

Probability Stuff

- RV → All Random Variables have domains
- RV is some aspect of the world we are uncertain of
- Distributions must be $\forall x P(x) \geq 0$ and $\sum_x P(x) = 1$
- Joint Distributions have ≥ 1 RV, like temp and rain
- Probabilistic Inference → Computing a probability from other probabilities (conditional from joint) → New observations change and update belief
- RVs are independent if $P(x, y) = P(x)P(y)$ $P(x, y) = P(x|y)P(y) = P(y|x)P(x) \rightarrow P(x|y) = \frac{P(y|x)}{P(y)}$

Naive Bayes → $P(\text{class} | \text{features}) = \frac{P(\text{features} | \text{class})}{P(\text{features})} \cdot P(\text{class})$ → Assumes all features are independent of each other

- Bayes Net Models often used for classification → Input and output are both random variables
- Empirical Risk Minimization → Optimization Problem to find best results on test data from training data
- Also concerned about over and underfitting → Aim to reduce variance and limit hypothesis complexity → Relative Frequency Params = count
- Parameter Estimation of PVs → Using training data, maximize likelihood of the data (learning!)
- Laplace Smoothing → $P(\text{feature} | \text{class}) = \frac{\text{count}(\text{feature} | \text{class}) + 1}{\text{count}(\text{class}) + n \text{ features}}$ → Can also extend to k $P_{\text{LAP}}(x|y) = \frac{c(x, y) + k}{c(y) + k(\text{count}(\text{class}))}$ Bad when $|X|$ or $|Y|$ are really big
- Tune Parameters on Training Set, tune hyperparameters on different data (held-out) then re-train and test → Choose best value after test
- Baselines help determine how hard a task is and define good accuracy | Confidence of Probabilistic Classifier: $\max_y P(y | x)$ ↗ not guaranteed
- Calibration → Weak Calibration: High Confidence = High Accuracy ; Strong Calibration: Confidence predicts accuracy rate

Linear Classifiers → Inputs are feature values, features have weights, sum is activation = $\sum_i w_i \cdot f_i(x) = [w \cdot f(x)]$ Dot Product!

- Start with weights = 0
 - For each training instance: **Binary**
 - Pick up training examples one by one
 - Predict with current weights
- Perceptron** $y = \arg \max_y w_y \cdot f(x)$
- If correct, no change! **MultiClass**
 - If wrong: lower score of wrong answer, raise score of right answer
- $w_y = w_y - f(x)$
- $w_{y^*} = w_{y^*} + f(x)$

Perceptrons have Separability, Convergence, and a mistake bound
true if some params set training set perfect

Even the best linear Perceps will make mistakes, so go probabilistic
Using Sigmoid $\Phi(z) = \frac{1}{1+e^{-z}}$

Best w? Batch GA

Maximum likelihood estimation: Stochastic vs random

$$\max_w I(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}, w)$$

$$w \leftarrow w + \eta \nabla I(w)$$

with:

$$P(y^{(i)} = +1 | x^{(i)}, w) = \frac{1}{1 + e^{-w^T f(x^{(i)})}}$$

$$P(y^{(i)} = -1 | x^{(i)}, w) = \frac{1}{1 + e^{w^T f(x^{(i)})}}$$

= Logistic Regression

Best w? Batch GA

Maximum likelihood estimation:

$$\max_w I(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}, w)$$

with:

$$P(y^{(i)} | x^{(i)}, w) = \frac{e^{w^T f(x^{(i)})}}{\sum_j e^{w^T f(x^{(j)})}}$$

= Multi-Class Logistic Regression

20

Gradient Ascent → Type of continuous optimization through partial derivatives and go in direction of steepest gradient

- Neural Networks: Directed graph of neurons; output of layer becomes next input
- At test time, given x_1 and x_0 , compute output. Training → give many x_i and x_0 with outputs and estimate weights
- The neuron is a function that takes a vector of inputs and transforms them into an output
- We learn the weights through supervised learning | Loss func. = $\frac{1}{2} \sum_j (y_j - f(x_j))^2$ There is always a hidden bias
 - We can also use gradient descent to minimize loss function → $w_{\text{new}} = w_{\text{old}} - \alpha \nabla E(w_{\text{old}})$, where α = learning rate

The complete recipe for training a single sigmoid neuron

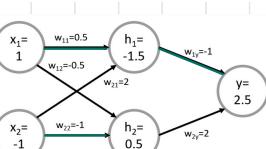
- Start with some initial weights $w = [w_0, w_1, \dots]$
- Repeat until convergence
 - For each input x and the current weights w to compute the output of the neuron $f(w^T x)$
 - Compute the loss $E(w)$ between $f(w^T x)$ and y
 - Compute the gradient of the loss wrt each w_i
 - slightly more complicated equation
- Update each weight w_i based on the gradient

Same idea as single neuron

- Start with some initial weights $w = [w_0, w_1, \dots]$
- Repeat until convergence
 - For each input x , update the current weights w to compute the output of the network $f(w)$
 - Compute the loss $E(w)$ between $f(w)$ and y
 - Compute the gradient of the loss wrt each w_i
 - slightly more complicated equation
- Update each weight w_i based on the gradient

Gradient of a weight w depends only on

- output of the neuron it is connected to
- gradient from a neuron in the next layer
- weight of the edge connected these two neurons



GD will try to increase highlighted weights if desired output is 0

Only unknowns are the weights. Individual needs to design # of layers, neurons / layer, neuron connections, activation function, loss function, and what GD method to use + learning hyperparameter(s)

A probabilistic model is a joint distribution over a set of random variables

Probabilistic models:

• Random variables with domains

• Assignments are called outcomes

• Joint distributions say whether assignments (outcomes) are likely

• Normalization is 1.0

• Ideally: only certain variables directly interact

• Constraint satisfaction problems:

• Constraints with domains

• Constraints: state whether assignments are valid

• Ideally: only certain variables directly interact

Distribution over T,W

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

Constraint over T,W

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

An event is a set of outcomes

Marginal Distributions eliminate ≥ 1 RV

$$P(a|b) = \frac{P(a,b)}{P(b)}$$

Inference is $O(d^n)$ Time and Space

Normalization Trick → Get table for condition

Smooth / Regularize Estimates
Ensure Generalization

Relative Frequency Params = count

Bayes Nets → A technique for describing complex joint distributions using simple, local distributions (Graphical Models)

Unconditional independence is very rare, conditional independence is our most basic form of knowledge on uncertain environments

$X \perp\!\!\!\perp Y|Z$ iff $\forall x,y,z: P(x,y|z) = P(x|z)P(y|z)$ iff $P(x|z,y) = P(x|z)$ Bayes Net nodes can be assigned or unassigned

If there are no interactions between variables → Absolute independence | Chain Rule: $P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots$

Bayes' Net Semantics



- A set of nodes, one per variable X
- A directed, acyclic graph
- A conditional distribution for each node
 - A collection of distributions over X_i for each combination of parents' values
 $P(X_i|x_1, \dots, x_n)$
 - CPT: conditional probability table
 - Description of a noisy "causal" process



$$P(X|A_1, \dots, A_n)$$

A Bayes net = Topology (graph) + Local Conditional Probabilities

Probabilities in BNs



- Bayes' nets implicitly encode joint distributions
 - As a product of local conditional distributions
 - To see what probability a BN gives to a full assignment, multiply all the relevant conditional distributions together:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i|\text{parents}(X_i))$$

- Example: $P(+\text{cavity}, +\text{cold}, -\text{toothache})$

Probabilities in BNs



- Why are we guaranteed that setting $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i|\text{parents}(X_i))$ results in a proper joint distribution?
- Chain rule (valid for all distributions): $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i|x_1, \dots, x_{i-1})$
- Assume conditional independencies: $P(x_i|x_1, \dots, x_{i-1}) = P(x_i|\text{parents}(X_i))$
- Consequence: $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i|\text{parents}(X_i))$
- Not every BN can represent every joint distribution
 - The topology enforces certain conditional independencies

When BNs reflect true causal patterns, they are simpler, easier to think about | Topology of BN encodes conditional indep.

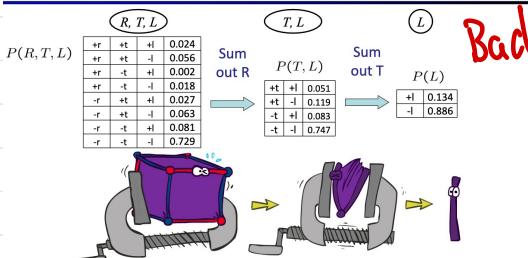
$[X] \rightarrow [Y] \rightarrow [Z]$ X and Z are not necessarily independent, we can use D-separation | Given Y , X is independent of Z

Once again, X is not guaranteed to be independent of Z unless we are given Y (Common Cause)

Dashed nodes are evidence nodes
If two nodes are connected by an undirected path and are not blocked by a shaded node, they are conditionally independent

Given a BN, run the D-separation algorithm to build a list of conditional independencies → Set of distributions represented

Multiple Elimination



Evidence

- If evidence, start with factors that select that evidence
- No evidence uses these initial factors:

$P(R)$	$P(T R)$	$P(L T)$
++ 0.1	++ 0.8	++ 0.3
- - 0.9	- - 0.2	- - 0.7
	- - 0.1	- - 0.1
	- - 0.9	- - 0.9

- Computing $P(L|r)$ the initial factors become:

$P(+r)$	$P(T r)$	$P(L T)$
++ 0.1	++ 0.8	++ 0.3
- - 0.9	- - 0.2	- - 0.7
	- - 0.1	- - 0.1
	- - 0.9	- - 0.9

- We eliminate all vars other than query + evidence



Marginalizing Early! (aka VE)

$P(R)$	Join R	$P(R, T)$	Sum out R	$P(T)$	Join T	$P(T, L)$	Sum out T
$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$
$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$
$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$
$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$	$\begin{array}{cc} ++ & ++ \\ ++ & - - \end{array}$

C	P(C)
red	0.6
green	0.1
blue	0.3

Call rand() n times:
0 ≤ rand ≤ 0.6 → R
0.6 ≤ rand ≤ 0.7 → G
0.7 ≤ rand ≤ 1.0 → B

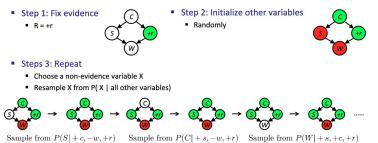
Rejection Sampling:
Reject / Throw out samples / Stop compute when a condition is not met. If $P(C|ts) < 0$, stop early if we get ts.

From uniform dist.

Prior Sampling

Likelihood weighting → Only generate samples with a certain condition by weighting inconsistencies with evidence higher than things that are inconsistent, calculate weighted average

Gibbs Sampling Example: $P(S \mid +r)$



Decision Networks

Umbrella = leave

$$\text{EU(leave)} = \sum_w P(w) U(\text{leave}, w)$$

$$= 0.7 \cdot 100 + 0.3 \cdot 0 = 70$$

Umbrella = take

$$\text{EU(take)} = \sum_w P(w) U(\text{take}, w)$$

$$= 0.7 \cdot 20 + 0.3 \cdot 70 = 35$$

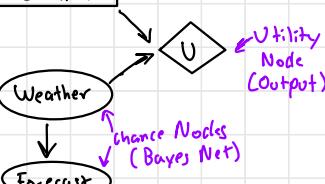
Optimal decision = leave

$$\text{MEU}(a) = \max_a \text{EU}(a) = 70$$

8

Decision Networks \rightarrow We want to get maximum expected utility MEU

Umbrella \leftarrow Actions (observed Evidence)



Decision Networks

Action selection

- Instantiate all evidence
- Set action node(s) each to true
- Calculate posterior for all parents of utility node, given the evidence
- Calculate expected utility for each action
- Choose maximizing action

7

Important to understand the value of acquiring evidence. After gaining information, we can calculate the gain in utility. Find the MEU given info, then subtract original MEU to see the gain. Cost analysis comes from this.

No such thing as a noisy variable

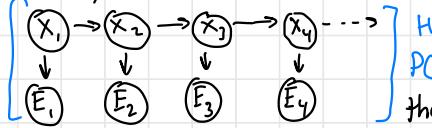
Hidden Markov Models \rightarrow A markov chain over states X , while observing outputs at each step

$(X_1 \rightarrow X_2 \rightarrow X_3 \rightarrow X_4 \rightarrow \dots)$ Parameters (transition probabilities) specify how the state evolves over time Forward Sim.
 $P(X_t)$ and $P(X_t \mid X_{t-1})$ First order Markov Property \rightarrow Each time step only depends on the previous one

$$P(x_t) = \sum_{x_{t-1}} P(x_t \mid x_{t-1}) P(x_{t-1})$$

The chain is just a Bayes Net and we can use BN properties if we truncate at a certain length

Stationary Distribution $P_{\infty}(X) = P_{00\dots}(X) = \sum_i P(X_1|i) P_{\infty}(x)$



Hidden Markov Model Defined by Initial Distribution $P(X_1)$, transitions $P(X_t \mid X_{t-1})$ and emissions $P(E_t \mid X_t)$. HMM future states depend on the past and current observation is independent of everything given state

Observation

Assume we have current belief $P(X \mid \text{previous evidence})$:

$$B'(X_{t+1}) = P(X_{t+1} | e_{1:t})$$

Then, after evidence comes in:

$$\begin{aligned} P(X_{t+1} | e_{1:t+1}) &= P(X_{t+1}, e_{t+1} | e_{1:t}) / P(e_{t+1} | e_{1:t}) \\ &\propto_{e_{t+1}} P(X_{t+1}, e_{t+1} | e_{1:t}) \\ &= P(e_{t+1} | e_{1:t}, X_{t+1}) P(X_{t+1} | e_{1:t}) \\ &= P(e_{t+1} | X_{t+1}) P(X_{t+1} | e_{1:t}) \end{aligned}$$

Or, compactly:

$$B(X_{t+1}) \propto_{e_{t+1}} P(e_{t+1} | X_{t+1}) B'(X_{t+1})$$



Online Belief Updates

Every time step, we start with current $P(X \mid \text{evidence})$

We update for time:

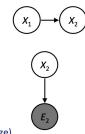
$$P(x_t | e_{1:t-1}) = \sum_{x_{t-1}} P(x_{t-1} | e_{1:t-1}) \cdot P(x_t | x_{t-1})$$

We update for evidence:

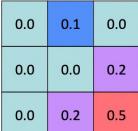
$$P(x_t | e_{1:t}) \propto_{e_t} P(x_t | e_{1:t-1}) \cdot P(e_t | x_t)$$

The forward algorithm does both at once (and doesn't normalize)

43



Particle Filtering used when $|X|$ is too big for exact inference, so we take an approximate inference (Particle = Sample)



We can track the movement of the particle by sampling $P(x'|x) \sim x'$. If we have enough particles, the distribution should match almost exactly

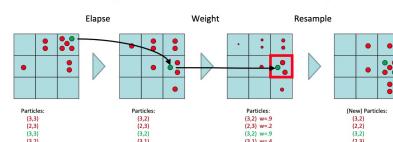
Weights picked by product of likelihood of the evidence given the state represented by the particle and ancestor weight

Instead of a massive state space, budget some # of particles and store their locations To get probability distribution iterate over non-empty states. Choose $P << |X|$

To get probability distribution iterate over non-empty states. Choose $P << |X|$

Recap: Particle Filtering

Particles: track samples of states rather than an explicit distribution



Search Problems have

- State Space
 - Successor Function
 - Start State and goal test
 - Solution is a set of steps to go from Start \rightarrow Goal
- State Space Graph
- Nodes are world states
 - Arcs are successors/transitions
- Search Trees are path representations from State Space Graphs

Greedy Search

- Expand the node closest to goal per heuristic

A* Search (Optimal)

- Combination of UCS and Greedy
 - UCS does backward costs, greedy does forward
 - A* is the sum of these
 - Stop when goal dequeued
 - Mainly expands towards the goal, but hedges bets to ensure optimality
- Heuristic Admissibility
- Admissible if heuristic value is \leq true cost

Depth First Search

- LIFO Stack implementation fringe
- Time: If depth is finite, $O(b^n)$
- Space: $O(bm)$
- Complete only if we prevent cycles
- Not optimal, only finds leftmost solutions

Breadth First Search

- FIFO Queue Fringe
- Process all nodes above shallowest solutions
- Time and Space: $O(b^d)$
- Complete because solution exists
- Optimal only if costs are all equal

Iterative Deepening combines DFS Space with BFS time
Uniform Cost Search

- Expand Cheapest Node First
- Fringe is a priority queue based on cumulative cost
- Processes all nodes with cost less than cheapest solution
- $T + S: O(b^{C/\epsilon}) \rightarrow C = \text{solution cost}, \epsilon = \text{minimum arc cost}$
- Complete and Optimal
- Explores in all directions, no information about goal location

Hueristic Estimates how close a state is to the goal

- Specific to problem

An arc is consistent if for all vars in the tail domain, there is a value in the head domain that doesn't violate any constraints

Graph Search

- Idea is to never expand a state more than once

Estimated heuristic costs should always \leq true cost

Tree Search

- A* Optimal if heuristics are admissible, UCS is the case $h=0$

Graph Search **Consistency = Admissibility**

- A* Optimal if heuristic is consistent, UCS optimal ($h=0$)

Constraint Satisfaction Problems Least Cons. Val.

Search Problem Subsets + Min. Rem. Values

- States with variables, domain, goal test checks for variable combinations

• Constraint Graph

Nodes = Vars; Arcs = Limits

Backtracking Search

• One var at a time, check cons. as you go

• Filtering gets rid of bad domains via forward checking

Arc Consistency is a good way too