## 15-780: Graduate AI Lecture 18. Learning

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

#### Admin

- HW4 questions?
- HW4 extension: was due Tue, now Thu
- Project interim reports
  - reminder: due 4/14
- Midterm 4/9
  - prep sessions TBA

## Review

## Probability

- Expectation
- Conditional expectation
  - law of iterated expectations
- Sample vs. population quantities
- Estimators; consistency

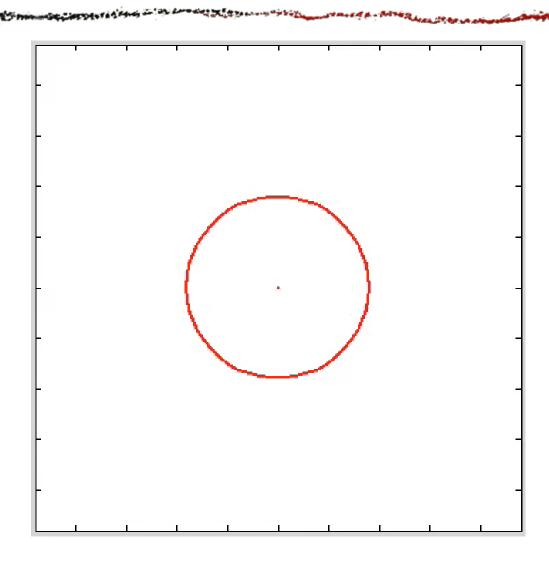
#### Markov-Chain Monte Carlo

• For computing:

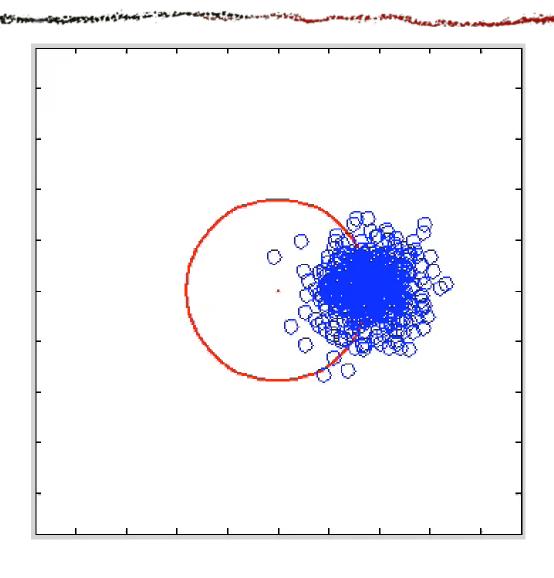
$$E_P(g(X)) = \int g(x)P(x)dx = \int f(x)dx$$

- Chief difficulty: finding a good importance distribution Q(x)
- Metropolis-Hastings: optimal distribution (Q = P) using randomized search

### Markov chain



## Stationary distribution



## Stationary distribution

$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \quad \Rightarrow \quad Q(\mathbf{x}_{t+1}) = \mathbb{P}(\mathbf{x}_{t+1})$$

$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \quad \Rightarrow \quad Q(\mathbf{x}_{t+1}) = \int \mathbb{P}(\mathbf{x}_{t+1}, \mathbf{x}_{t}) d\mathbf{x}_{t}$$

$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \quad \Rightarrow \quad Q(\mathbf{x}_{t+1}) = \int \mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_{t}) \mathbb{P}(\mathbf{x}_{t}) d\mathbf{x}_{t}$$

$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \quad \Rightarrow \quad Q(\mathbf{x}_{t+1}) = \int \mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_{t}) Q(\mathbf{x}_{t}) d\mathbf{x}_{t}$$

## Stationary distribution

$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \implies Q(\mathbf{x}_{t+1}) = \mathbb{P}(\mathbf{x}_{t+1})$$

$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \implies Q(\mathbf{x}_{t+1}) = \int \mathbb{P}(\mathbf{x}_{t+1}, \mathbf{x}_{t}) d\mathbf{x}_{t}$$

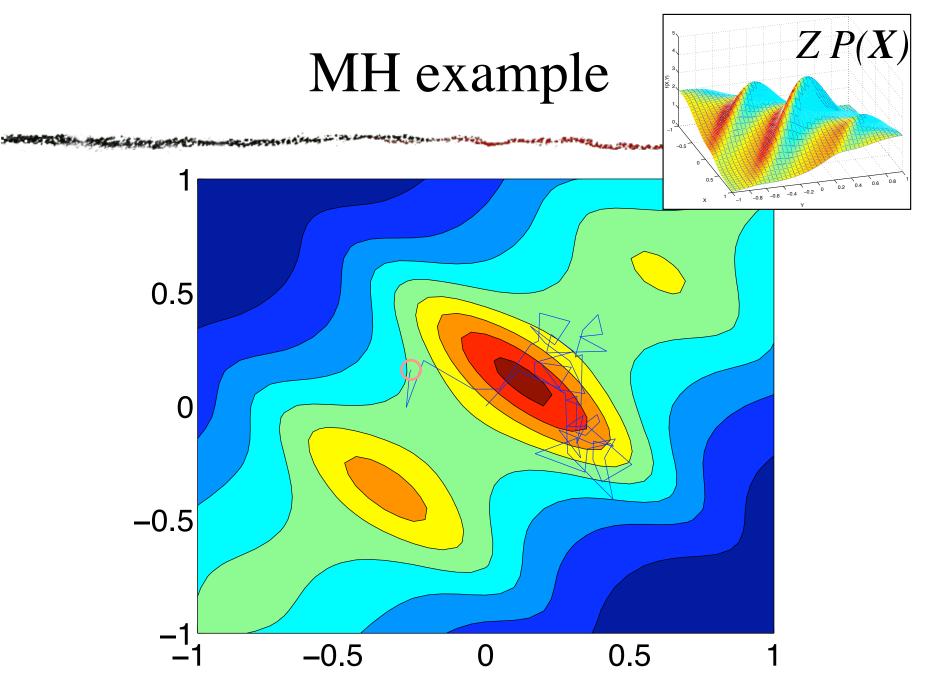
$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \implies Q(\mathbf{x}_{t+1}) = \int \mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_{t}) \mathbb{P}(\mathbf{x}_{t}) d\mathbf{x}_{t}$$

$$Q(\mathbf{x}_{t}) = \mathbb{P}(\mathbf{x}_{t}) \implies Q(\mathbf{x}_{t+1}) = \int \mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_{t}) Q(\mathbf{x}_{t}) d\mathbf{x}_{t}$$

## MH algorithm

note: we don't need to know Z

- $Sample x' \sim Q(x' \mid x)$
- $\circ \ \textit{Compute } p = \frac{P(x')}{P(x)} \frac{Q(x \mid x')}{Q(x' \mid x)}$
- With probability min(1, p), set x := x'
- Repeat for T steps; sample is  $x_1, ..., x_T$  (will usually contain duplicates)



#### MH considerations

- Acceptance rate
- Mixing rate = 1 / mixing time
- Annealing (start at high temperature, reduce to T=1)

# MH proof

## MH proof

- Write T(x'|x) for transition probability
- Write p(x'|x) for acceptance probability

$$\min\left(1, \frac{P(x')}{P(x)} \frac{Q(x \mid x')}{Q(x' \mid x)}\right)$$

 $\circ$  If  $x' \neq x$ , then

$$T(x'|x) = Q(x'|x) p(x'|x)$$

#### Detailed balance

$$P(x)T(x' \mid x) = P(x')T(x \mid x') \qquad \forall x, x'$$

- Proof based on detailed balance
- If we can show detailed balance, P(x) is our stationary distribution:
  - take integral dx on both sides
  - use law of total probability on RHS

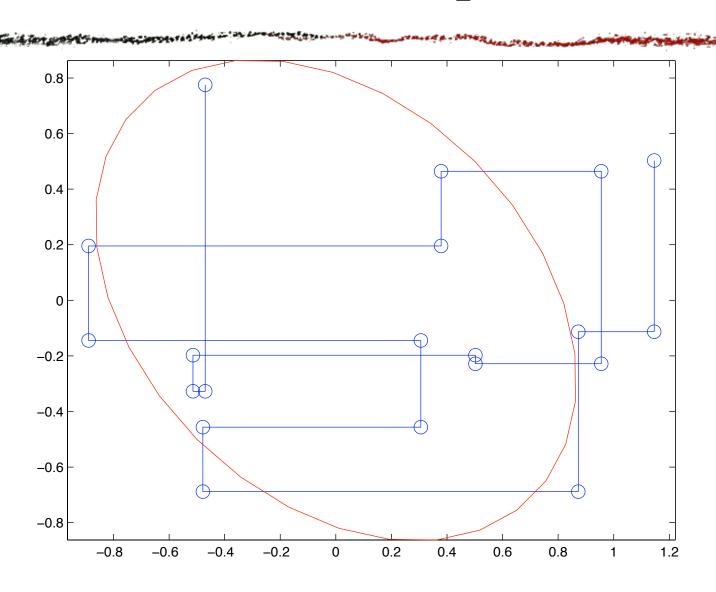
170(x) T(x) (x) = (x) T(x (x))dx = T(x') STOCKTOK P(x) T(x (x) = P(x) Q(x'1x)p(x'1x) = P(x) Q(x'1x)win(1, P(x)) Q(x|x') = 5P(x)Q(x'\x) case ( P(x')Q(x\x') case Z

## Gibbs

## Gibbs sampler

- Special case of MH
- Divide X into blocks of r.v.s B(1), B(2), ...
- Proposal Q:
  - pick a block i uniformly (or round robin, or any other schedule)
  - $\circ$  sample  $X_{B(i)} \sim P(X_{B(i)} \mid X_{\neg B(i)})$

## Gibbs example



## Why is Gibbs useful?

• For Gibbs, 
$$p = \frac{P(x_i', x_{\neg i}')}{P(x_i, x_{\neg i})} \frac{P(x_i \mid x_{\neg i}')}{P(x_i' \mid x_{\neg i})}$$

#### Gibbs derivation

$$\frac{P(x_{i}', x_{\neg i}')}{P(x_{i}, x_{\neg i})} \frac{P(x_{i} \mid x_{\neg i}')}{P(x_{i}' \mid x_{\neg i})}$$

$$= \frac{P(x_{i}', x_{\neg i})}{P(x_{i}, x_{\neg i})} \frac{P(x_{i} \mid x_{\neg i})}{P(x_{i}' \mid x_{\neg i})}$$

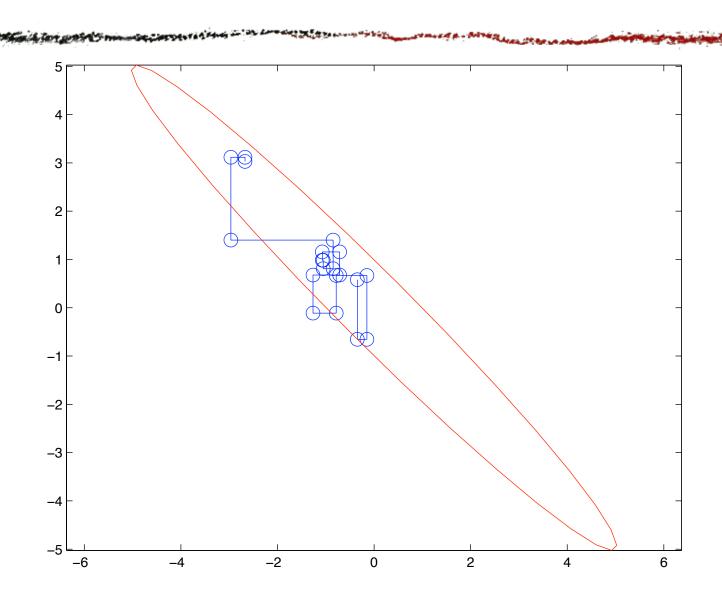
$$= \frac{P(x_{i}', x_{\neg i})}{P(x_{i}, x_{\neg i})} \frac{P(x_{i}, x_{\neg i})/P(x_{\neg i})}{P(x_{i}', x_{\neg i})/P(x_{\neg i})}$$

$$= 1$$

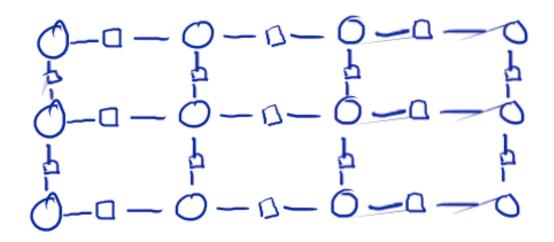
## Gibbs in practice

- Above fact about p means Gibbs is often easy to implement
- Often works well
  - *if* we choose good blocks (but there may be no good blocking!)
- Fancier version: adaptive blocks, based on current x

## Gibbs failure example



## Sequential sampling



- In an HMM or DBN, to sample  $P(X_T)$ , start from  $X_1$  and sample forward step by step
  - $\circ X_{t+1} \sim P(X_{t+1} \mid X_t)$
- $P(X_{1:T}) = P(X_1) P(X_2 | X_1) P(X_3 | X_2) \dots$

#### Particle filter

- Can sample  $X_{t+1} \sim P(X_{t+1} \mid X_t)$  using any algorithm from above
- If we use parallel importance sampling to get N samples at once from each  $P(X_t)$ , we get a particle filter
- Write  $\mathbf{x}_{t,i}$  (i = 1...N) for sample at time t

#### Particle filter

- Want one sample from each of  $P(X_{t+1} | x_{t,i})$
- $\circ$  Have only  $ZP(X_{t+1} \mid x_{t,i})$
- For each i, pick  $\mathbf{x}_{t+1,i}$  from proposal Q(x)
- Compute unnormalized importance weight

$$\hat{w}_i = ZP(\mathbf{x}_{t+1,i} \mid \mathbf{x}_{t,i})/Q(\mathbf{x}_{t+1,i})$$

#### Particle filter

• Normalize weights:

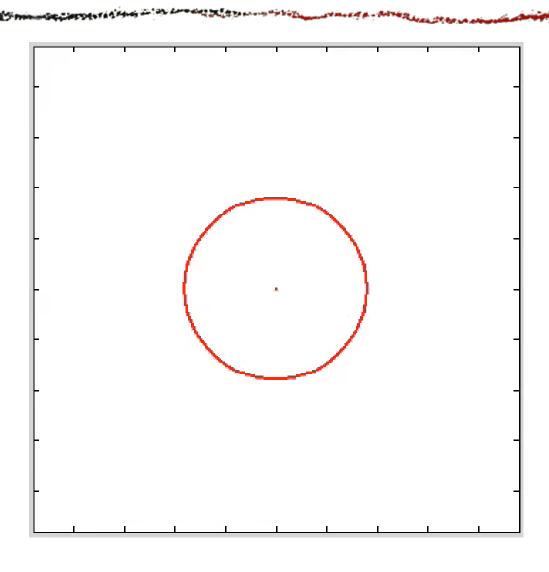
$$\bar{w} = \frac{1}{N} \sum_{i} \hat{w}_{i} \qquad w_{i} = \hat{w}_{i} / \bar{w}$$

- Now,  $(w_i, \mathbf{x}_{t+1,i})$  is an approximate weighted sample from  $P(\mathbf{X}_{t+1})$
- To get an unweighted sample, resample

## Resampling

- Sample N times (with replacement) from  $\mathbf{x}_{t+1,i}$  with probabilities  $w_i/N$ 
  - alternate method: deterministically take  $floor(w_i)$  copies of  $\mathbf{x}_{t+1,i}$  and sample only  $from [w_i floor(w_i)]$
- Each  $\mathbf{x}_{t+1,i}$  appears  $w_i$  times on average, so we're still a sample from  $P(\mathbf{X}_{t+1})$

## Particle filter example



# Learning

## Learning

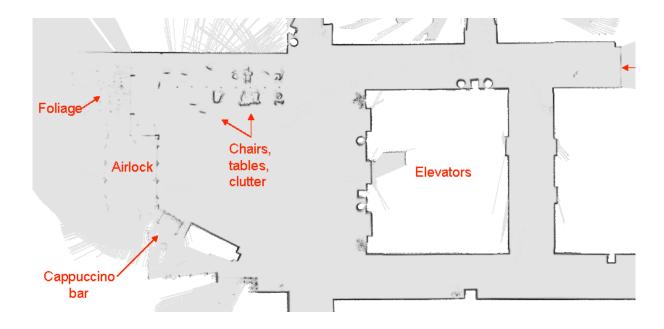
- So far we've assumed our model of the world (factor graph, or SAT formula, or MILP, etc.) was given
- In fact, one of the most important attributes of an intelligent agent is that it learns from experience

## Learning

- Basic learning problem: given some experience, find a new or improved model
- Experience: a sample  $x_1, ..., x_N$
- *Model:* want to predict  $x_{N+1}$ , ...

## Example

- Experience = range sensor readings & odometry from robot
- $\circ$  *Model* = map of the world



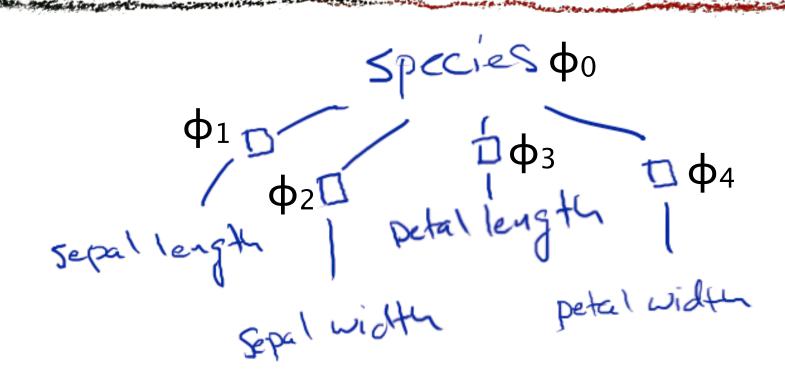
## Example

- Experience = physical measurements of surveyed specimens & expert judgements of their true species
- Model = factor graph relating species to measurements

## Sample data

sepal length	sepal width	petal length	petal width	species
5.1	3.5	1.4	0.2	Iris setosa
5.6	3.0	4.5	1.5	Iris versicolor
4.9	3.0	1.4	0.2	Iris setosa
6.4	2.8	5.6	2.1	Iris virginica
5.8	2.7	4.1	1.0	Iris versicolor

# Factor graph



- One of many possible factor graphs
- $\circ$  Values of  $\Phi s$  not shown, but part of model

### In general

 $\circ$  For our purposes, a model is exactly a joint distribution P(X) over possible samples

## Comparing models

- When is a model P(X) better than another model P'(X)?
- Need to make future decisions based on model; better model = better decisions
- E.g., suppose robot runs on I. versicolor,
   but I. setosa is poisonous to it
- Knowing  $P(I. setosa \mid petal \ length = 4.2)$ lets us weigh risks of eating specimen

### The problem

- We don't know what future examples we'll see, or what decisions we'll have to make about them
- So, we'll use various proxies

#### Conditional model

- Split variables into (X, Y)
- Suppose we always observe X
- Two ways P(X, Y) and P'(X, Y) can differ:
  - $\circ P(X) \neq P'(X)$ , and/or
  - $\circ P(Y \mid X) \neq P'(Y \mid X)$
- First way doesn't matter for decisions
- $\circ$  Conditional model: only specifies  $P(Y \mid X)$

# Conditional model example

- $\circ$  Experience = samples of (X, Y)
- $\circ X = features \ of \ object$
- Y = whether object is a "framling"
- Model = rule for deciding whether a new object is a framling

# Sample data & possible model

tall	pointy	blue	framling
T	T	F	T
T	F	F	T
$oxed{F}$	T	F	$oxed{F}$
T	T	T	$oxed{F}$
T	F	F	T

$$H = tall \land \neg blue$$

## Hypothesis space

- Hypothesis space  $\mathcal{H}=$  set of models we are willing to consider
  - for philosophical or computational reasons
- E.g., all factor graphs of a given structure
- Or, all conjunctions of up to two literals

## A simple learning algorithm

- Conditional learning: samples  $(x_i, y_i)$
- Let *H* be a set of propositional formulae

$$\circ \mathcal{H} = \{ H_1, H_2, \dots \}$$

- *H* is consistent if  $H(x_i) = y_i$  for all i
- *Version space*  $V = \{ all \ consistent \ H \} \subseteq \mathcal{H}$
- Version space algorithm: predict  $y = majority \ vote \ of \ H(x) \ over \ all \ H \in V$

## Framlings

tall	pointy	blue	framling
T	T	F	T
T	F	F	T
$oxed{F}$	T	F	F
T	T	T	F
T	$\overline{F}$	$\overline{F}$	T

∘  $\mathcal{H}$  = { conjunctions of up to 2 literals } = { T, F, tall, pointy, blue, ¬tall, ¬pointy, ¬blue, tall ∧ pointy, tall ∧ blue, pointy ∧ blue, ¬tall ∧ pointy, ... }

### Analysis

- Mistake = make wrong prediction
- If some  $H \in \mathcal{H}$  is always right, eventually we'll eliminate all competitors, and make no more mistakes
- If no  $H \in \mathcal{H}$  is always right, eventually V will become empty
  - e.g., if label noise or feature noise

# Analysis

- $\circ$  Suppose  $|\mathcal{H}| = N$
- How many mistakes could we make?

### Analysis

- $\circ$  Suppose  $|\mathcal{H}| = N$
- How many mistakes could we make?
- Since we predict w/ majority of V, after any mistake, we eliminate half (or more) of V
- Can't do that more than log<sub>2</sub>(N) times

#### Discussion

- *In example,* N = 20,  $log_2(N) = 4.32$
- Made only 2 mistakes
- Mistake bound: limits wrong decisions, as desired
- $\circ$  But, required strong assumptions (no noise, true H contained in  $\mathcal{H}$ )
- Could be very slow!