15-780: Graduate AI Lecture 15. Uncertainty, Inference

Geoff Gordon (this lecture) Tuomas Sandholm TAs Sam Ganzfried, Byron Boots

Review

Spatial planning

- C-space
- Ways of splitting up C-space
 - Visibility graph
 - Voronoi
 - Cell decomposition
 - Variable resolution or adaptive cells (quadtree, parti-game)

RRTs

- Build tree by randomly picking a point, extending tree toward it
- Optionally:
 - build forward and backward at once
 - cross-link within tree
- Plan within tree to get close enough to goal

RRTs

- Tend to break up large Voronoi regions
- So, RRT search is coarse to fine
- First path found usually suboptimal; if we continue to grow tree and cross-link, get optimality in limit

Probability

- Random variables
- Events (atomic, composite, AND/OR/NOT)
- Distributions: joint, marginal
- Law of total probability
 - $P(X) = P(X, Y=y_1) + P(X, Y=y_2) + \dots$

Probability

Conditional: incorporating observations

15-780

		\boldsymbol{A}	A-	B+	
V4	A	.21	.17	.07	
HW4	<i>A-</i>	.17	.15	.06	
	B+	.07	.06	.04	

)	$oxed{A}$.41
5-780	A-	.35
I	B+	.24

$$P(15-780=A \mid HW4=B+) = .04 / (.07+.06+.04)$$

Conditioning

- In general, divide a row or column by sum
 - P(X | Y=y) = P(X, Y=y) / P(Y=y)
- \circ P(Y=y) is a marginal probability
- \circ $P(X \mid Y=y)$ is a row or column of table
- Thought experiment: what happens if we condition on an event of zero probability?

Notation

- \circ $P(X \mid Y)$ is a function: $x, y \to P(X=x \mid Y=y)$
- As is standard, expressions are evaluated separately for each realization:
 - $P(X \mid Y) P(Y)$ means the function $x, y \rightarrow P(X=x \mid Y=y) P(Y=y)$

Independence

- X and Y are **independent** if, for all possible values of y, $P(X) = P(X \mid Y=y)$
 - equivalently, for all possible values of x, P(Y) = P(Y | X=x)
 - \circ equivalently, P(X, Y) = P(X) P(Y)
- Knowing X or Y gives us no information about the other

Independence: probability = product of marginals

AAPL price

_					
\mathcal{L}		ир	same	down	
Weather	sun	.09	.15	.06	0.3
W	rain	.21	.35	.14	0.7
'					

0.3 0.5 0.4

Bayes Rule

- *For any X, Y, C*
 - $\circ P(X \mid Y, C) P(Y \mid C) = P(Y \mid X, C) P(X \mid C)$
- Simple version (without context)
 - $\circ P(X \mid Y) P(Y) = P(Y \mid X) P(X)$
- \circ Proof: both sides are just P(X, Y)
 - P(X | Y) = P(X, Y) / P(Y) (by def'n of conditioning)

Continuous distributions

- What if X is real-valued?
 - Atomic event: $(X \le x)$, $(X \ge x)$
 - any* other subset via AND, OR, NOT
 - X=x has zero* probability
- \circ P(X=x, Y=y) means
 - $\circ \lim P(x \le X \le x + \varepsilon, Y = y) / \varepsilon \qquad \varepsilon \to 0 + \varepsilon$

Factor Graphs

Factor graphs

- Strict generalization of SAT to include uncertainty
- Factor graph = (variables, constraints)
 - variables: set of discrete r.v.s
 - constraints: set of (possibly) soft or probabilistic constraints called factors or potentials

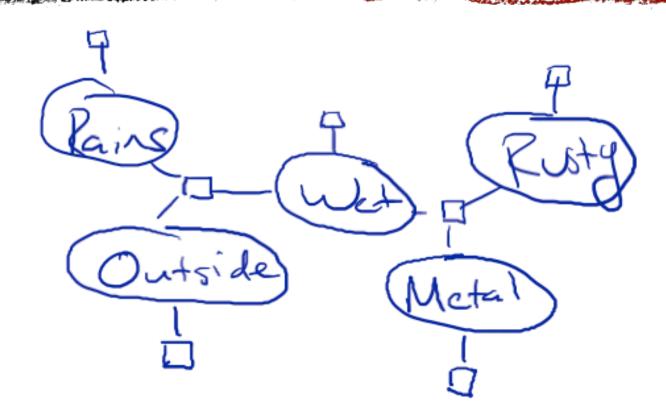
Factors

- *Hard constraint:* $X + Y \ge 3$
- ∘ Soft constraint: $X + Y \ge 3$ is more probable than X + Y < 3, all else equal
- Domain of factor: set of relevant variables
 - \circ {X, Y} in this case

Factors

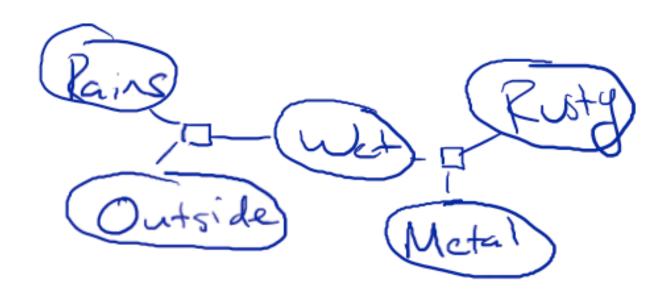
	Hard				Se	oft			
_			X					X	
		0	1	2			0	1	2
	0	0	0	0		0	1	1	1
Y	1	0	0	1	Y	1	1	1	3
	2	0	1	1		2	1	3	3

Factor graphs



 Variables and factors connected according to domains

Factor graphs



 Omit single-argument factors from graph to reduce clutter

Factor graphs: probability model

$$P(R_{a}, W, O, M, R_{b}) =$$

$$\phi(R_{a}, W, O, M, R_{b}) / Z$$

$$\phi(R_{a}, W, O, M, R_{b}) =$$

$$\phi(R_{a}, W, O, M, R_{b}) =$$

$$\phi(R_{a}, W, O) \phi_{z}(O, M, R_{b})$$

$$\phi_{z}(R_{a}) \phi_{y}(W) \phi_{z}(O)$$

$$\phi_{z}(R_{b}) \phi_{y}(W) \phi_{z}(O)$$

$$\phi_{b}(M) \phi_{z}(R_{b})$$

Normalizing constant

Also called partition function

Factors

Ra	0	W	$ \Phi_1 $
T	T	T	3
T	T	$oxed{F}$	1
$\mid T \mid$	$oldsymbol{F}$	T	3
T	F	$oxed{F}$	3
$oxed{F}$	T	T	3
$oxedsymbol{F}$	T	$oldsymbol{F}$	3
$oxed{F}$	$oxed{F}$	T	3
$oxed{F}$	$oldsymbol{F}$	$oldsymbol{F}$	3

$oxed{W}$	M	Ru	$ \Phi_2 $
$oxed{T}$	T	T	3
$oxed{T}$	T	$oxed{F}$	1
T	$oxed{F}$	T	3
$oxed{T}$	$oxed{F}$	$oxed{F}$	3
$oxed{F}$	T	T	3
$oxed{F}$	T	$oxed{F}$	3
$oxed{F}$	$oxed{F}$	T	3
$oxed{F}$	$oxed{F}$	$oxed{F}$	3

Unary factors

Ra	Φ_3
T	1
$oxed{F}$	2

W	$oxedsymbol{\Phi}_4$
T	1
$oxed{F}$	1

O	$ \Phi_5 $
T	5
$oldsymbol{F}$	2

M	$ \Phi_6 $
T	10
$oxedsymbol{F}$	1

Ru	$ \Phi_7 $
T	1
$oxedsymbol{F}$	3

Inference Qs

- Is Z > 0?
- \circ What is P(E)?
- What is $P(E_1 \mid E_2)$?
- Sample a random configuration according to P(.) or P(. | E)
- \circ Hard part: taking sums of Φ (such as Z)

Example

 \circ What is P(T, T, T, T, T)?

Example

- What is $P(Rusty=T \mid Rains=T)$?
- \circ This is P(Rusty=T, Rains=T) / P(Rains=T)
- \circ P(Rains) is a marginal of P(...)
- So is P(Rusty, Rains)
- Note: Z cancels, but still have to sum lots of entries to get each marginal

Relationship to SAT, ILP

- Easy to write a clause or a linear constraint as a factor: 1 if satisfied, 0 o/w
- Feasibility problem: is Z > 0?
 - more generally, count satisfying assignments (determine Z)
 - NP or #P complete (respectively)
- Sampling problem: return a satisfying assignment uniformly at random

Inference

Notation

- Boldface = vector of r.v.s or values
 - $X = (X_1, X_2, X_3, X_4, X_5)$
 - $\bullet \ x = (x_1, x_2, x_3, x_4, x_5)$
- \circ *Set index* = *subvector*
 - If $D = \{1, 3, 4\}$ then $X_D = (X_1, X_3, X_4)$
- \circ In particular, $X_{D(\Phi)} = input \ to \ factor \ \Phi$

Finding x w/ P(X=x)>0

 $\circ \{x \mid P(X=x)>0\} = support \ of \ P(X)$

- Replace each factor Φ with $\Phi' = I(\Phi > 0)$
- Now we have a CSP (or SAT); do DPLL

Counting support

- DPLL also works to count $\{x \mid P(X=x)>0\}$
- Or, we can get smarter: dynamic programming
- $Example: (A \lor B) \land (B \lor C) \land (C \lor D)$

DP: $(A \lor B) \land (B \lor C) \land (C \lor D)$

$oxed{A}$	B	Φ_1
T	T	1
T	F	1
$oxed{F}$	T	1
$oxed{F}$	$oxed{F}$	0

B	C	Φ_2
T		1
T	F	1
F		1
$oxed{F}$	$\int F$	0

$$egin{array}{c|cccc} C & D & oldsymbol{\Phi}_3 \\ \hline T & T & 1 \\ \hline T & F & 1 \\ \hline F & T & 1 \\ \hline F & F & 0 \\ \hline \end{array}$$

- Consider B=T and B=F separately
- \circ B=F: 1; B=T: 2*2 = 4
- Subtotal: |supp(ABC)| = 5; note C=F in 2

DP: $(A \lor B) \land (B \lor C) \land (C \lor D)$

$oxed{A}$	B	Φ_1
T	T	1
T	F	1
$oxed{F}$	T	1
$oxed{F}$	F	0

$oxed{B}$	C	Φ_2
$\mid T \mid$	T	1
T	F	1
$oxed{F}$	T	1
$oxed{F}$	$oxed{F}$	0

C	D	Φ_3
T	T	1
T	F	1
$oxed{F}$	T	1
$oxed{F}$	F	0

- Consider C=T and C=F separately
- \circ C=F: 2*1=2; C=T: 3*2=6
- Total: |support| = 8

Quiz

• What is size of support for $(A \lor B) \land (B \lor C) \land (C \lor D) \land (D \lor E)$

DP for finding Z

- Same trick works for finding Z
 - not surprising: Z is a sum over support

Ra	0	W	Φ_l
T	T	T	3
T	T	$oxed{F}$	1
T	F	T	3
T	F	F	3
$oxed{F}$	T	T	3
$oxed{F}$	T	$oxed{F}$	3
$oxed{F}$	F	T	3
$oxed{F}$	F	F	3

Ra	$ \Phi_3 $	0	Φ_5
T	1	T	5
F	2	F	2

$oxed{W}$	Φ_4
T	1
$oxed{F}$	1

Ra	O	W	$ \Phi_1 $	Φ_{β}	Φ_4	$ \Phi_5 $	*
T	T	T	3	1	1	5	15
$\mid T \mid$	T	F	1	1	1	5	5
T	F	T	3	1	1	2	6
$\mid T \mid$	F	F	3	1	1	2	6
$oxed{F}$	T	T	3	2	1	5	30
$oxed{F}$	T	F	3	2	1	5	30
$oxed{F}$	F	T	3	2	1	2	12
$oxed{F}$	F	F	3	2	1	2	12

Ra	Φ_3	0	Φ_5
T	1	T	5
$oxed{F}$	2	$oxed{F}$	2

$oxed{W}$	Φ_4
T	1
$oxed{F}$	1

Ra	0	W	*
$oxed{T}$	T	T	15
$oxed{T}$	T	$oldsymbol{F}$	5
$\mid T \mid$	$oldsymbol{F}$	T	6
$oxedsymbol{T}$	$oldsymbol{F}$	$oldsymbol{F}$	6
$oxed{F}$	T	T	30
$oxed{F}$	T	$oldsymbol{F}$	30
$oxed{F}$	F	T	12
$oxedsymbol{F}$	F	$oxed{F}$	12

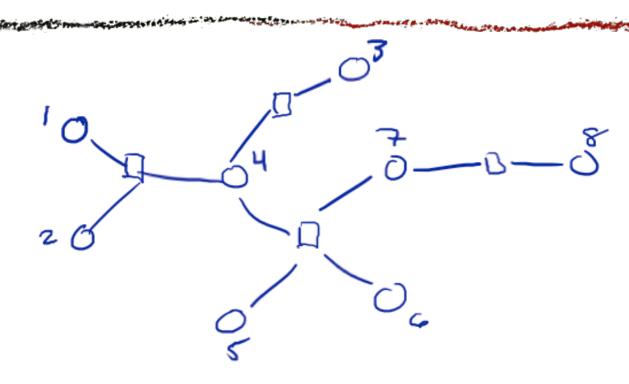
$oxed{W}$	λ
T	63
$oxed{F}$	53

APPROPRIE	-	4 for sens	e mining	Cheroderich	والمراجعة المارية		والمستعملية المتأثثة المتارية	red & parental	أحضرتك ويدين وااله	-	distribution of the same of th
$oxed{W}$	M	Ru	Φ_2	Φ_6	Φ_7	λ	*	M	$ \Phi_6 $	Ru	$ \Phi_7 $
T	T	T	3	10	1	63	1890	lacksquare	10	T	1
	T	$\mid F \mid$	1	10	3	63	1890		10		1
T	F		3	1	1	63	189	$\mid F \mid$	$\mid 1 \mid$	$\mid F \mid$	3
T	F	F	3	1	3	63	567				
$oxed{F}$	T	T	3	10	1	53	1590		$\mid W \mid$	$ \lambda $	
$oxed{F}$	T	F	3	10	3	53	4770				
$oxed{F}$	F	$\mid T \mid$	3	1	1	53	159			63	
$oxed{F}$	F	F	3	1	3	53	477		$\mid F \mid$	53	

$oxed{W}$	M	Ru	*
T	T	T	1890
T	T	$oxed{F}$	1890
T	F	T	189
T	$oldsymbol{F}$	$\mid F \mid$	567
$oxed{F}$	T	T	1590
$oxed{F}$	T	$oxed{F}$	4770
$oxed{F}$	F	T	159
$oxed{F}$	F	$oxed{F}$	477

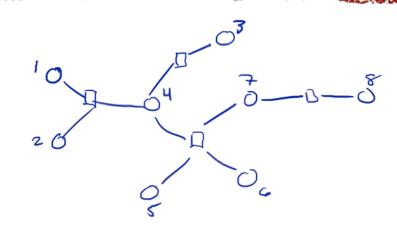
$$Z = 11532$$

Variable elimination



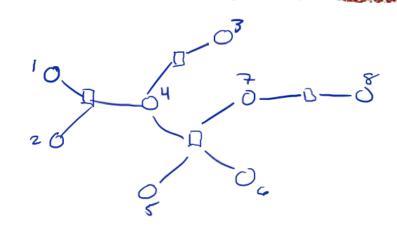
 Basic step: move a sum inward as far as possible, then do sum for each possible way of setting neighbors

Eliminating partway



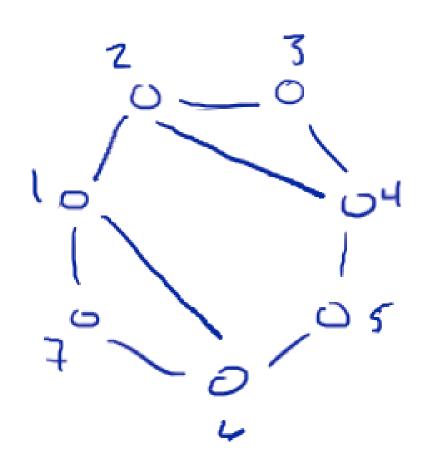
- For Z, want to sum over all X
- \circ For marginal $P(X_{45})$, eliminate X_{123678}
- Sum out 123, then 876, to get $ZP(X_{45})$
- Sum out 45 to get Z

Eliminating partway



- Conditional $P(X_{45} | X_6) = P(X_{456}) / P(X_6)$
- Sum out 123, then 87, to get $ZP(X_{456})$
- Sum out 45 to get $ZP(X_6)$
- Divide

A more difficult example



Treewidth

- Elimination order E: sum x_1 , then x_5 , ...
- treewidth(E) =
 (size of largest factor formed) 1
- \circ treewidth = min_E treewidth(E)
- Variable elimination uses space, time exponential in treewidth
- Worse: even computing treewidth is NPcomplete

Belief propagation

- Suppose we want all 1-variable marginals
- Could do N runs of variable elimination
- Or: the BP algorithm simulates N runs for the price of 2
- For details: Kschischang et al. reading

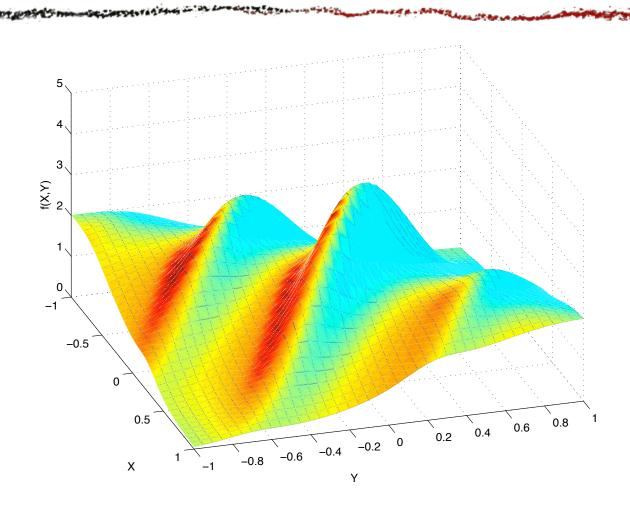
Approximate Inference

Most of the time...

- Treewidth is big
- Variables are high-arity or continuous
- Can't afford exact inference

- Partition function = numerical integration (and/or summation)
- We'll look at randomized algorithms

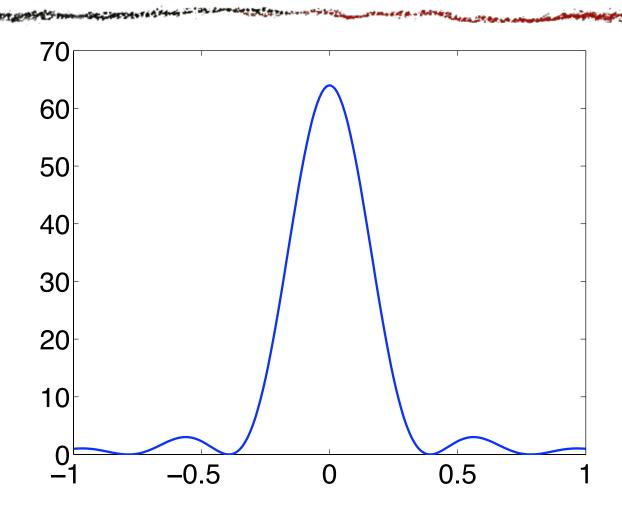
Numerical integration



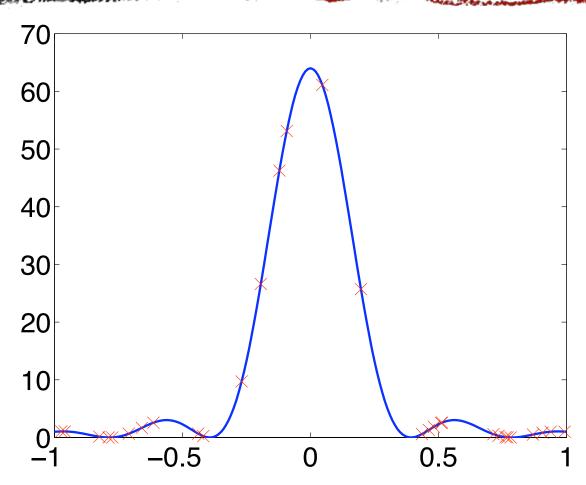
Integration in 1000s of dims



Simple 1D problem



Simplest randomized algorithm



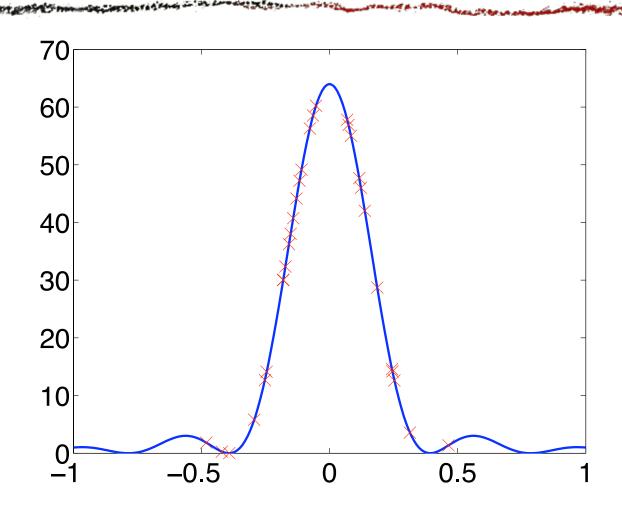
• *Uniform sampling:* $sum(f(x_i))/N$

Uniform sampling

$$E(f(X)) = \int P(x)f(x)dx$$
$$= \frac{1}{V} \int f(x)dx$$

- \circ So, VE(f(X)) is desired integral
- But standard deviation can be big
- Can reduce it by averaging many samples
- But only at rate 1/sqrt(N)

Nonuniform (importance) sampling



Importance sampling

- Instead of $x \sim uniform$, use $x \sim Q(x)$
- \circ Q = importance distribution
- Should have Q(x) large where f(x) is large
- Problem:

$$E_Q(f(X)) = \int Q(x)f(x)dx$$

Importance sampling

$$h(x) \equiv f(x)/Q(x)$$

$$E_Q(h(X)) = \int Q(x)h(x)dx$$

$$= \int Q(x)f(x)/Q(x)dx$$

$$= \int f(x)dx$$

Importance sampling

- So, take samples of h(X) instead of f(X)
- $w_i = 1/Q(x_i)$ is importance weight
- $\circ Q = uniform \ yields \ uniform \ sampling$

Variance

- How does this help us control variance?
- \circ Suppose f big ==> Q big
- \circ And Q small ==>f small
- Then h = f/Q never gets too big
- Variance of each sample is lower ==>
 need fewer samples
- A good Q makes a good IS

Importance sampling, part II

Suppose we want

$$\int f(x)dx = \int P(x)g(x)dx = E_P(g(X))$$

- Pick N samples x_i from proposal Q(X)
- Average w_i $g(x_i)$, where $w_i = P(x_i)/Q(x_i)$ is importance weight

$$E_Q(Wg(X)) = \int Q(x)[P(x)/Q(x)]g(x)dx = \int P(x)g(x)dx$$

Parallel importance sampling

Suppose we want

$$\int f(x)dx = \int P(x)g(x)dx = E_P(g(X))$$

 But P(x) is unnormalized (e.g., represented by a factor graph)—know only Z P(x)

Parallel IS

- Pick N samples x_i from proposal Q(X)
- If we knew $w_i = P(x_i)/Q(x_i)$, could do IS
- Instead, set $\hat{w}_i = ZP(x_i)/Q(x_i)$

Parallel IS

$$E(\hat{W}) = \int Q(x)(ZP(x)/Q(x))dx$$
$$= Z \int P(x)dx$$
$$= Z$$

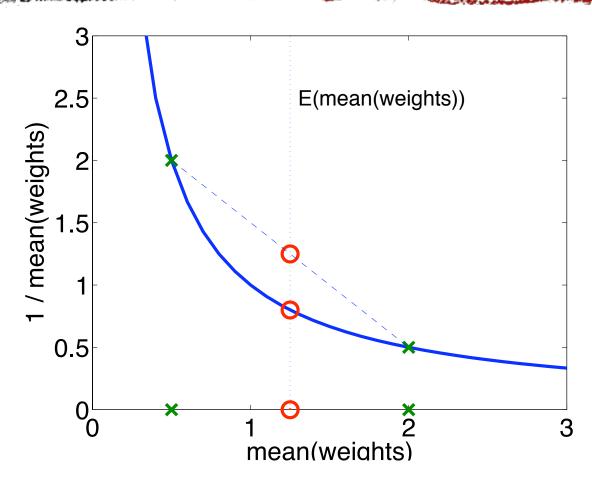
• So,
$$\bar{w} = \frac{1}{N} \sum_{i} \hat{w}_{i}$$
 is an unbiased estimate of Z

Parallel IS

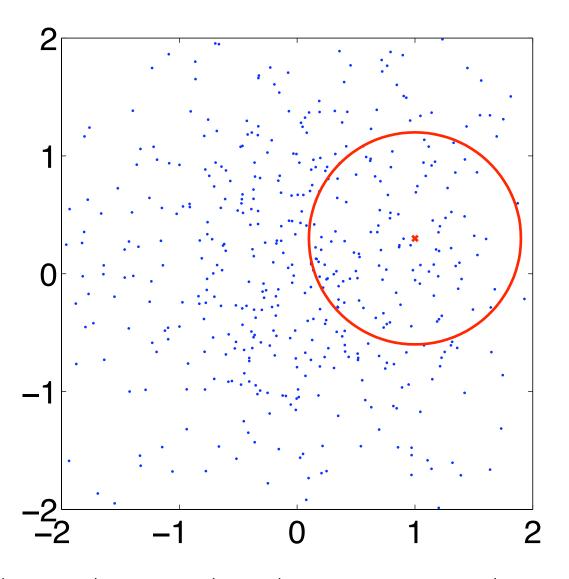
- So, \hat{w}_i/\bar{w} is an estimate of w_i , computed without knowing Z
- Final estimate:

$$\int f(x)dx \approx \frac{1}{n} \sum_{i} \frac{\hat{w}_{i}}{\bar{w}} g(x_{i})$$

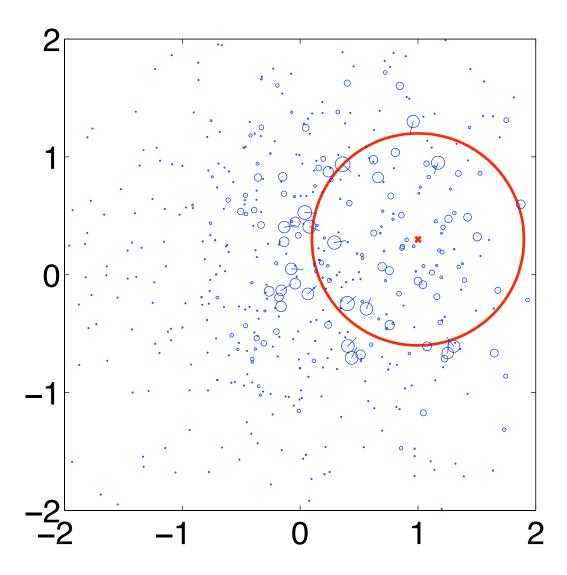
Parallel IS is biased



$$E(\overline{W}) = Z$$
, but $E(1/\overline{W}) \neq 1/Z$ in general



$$Q: (X, Y) \sim N(1, 1)$$
 $\theta \sim U(-\pi, \pi)$
 $f(x, y, \theta) = Q(x, y, \theta)P(o = 0.8 \mid x, y, \theta)/Z$



Posterior $E(X, Y, \theta) = (0.496, 0.350, 0.084)$

Back to high dimensions

- Picking a good importance distribution is hard in high-D
- Major contributions to integral can be hidden in small areas
 - \circ recall, want (f big ==> Q big)
- Would like to search for areas of high f(x)
- But searching could bias our estimates

MCMC

Markov-Chain Monte Carlo

- Design a randomized search procedure M which tends to increase f(x) if it is small
- Run M for a while, take resulting x as a sample
- Importance weight Q(x)?

Markov-Chain Monte Carlo

- Design a randomized search procedure M which tends to increase f(x) if it is small
- Run M for a while, take resulting x as a sample
- Importance weight Q(x)? Stationary distribution of M

Designing a search chain

$$\int f(x)dx = \int P(x)g(x)dx = E_P(g(x))$$

- Would like Q(x) = P(x)
- Turns out we can get this exactly, using Metropolis-Hastings

Metropolis-Hastings

- Way of designing chain w/Q(x) = P(x)
- Basic strategy: start from arbitrary x
- Repeatedly tweak x to get x'
- If $P(x') \ge P(x)$, move to x'
- \circ If P(x') << P(x), stay at x
- In intermediate cases, randomize

Proposal distribution

- Left open: what does "tweak" mean?
- Parameter of MH: Q(x' | x)
 - one-step proposal distribution
- Good proposals explore quickly, but remain in regions of high P(x)
- Optimal proposal?

MH algorithm

- Sample $x' \sim Q(x' \mid x)$
- Compute $p = \frac{P(x')}{P(x)} \frac{Q(x \mid x')}{Q(x' \mid x)}$
- With probability p, set x := x'
- Repeat for T steps (usually < T distinct samples)

MH notes

- Only need P(x) up to a constant factor
- Efficiency determined by:
 - how fast Q(x' | x) moves us around
 - how high acceptance probability p is
- Tension between fast Q and high p