want to enumerate them all...

. . .

First-order logic: examples

Alice and Bob both know arithmetic.

Knows(alice, arithmetic) \land Knows(bob, arithmetic)

All students know arithmetic.

 $\forall x \operatorname{Student}(x) \rightarrow \operatorname{Knows}(x, \operatorname{arithmetic})$

Syntax of first-order logic

Terms (refer to objects):

- Constant symbol (e.g., arithmetic)
- Variable (e.g., x)
- Function of terms (e.g., Sum(3, x))

Formulas (refer to truth values):

- Atomic formulas (atoms): predicate applied to terms (e.g., Knows(x, arithmetic))
- Connectives applied to formulas (e.g., $Student(x) \rightarrow Knows(x, arithmetic)$)
- Quantifiers applied to formulas (e.g., $\forall x \operatorname{Student}(x) \to \operatorname{Knows}(x, \operatorname{arithmetic})$)

Quantifiers

Universal quantification (\forall) :

```
Think conjunction: \forall x P(x) is like P(A) \land P(B) \land \cdots
```

Existential quantification (\exists) :

Think disjunction: $\exists x P(x)$ is like $P(A) \lor P(B) \lor \cdots$

Some properties:

- $\neg \forall x P(x)$ equivalent to $\exists x \neg P(x)$
- $\forall x \exists y \operatorname{Knows}(x, y)$ different from $\exists y \forall x \operatorname{Knows}(x, y)$

Natural language quantifiers

Universal quantification (\forall):

Every student knows arithmetic.

 $\forall x \operatorname{Student}(x) \rightarrow \operatorname{Knows}(x, \operatorname{arithmetic})$

Existential quantification (\exists) :

Some student knows arithmetic.

 $\exists x \operatorname{Student}(x) \land \operatorname{Knows}(x, \operatorname{arithmetic})$

Note the different connectives!

Models in first-order logic



Definition: model in first-order logic-

A model w in first-order logic maps:

• constant symbols to objects

 $w(alice) = o_1, w(bob) = o_2, w(arithmetic) = o_3$

• predicate symbols to tuples of objects

 $w(Knows) = \{(o_1, o_3), (o_2, o_3), \dots\}$

Graph representation of a model

If only have unary and binary predicates, a model w can be represented as a directed graph:



- Nodes are objects, labeled with constant symbols
- Directed edges are binary predicates, labeled with predicate symbols; unary predicates are additional node labels

A restriction on models

John and Bob are students.

$\mathsf{Student}(\mathsf{john}) \land \mathsf{Student}(\mathsf{bob})$



- Unique names assumption: Each object has at most one constant symbol. This rules out w_2 .
- Domain closure: Each object has **at least one** constant symbol. This rules out w_3 . Point:

constant symbol



Propositionalization

If one-to-one mapping between constant symbols and objects (unique names and domain closure),

first-order logic is syntactic sugar for propositional logic:

Knowledge base in first-order logic Student(alice) \land Student(bob) $\forall x$ Student(x) \rightarrow Person(x) $\exists x$ Student(x) \land Creative(x)

Knowledge base in propositional logic-

 $\begin{array}{l} \mathsf{Studentalice} \land \mathsf{Studentbob} \\ (\mathsf{Studentalice} \rightarrow \mathsf{Personalice}) \land (\mathsf{Studentbob} \rightarrow \mathsf{Personbob}) \\ (\mathsf{Studentalice} \land \mathsf{Creativealice}) \lor (\mathsf{Studentbob} \land \mathsf{Creativebob}) \end{array}$

Point: use any inference algorithm for propositional logic!

Substitution

 $\mathsf{Subst}[\{x/\mathsf{alice}\}, P(x)] = P(\mathsf{alice})$

```
\mathsf{Subst}[\{x/\mathsf{alice}, y/z\}, P(x) \land K(x, y)] = P(\mathsf{alice}) \land K(\mathsf{alice}, z)
```



Definition: Substitution-

A substitution θ is a mapping from variables to terms.

Subst $[\theta, f]$ returns the result of performing substitution θ on f.

Unification

```
Unify[Knows(alice, arithmetic), Knows(x, arithmetic)] = {x/alice}
```

```
\mathsf{Unify}[\mathsf{Knows}(\mathsf{alice}, y), \mathsf{Knows}(x, z)] = \{x/\mathsf{alice}, y/z\}
```

```
\mathsf{Unify}[\mathsf{Knows}(\mathsf{alice}, y), \mathsf{Knows}(\mathsf{bob}, z)] = \mathsf{fail}
```

 $\mathsf{Unify}[\mathsf{Knows}(\mathsf{alice}, y), \mathsf{Knows}(x, F(x))] = \{x/\mathsf{alice}, y/F(\mathsf{alice})\}$



Unification takes two formulas f and g and returns a substitution θ which is the most general unifier:

```
\mathsf{Unify}[f,g] = \theta such that \mathsf{Subst}[\theta,f] = \mathsf{Subst}[\theta,g]
```

or "fail" if no such θ exists.

Modus ponens



Get most general unifier θ on premises:

• $\theta = \text{Unify}[a'_1 \wedge \cdots \wedge a'_k, a_1 \wedge \cdots \wedge a_k]$

Apply θ to conclusion:

• Subst $[\theta, b] = b'$



Theorem: completeness-

Modus ponens is complete for first-order logic with only Horn clauses.

Theorem: semi-decidability-

First-order logic (even restricted to only Horn clauses) is semi-decidable.

- If $KB \models f$, forward inference on complete inference rules will prove f in finite time.
- If $KB \not\models f$, no algorithm can show this in finite time.