

Probability

- marginal: $\sum_n P(A=a, B) = P(B)$
- Sum $\sum_n P(A=a) = 1$
 - conditional $P(A|C) = \frac{P(A,C)}{P(C)}$
 - independence $P(A,B) = P(A) \cdot P(B)$
 - All of these still apply if all conditioned on D (i.e. $\sum_n P(A=a|D) = 1$)
- Always formally write out what prob. you're solving for

VPI / Decision Networks

observing noisy variables is less valuable

$$VPI(x) \triangleq MEUC(x) - MEUC(\emptyset)$$

$$VPI(Y|X) \triangleq MEUC(X,Y) - MEUC(X)$$

$$VPI(X,Y) = VPI(X) + VPI(Y|X)$$

① If have data what will I do

Value iteration \leftrightarrow variable elimination \Rightarrow computes full thing and then predicts \rightarrow too much time

Product Rule \Rightarrow doesn't have to break it down all the way

i.e. $P(A,B,C) = P(C|A)P(A)P(B)$

Inference by enumeration

- select entries consistent w/ evidence
- sum out H to get joint of query & evidence
- normalize

Product Rule $P(y)P(x|y) = P(x,y)$

Chain Rule $P(x_1, x_2, x_3, \dots, x_n) = \prod P(x_i | x_1, \dots, x_{i-1})$

(i.e. $P(x_1)P(x_2|x_1)P(x_3|x_1, x_2) \dots$)

Sampling = inference approx p then show converges in limit (consistent)

Independence $P(x,y) = P(x)P(y)$

$P(x|y) = P(x)$

Bayes Rule

$P(x,y) = P(x|y)P(y) = P(y|x)P(x)$

Conditional Independence

$P(x,y|z) = P(x|z)P(y|z)$
 $P(x|y,z) = P(x|z)$

Variable Elimination

- if variable appears in front of conditioning bar \rightarrow in resulting factor; else behind; can't appear in front for more than one factor
- join & eliminate till you get query
- if evidence, start w/ factors that select the evidence
- result is selected joint of query & evidence
- normalize \rightarrow by what you want to condition by
- Local CPTs instantiated w/ evidence
- join/elim hidden variables
- polytree always has efficient ordering \rightarrow cutset to polytree

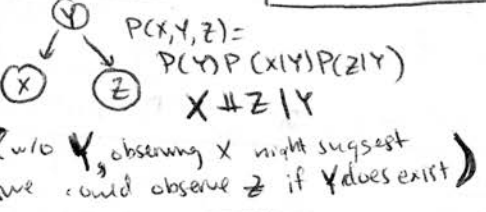
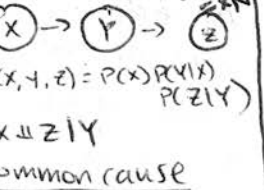
Error Sampling (~ direct evaluation)

Start from top & move down, sampling at each node; total/normalize

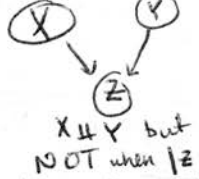
Rejection Sampling - if query doesn't appear, toss out

Bayes' Net # Not every BN can represent every joint distribution

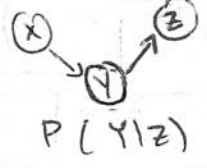
- arrows encode conditional independence (NOT necessarily causal structure; could be correlation)
- joint distribution over N variables 2^N
- N -node net \rightarrow save space $O(N \cdot 2^{k+1})$



Common effect



CPT from midterm

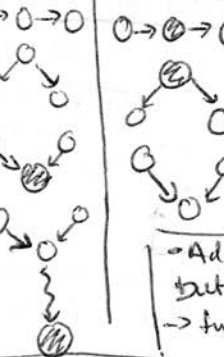


Probability tables
 $P(x)$
 $P(y|x)$
 $P(z|y)$
 $P(y|z) \Rightarrow x, y, z$

choose CPTs based on how you marginalize

Probability of full assignment = product of local conditional distributions

Active / Inactive



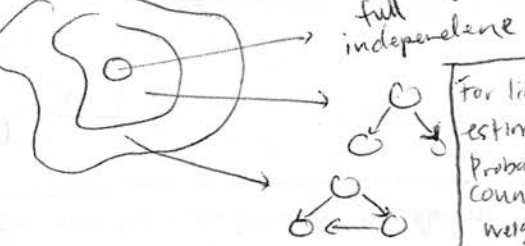
No active paths btw x & y (considering all undirected paths) \Rightarrow independence

gives precise indep info so can build joint dist using this

Graph structure encodes certain (conditional) independencies

Adding arcs increases set of distributions but has costs (can represent more)

full conditioning can encode any distribution



For likelihood, estimate probability by counting weights

Factor Zoo

- # of capitals = dimensionality of table
- Joint: $P(x,y) \rightarrow$ sum to 1
- Selected joint: $P(x,y) \rightarrow$ fixed x , all $y \rightarrow$ sums to $P(x)$
- Single conditional: $P(y|x) \rightarrow$ fixed x , all $y \rightarrow$ sums to 1
- Family of conditionals: $\rightarrow P(x|y)$ All x , all $y \rightarrow$ sums to $|Y|$ (i.e. $P(x|y) \rightarrow 1, |Y|$ times)
- Specified family $P(y|x) \rightarrow$ fixed y , all $x \rightarrow$ sums to unknown

ISVM \Rightarrow max margin $\min_w \frac{1}{2} ||w||^2$
 $\forall_i, y_i, w_i^* \cdot f(x_i) \geq w_i^* y_i + 1$

w/ factors \Rightarrow # of entries = $2^{3 \times 20}$

Likelihood Weighting

downside: only reflects in downstream stuff

\rightarrow set evidence & sample rest (weight by probability of evidence)

$$SWS(Z, e) \cdot W(Z, e) = \prod_{i=1}^m P(z_i | Parents(Z_i)) \prod_{i=1}^m P(e_i | Parents(e_i))$$

won't know some independencies unless inspect specific distribution (not from BN alone)

Reflex Agent

→ preplanned behavior (acts on it)

DFS

- stack
- "leftmost" solution
- $O(b^m)$ time

fringe = $O(b \ln)$

More Csp's...?
Local search
→ improve single var until can't be better
(local) → not complete (faster more or optimal mem efficient) → no change

BFS

- Queue
- optimal if all costs = 1
- shallowest level (for solution)

$O(b^d)$ search time

$O(b^d)$ fringe (opens everything)

Iterative Deepening

DFS memory but BFS time

DFS depth 1 2 3

(ov: computationally

Uniform Search Cost

- priority queue
- optimal

explore in every direction

cheapest solution = C^*

ans @ least ϵ ... effective depth

time & depth $O(b^{C/\epsilon})$

Greedy

- looks only at heuristic (state that looks @ the closest state)
- badly-guided DFS (lol)

$h(n) = \text{forward cost}$

A* (A^* optimal → admissible (tree))

- $f(n) = h(n) + g(n)$
- stop when dequeue goal
- expands toward the goal (not equally in all directions)

A^* optimal (consistent graph search)

CONSISTENT

$h(A) - h(B) \leq \text{cost}(A \rightarrow B)$

$h(A) \leq \text{cost}(A \rightarrow C) + h(C)$

f never decreases along path

Search state

* Based on what info is available

→ details necessary for planning

example = positions (x,y) XY etc. (counting)

Search tree properties

- b is the branching factor
- m = max depth
- $O(b^m)$ = num nodes in entire tree

Constraint satisfaction problems ($O(d^m)$ domain)

- unary constraint - explicit: variable \in values?
- binary constraint - implicit: $VA \neq WB$
- backtracking (DFS w/ 2 improvements)
 - 1 variable @ a time
 - 2 Check constraints as you go (consider values w/ no conflict)

speed up: 1 filtering: - forward checking (cross off values after current assignment, based on constraints) * X lose value, neighbors rechecked

(a.k.a. arc consistency) * delete from tail

X → Y consistent iff every x in tail, some y in head

modify forward checking w/ constraint propagation

queue arcs ⇒ modify then enqueue neighbors

$O(n^2 d^3) \rightarrow O(n^2 d^2)$ [can be reduced]

Ordering CSP ⇒ BREAK CYCLE

- 1 One solution
- 2 Multiple solutions
- 3 No solutions but doesn't know it yet

Minimum Remaining Values

- 1 choose variable w/ fewest remaining values in domain

K-consistency $K=2$ arc consistency (binary CSP)

Tree structured CSPs $O(nd^2)$ (no loops)

Pick root, order $a \rightarrow a \rightarrow a \rightarrow a$

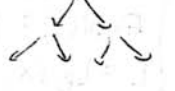
- 1 Remove backward
- 2 Assign forward (if no sol, arc c was used)

Nealy Tree

assign all possible values for cutset the tree the rest

$O(d^c (n-c) d^2)$ * cutset back [run-time] tracking ⇒ exponential

Search tree



Search graph = condensed search tree

func Tree-search (prob, strategy)

init search tree loop do

if there are no candidates for expansion ⇒ failure

choose a leaf node for expansion

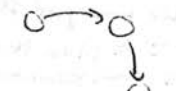
if node contains a goal state return solution

else expand the node & add the resulting nodes to the search tree

end

→ add "visited" closed set for graph search (to avoid duplicate visitations)

Search Graph



Search graph = condensed search tree

func Tree-search (prob, strategy)

init search tree loop do

if there are no candidates for expansion ⇒ failure

choose a leaf node for expansion

if node contains a goal state return solution

else expand the node & add the resulting nodes to the search tree

end

→ add "visited" closed set for graph search (to avoid duplicate visitations)

Heuristics

Admissible

$0 \leq h(n) \leq h^*(n)$

"optimistic"

→ relaxed problems

dominance: $h_a(n) \geq h_b(n)$ if true $\forall n$

$h(n) = \max(h_a(n), h_b(n))$

→ closer = better directed search

Admissible \leq actual optimal cost current node to goal

- for graph search need consistent consistency ⇒ admissibility

* Tree must be linear

Iterative algorithm for CSPs

live on the edge

while not solved: randomly select conflicted variable

variable selection = min-conflicts heuristic

evolute fewest constraints

Performance: CPU time

convert to tree if possible for analysis

Least constraining value

- 1 One that rules out fewest values in remaining variables

Iterative "local"

- no solution = run forever (t)

Genetic Algorithms

→ Best N hypotheses (fitness test)

→ pairwise crossover opera

Entropy (Information Encoding)

- code of length $\log(1/p)$
- more unlikely \Rightarrow more bits
- more likely \Rightarrow less bits

$$E_p(\log_2 1/p_i) = \sum_{i=1}^n -p_i \log_2 p_i$$

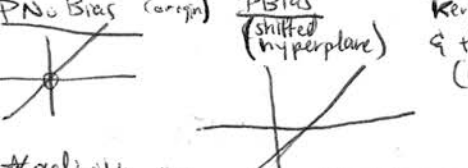
- more uniform \Rightarrow higher entropy
- more values \Rightarrow higher entropy
- more peaked \Rightarrow lower entropy
- rare values almost "don't count"

continuous \rightarrow discrete variables

- pick a point a to split on to give highest information gain (effective "discrete")

* if got to data missing in training (for particular attribute), then pick "most likely" \Rightarrow branch w/ node w/ greatest # of items

Classification + Kernels



* activity on one side of decision tree shouldn't reflect on other side b/c they already diverged

Bag of Words Model

feature for each word in the document; value is the word @ that position in document (library ALL words)

same for each feature; where word appears is irrelevant

- A - B - C - D - E

This is not a tree can't run an consistency; but if cutset & assign a , then can apply an consistency

Handout Question w/ MRV & LCV

LCV doesn't affect nodes expand b/c it just orders values (if we got to go through all)

MRV does affect b/c increase computation b/c of incompatible values remaining (used) w/ restricted variables

Information Gain

Gain = entropy before split - entropy after split

use * expected entropy

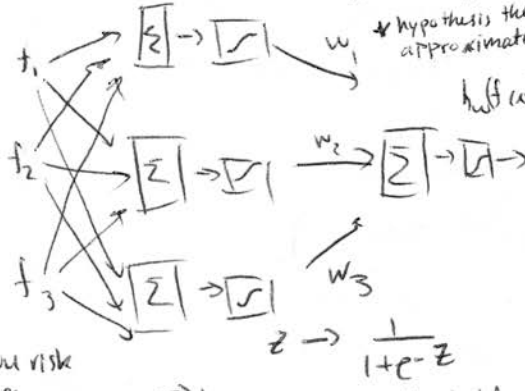
Pruning

- Build tree
- Start from bottom & prune according to $P_{chance} > \text{Max } P_{chance}$
- delete this one
- max P_{chance} = regularization param
- some chance nodes not pruned if redeemed later
- pruning intended to improve hold out accuracy

held-out data

Neural Networks

- universal approximators
- learning features



* to many neurons & you risk overfitting

* each weight is a learned parameter

\rightarrow by not having threshold be strict edge then won't get stuck when hill climbing

Quadratic

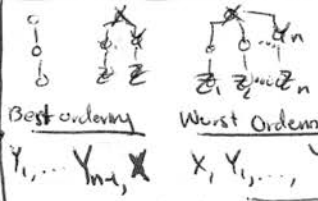
kernel maps to $[1, x_1, x_2, x_1^2, x_2^2]$

& then projects back (linear in quadratic space)

circles, ellipsoids & parabolas

* just b/c the features exist in weight, doesn't mean they need a value

* dual perception works best if more dimensionality of feature space than training data / points



Query: $P(Y_n | z_1, \dots, z_n)$

exponential growth

\Rightarrow more training data always good

\Rightarrow remove edges on Bayes means reduce hypothesis space

* admissible just needs to be less than true cost; underestimate

Game trees

$\Rightarrow \beta$ resets itself for the new branch & only looks @ α & β relevant to itself

RL learns about agent's own optimized behavior

if want to learn more about an adversary's behavior (which you know follows a certain rule) then state space search

EMV = expected monetary value (expectation)

equivalent monetary value "how much I do I need to get offered another are equivalent"

how much would I pay questions

insurance = equivalent value

b/c paying more means lottery is not worth it

If there is uncertainty in Pacman's position (non-deterministic behavior) \downarrow model w/ MDP

hold out is the norm

training error improves so does hold out, but eventually error flattens, then classifier probably overfits

How to determine CPTs

- Independence
- what do I need to marginalize

If assign in linear order where node assigned before all children, then one are consistency will prevent backtracking

If assign cutsets \Rightarrow so it's not cyclic anymore

first then repeatedly apply are consistency won't need to backtrack

so it's not cyclic anymore

Tabular Q-learning & full-depth minimax both compute exact value of all states

Naive Bayes is overconfident \Rightarrow lower entropy b/c variation in prob decreases

Feature based Q learning is approximation of the states (won't match minimax)

Non linearly separable data

1) Case-based reasoning

Classification (labeled data)

Parametric models

Nearest neighbor (sim. function)

- 1-NN \Rightarrow most similar data pt
- K-NN \Rightarrow pick K nearest neighbors & vote (need weighting scheme)

Dual Perceptron

$$w_j = 0 + f(x_j) - f(x_i) + \dots$$

$$w_j = \sum_i a_{i,j} f(x_i)$$

$$\alpha = \langle \alpha_{1,j}, \alpha_{2,j}, \dots \rangle \text{ dual representation}$$

$$\text{score} = w_j \cdot f(x)$$

$$= \left(\sum_i a_{i,j} f(x_i) \right) \cdot f(x)$$

$$= \sum_i a_{i,j} (f(x_i) \cdot f(x))$$

$$= \sum_i a_{i,j} K(x_i, x)$$

perceptron is explicit learning b/c train parameters

- 1) Fixed set of parameters
- 2) more data \Rightarrow better settings (tune the parameters) (i.e. perceptron)

Non-parametric

- complexity of classifier increases w/ data
- better in limit, worse in non-limit
- when similarity is needed
- Sim function: 1) dot product
- 2) invariant matrices (i.e. rotation)
- 3) data augmentation
- what you use to determine nearest neighbor

Infinite Dimension $K(x, x') = \exp(-\|x - x'\|)$

Dual Perceptron

how to do similarity function perceptron

one alpha vector per label; one alpha entry for each training data
can learn to ignore noisy data
allows to do nearest neighbor w/ perceptron \Rightarrow allowing learning to happen (tune α 's \Rightarrow parameters)

Kernels

\Rightarrow non-linear separators
e.g. quadratic kernel $K(x, x') = (x \cdot x' + 1)^2$
 $= \sum_{i,j} x_i x_j x'_i x'_j + 2 \sum_i x_i x'_i + 1$
so as not to represent according to vector full size of features

Kernel could be black box; don't even need to know weight vectors or feature vectors
Kernel = similarity
diff types of Kernel functions

Invariant Classification

- 1) Start w/ Error
- 2) $y = \text{argmax}_y \sum_i \alpha_{i,y} K(x_i, x)$
- $\alpha_{y,n} = \alpha_{y,n} - 1$
- $\alpha_{y^*n} = \alpha_{y^*n} + 1$

Perceptron: use hid out error to determine when to stop

Agglomerative Clustering

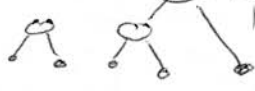
produce dendrogram

Clustering

no labels yet
un supervised learning
K-Means (class used euclidean dist)

- 1) initially, each instance its own cluster
- Repeat:
 - 2) two closest clusters
 - 3) Merge into new cluster
 - 4) stop when one cluster left

can pick which level of dendrogram to look @



It's ok to think of same looking points as different as long as classify

- "Closest":
- 1) closest pair
 - 2) farthest pair
 - 3) average of all pairs
 - 4) Ward's method (min variance, like k-means)

- 1) pick K arbitrary centers
- 2) assign data to nearest center
- 3) average & reassign center
- 4) repeat till assignments don't change

minimizes dist b/w diff pts by reassignment w/ Manhattan dist \Rightarrow median

Decision Trees

compact representation of probability table
automatically reduces size of hypothesis space
allows to be more expressive than typical perceptron \Rightarrow can join features

SIZE of CPT

Binary classification w/ 3 variables (binary)
 $2^3 \Rightarrow$ variable combos. each of these 8 have 2 labels ($2^2 = 2^8$)

bias \Rightarrow keep getting same consistent error in estimation (too "biased")
variance \rightarrow overfitting (too much hypothesis)
reduce bias \rightarrow simplicity (reduce hypothesis space)
regularization (smoothing) can keep hypothesis space large but can't get to extremes
variance \leftrightarrow consistency
consistency & simplicity \rightarrow no overfitting

have a g we need to approx from $(x, g(x))$
hypothesis space, H (defined by training data)
 \rightarrow classification? uses hypothesis space
 \rightarrow regression

Decision Tree Size

$4n: n = \#$ of attributes
i.e. w/ 1 attribute
 $\{y_0, y_1, y_0, y_1\}$ 4 diff possibilities

$\phi(\{x_i, z, \{a_i, z, \{c_i, z\}\}) = \sum \text{dist}(x_i, c_i)$
points assignments means
this is a minimization problem so it will converge
initialization makes a difference (variance-based, split/merge, etc.)

reassignment minimize according to metric used

Search Problem

start state \rightarrow goal state
 resolution = a plan that does this transformation

Game algorithms

strategy (policy) for move to move

terminal state: value known
 value: best utility from that state

minimax

\rightarrow minimal value: best achievable utility against a rational (optimal) adversary

like DFS: Time $O(b^m)$
 Space $O(bm)$

def value (state):

- if terminal, return utility
- if next agent max: return max_val
- if next agent min: return min_val

α = MAX's best option on path to root
 β = MIN's best option on path to root

def max-value (state, α, β):

init $v = -\infty$
 each successor
 $v = \max(v, \text{val}(\text{succ}, \alpha, \beta))$
 if $v \geq \beta$ return v
 $\alpha = \max(\alpha, v)$
 return v

alpha-Beta: no effect on minimax value for root

expectimax (coverage case) \leftarrow replace min-value

* careful w/ pruning for expectimax

* depending on adversary, expectimax can be optimistic

* minimax = pessimistic

* minimax, terminal function scale doesn't matter

insensitivity to monotonic transformations

* expectimax \Rightarrow magnitudes must be meaningful

Preference $A \succ B$

Indifference $A \sim B$

Axioms of Rationality

1. orderability
 $(A \succ B) \vee (B \succ A) \vee (A \sim B)$
2. Transitivity
 $(A \succ B) \wedge (B \succ C) \Rightarrow (A \succ C)$
3. continuity
 $A \succ B \succ C \Rightarrow \exists p [p, A; 1-p, C] \sim B$

$U(A) \geq U(B) \Leftrightarrow A \geq B$
 $U(L) = \sum_{x \in L} P(x) U(x)$

4. substitutability
 $A \sim B \Rightarrow [p, A; 1-p, C] \sim [p, B; 1-p, C]$

5. Monotonicity
 $A \succ B \Rightarrow [p, A; 1-p, B] \succ [q, A; 1-q, B]$

Money

if pay in exchange for

$EMV(L) = p \cdot X + (1-p) \cdot Y$

$L = [p, X; (1-p), Y]$

$U(L) = p \cdot U(X) + (1-p) \cdot U(Y)$

$U(L) < U(EMV(L))$

* γ helps algorithm converge

Bellman equations (expected utility when act optimally)

$V^*(s) = \max_a Q^*(s, a)$

$Q^*(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$

$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^*(s')]$

Value Iteration: * Repeat till converge

$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s)]$

* complexity of iteration: $O(s^2)$

* policy may converge before optimal values

consider what reward is motivating

MDP

* one way to solve is expectimax

* depends only on current state

* more, sooner

* γ helps algorithm converge

* Alpha-beta

* remember what each references

* all updates based on full iter

* efficiency $O(s^2)$ per iter

* Policy Evaluation

* Policy Extraction

* Policy Iteration

one more iteration, max(Q)

1. Policy evaluation

2. Policy improvement \Rightarrow use one-step look-ahead w/ resulting utilities

* converge much faster under some conditions

learning \rightarrow model-based

\rightarrow MDP model

- 1. Passive (observe)
- 2. Active

1. Passive (observe) \rightarrow no choice in policy (given) \Rightarrow effectively evaluation \Rightarrow find value

(A) Direct Evaluation

1. Act according to π

2. every time visit a state, write what discounted rewards turned out to be

3. average (equal weight) $\frac{1}{N} \sum a_i$

2. Active: get to choose actions; find policy

4. Q learning \Rightarrow will converge; off-policy learning

$Q_{k+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$

sample = $R(s, a, s') + \gamma \max_{a'} Q(s', a')$

$Q(s, a) \leftarrow (1-\alpha) Q(s, a) + \alpha(\text{sample})$

* explore enough

* eventually make α small enough; not too quickly else conv. to local

Exploration vs. Exploitation

E-greedy

$P = \epsilon \Rightarrow$ random

$1 - \epsilon \Rightarrow$ current policy

* to prevent thrashing overtime lower ϵ over time

or exploration func

eventually stop exploring states $[\delta f(Q(s, a'), w(s, a'))]$

$n = \#$ of times visited \rightarrow bonus to unknown

Policy = action for each state

Utility = sum of (discounted) rewards

Time-limited

$V_k(s)$ = optimal value k more time steps

(probabilities included)

* Repeat till converge

* complexity of iteration: $O(s^2)$

* policy may converge before optimal values

1. is faster than value iteration

2. is same pace as value iteration

\Rightarrow Solve MDP

\rightarrow decreasing learning rate can give converging values

Episode = start to finish

Iter = each step in episode

learn policies that maximize rewards, in values that predict them

Approximate Q

$w_i \leftarrow w_i + \alpha(\text{diff})f_i(s, s)$ - generalize states

weighted combo of feat = value

- con. share feat but diff on value

Naive Bayes Nets



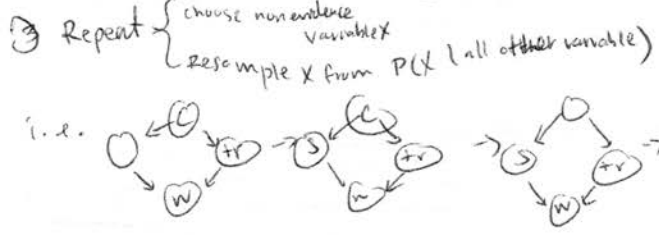
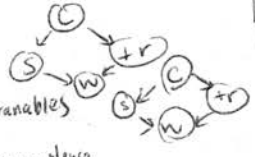
→ features are conditionally independent given label
 → info about f_1 doesn't give info about f_2 [i.e. aligned along rectangular axis]

Training (create for test) → small data (training) → less error
 → also more overfit

Gibbs Sampling

* sum of weights = how "effective" samples were obtained (cannot high)
 * start w/ arbitrary instantiation consistent w/ evidence; sample one variable @ a time, conditioned on all the rest, but keep evidence fixed. Repeat

① Fix evidence
 ② initialize other variables
 ③ Repeat { choose non-evidence variable x resample x from $P(x | \text{all other variables})$ }



$P(S|C, Tr, -W) = \frac{P(S|C) P(-W|S, Tr)}{\sum_s P(S|C) P(-W|s, Tr)}$

→ only CPTs that have resampled variable need to be considered, & joined together

Decision networks

$P(X_{1:T} | e_{1:T})$
 $\propto P(X_{1:T}, e_{1:T})$ (maximize path)

MEU = max EU(a)
 EU(leave|bad) = $\sum_w P(w|bad) U(leave, w)$
 MEU (F=bad) = max EU(a|bad)

VPI ($E_i | e$) = $(\sum_{e'} P(e' | e) MEU(e, e')) - MEU(e)$
 → for every evidence, insert a step like this

VPI properties:
 ① Non-negative
 ② Non-additive
 VPI ($E_j, E_k | e$) \neq VPI ($E_j | e$) + VPI ($E_k | e$)
 ③ Order - Independent
 VPI ($E_j, E_k | e$) = VPI ($E_j | e$) + VPI ($E_k | e, E_j$)
 = VPI ($E_k | e$) + VPI ($E_j | e, E_k$)

Model structure (Laplace Smoothing)

→ smoothing decreases maximum likelihood expectation
 → if model represents a larger family of models than the MLE for that model is greater

Particle Filtering (Likelihood Sampling)

→ calculating probabilities = normalizing the weights (counting weighted samples)
 → only count samples w/ the correct conditioning matching the query

① Elapse time
 $X_t = \text{sample}(P(X_t | x_{1:t}))$

② Observe
 $W(X_t) = P(e_t | x_t)$

③ Resample

State Trellis
 - graph of states & transitions
 - arc weight = $P(x_t | x_{t-1}) P(e_t | x_t)$
 - each transition $x_{t-1} \rightarrow x_t = \text{arc}$
 - product of weights = sequence's probability

$m_t[x_{1:t}] = \max_{x_{1:t-1}} P(x_{1:t-1}, x_{1:t} | e_{1:t})$

Dynamic Bayes Net
 * create super nodes
 * "unroll" for T time steps, eliminate variables, then eliminate till $P(X_{1:T} | e_{1:T})$

Naive Bayes

$P(Y, F_1, \dots, F_n) = P(Y) \prod P(F_i | Y)$
 $|Y| \times |F|^n$ values
 $n \times |F| \times |Y|$ parameters (for naive)

Viterbi Algorithm
 $P(X_{1:T} | e_{1:T}) \propto P(X_{1:T}, e_{1:T})$ (maximize path)

Inference

① get joint probabilities of label & evidence for each label
 ② sum to get probabilities of evidence
 ③ normalize

$\text{argmax } P(\theta | x) = \text{argmax } P(x | \theta) P(\theta)$

MIRA

weight vector = hyperplane
Binary perceptron
 $w = w + y^* \cdot f$
Multiclass
 $w_y = w_y - f(x)$
 $w_{y^*} = w_{y^*} + f(x)$

$\min_w \frac{1}{2} \sum_y \|w_y - w_{y^*}\|^2$
 $w_y = w_{y^*} - \sum f(x)$
 $w_{y^*} = w_{y^*} + \sum f(x)$

max step size
 $\epsilon = \min \left(\frac{w_{y^*} - w_y}{25 \cdot f} \right)$
 always changes for @ least 1 improvement

Time elapse

Observation
 $P(e_{t+1} | x_{t+1}) P(x_{t+1} | e_{1:t}) = P(x_{t+1} | e_{1:t+1})$

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

Forward Algorithm

Inferences
 $P(x_1, e_1) \propto P(x_1) P(e_1)$
 $P(x_2) = \sum_{x_1} P(x_1) P(x_2 | x_1) P(e_2)$ (base cases)

Stationary Dist
 $P_{\infty}(x) = P_{\text{start}}(x) = \sum_x P(x | x) P_{\infty}(x)$

when learned select distribution that maximizes likelihood

→ how much you pay for info = prob of not ideal • payoff for continuing knowing not ideal