# Review
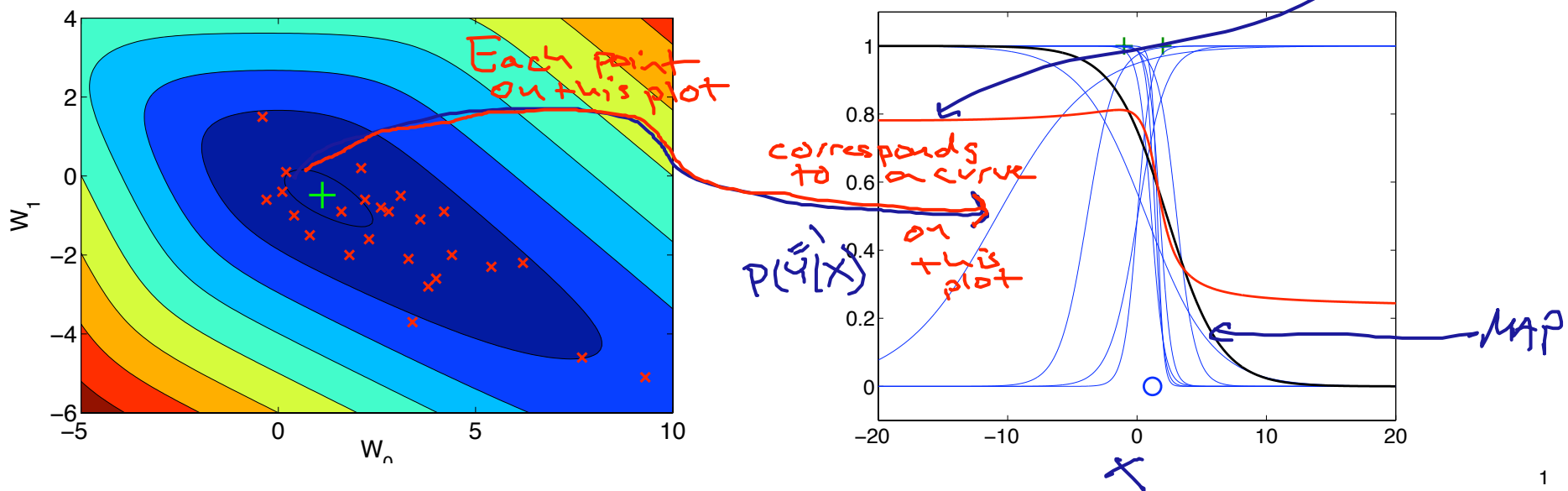
- Multiclass logistic regression

- Priors, conditional MAP logistic regression

- Bayesian logistic regression

  ‣ MAP is not always typical of posterior

  ‣ posterior predictive can avoid overfitting

# Review

- Finding posterior predictive distribution often requires numerical integration

  ▸ uniform sampling

  ▸ importance sampling

  ▸ parallel importance sampling

- These are all ***Monte-Carlo algorithms***

  ▸ another well-known MC algorithm coming up

*we've the house*

# Application: SLAM



Foliage
Airlock
Chairs, tables, clutter
Cappuccino bar
Window
Elevators

Eliazar and Parr, IJCAI-03

# Parallel IS

set $\hat{w}_i = Z \, P(x_i) / Q(x_i)$

$$\bar{w} = \frac{1}{N} \sum_{i=1}^{N} \hat{w}_i$$

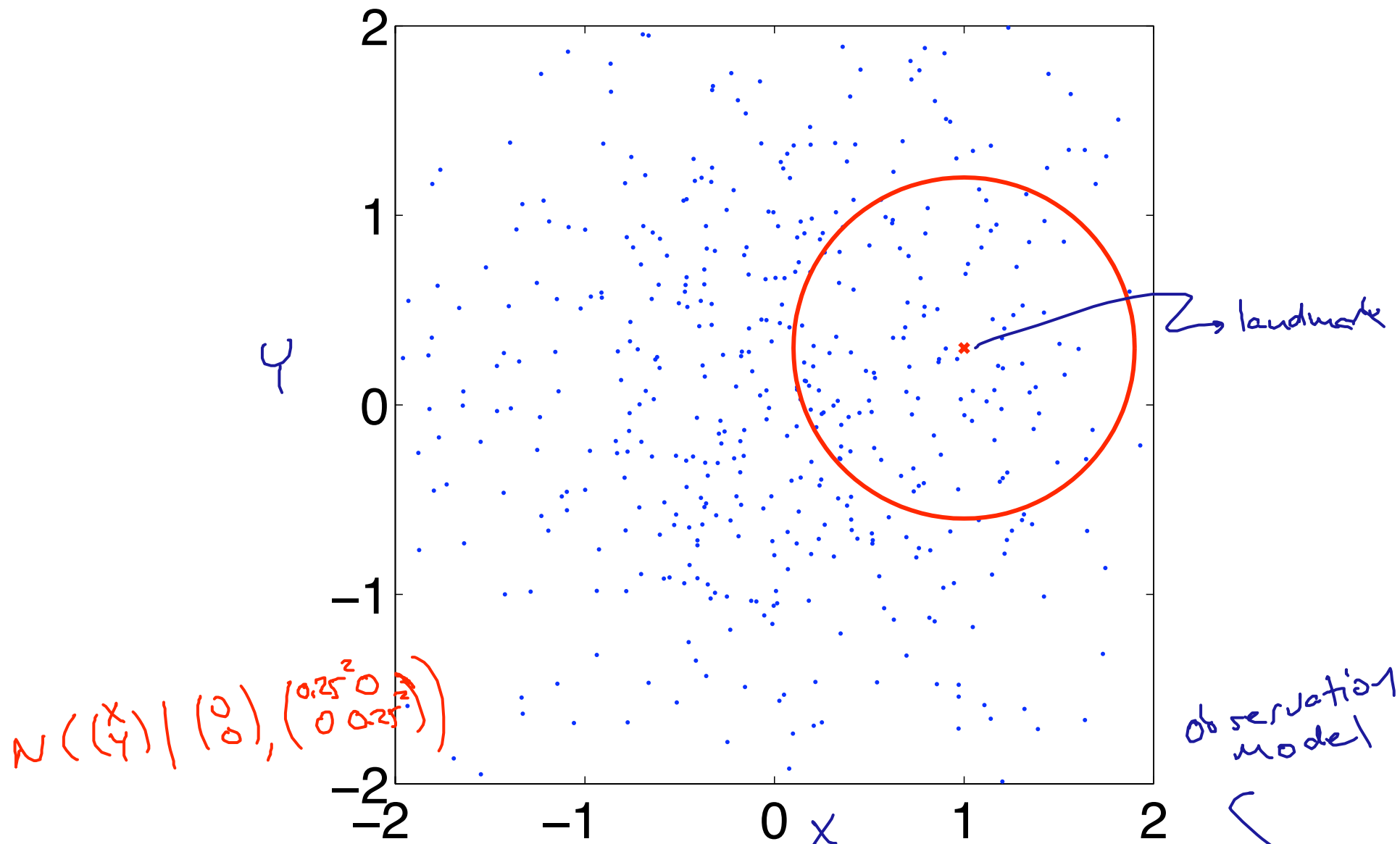$$E(\hat{w}_i) = \int Q(x) \, Z \, P(x) / Q(x) \, dx$$

$$= Z \int P(x) \, dx = Z$$

$$E(\bar{w}) = Z \quad (\text{lower variance})$$

$$E_P(g(x)) \approx \hat{G} = \frac{1}{N} \sum_{i=1}^{N} \frac{\hat{w}_i}{\bar{w}} g(x_i)$$
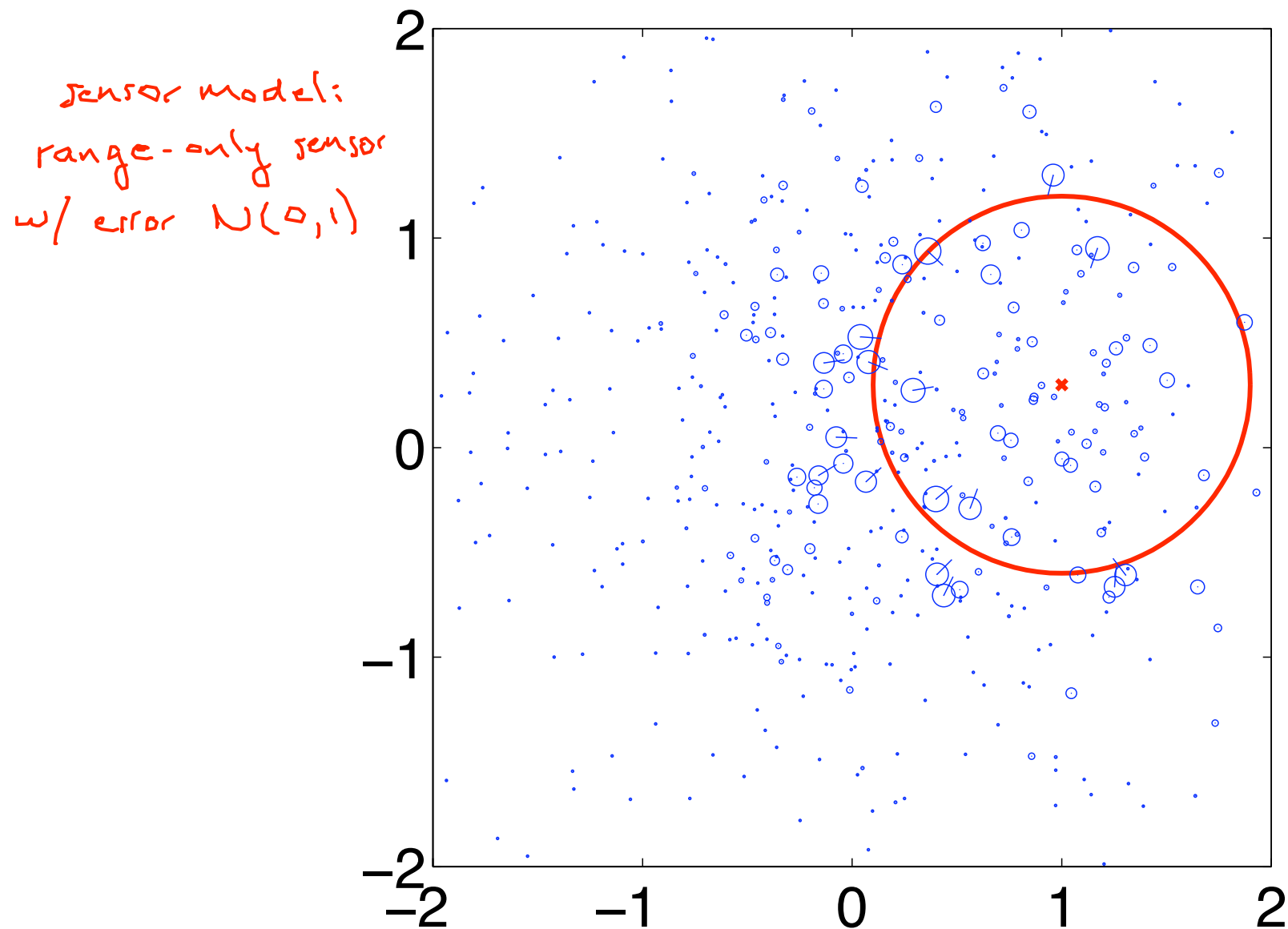
normalized importance wts

# Parallel IS is biased



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

5

$Y$

$X$

$N\left(\binom{X}{Y} \mid \binom{0}{0}, \binom{0.25^2\ 0}{0\ 0.25^2}\right)$

landmark

observation model

importance dist'n = prior

$Q : (X, Y) \sim N(1, 1) \qquad \theta \sim U(-\pi, \pi)$

$f(x, y, \theta) = Q(x, y, \theta) P(o = 0.8 \mid x, y, \theta) / Z$

prior
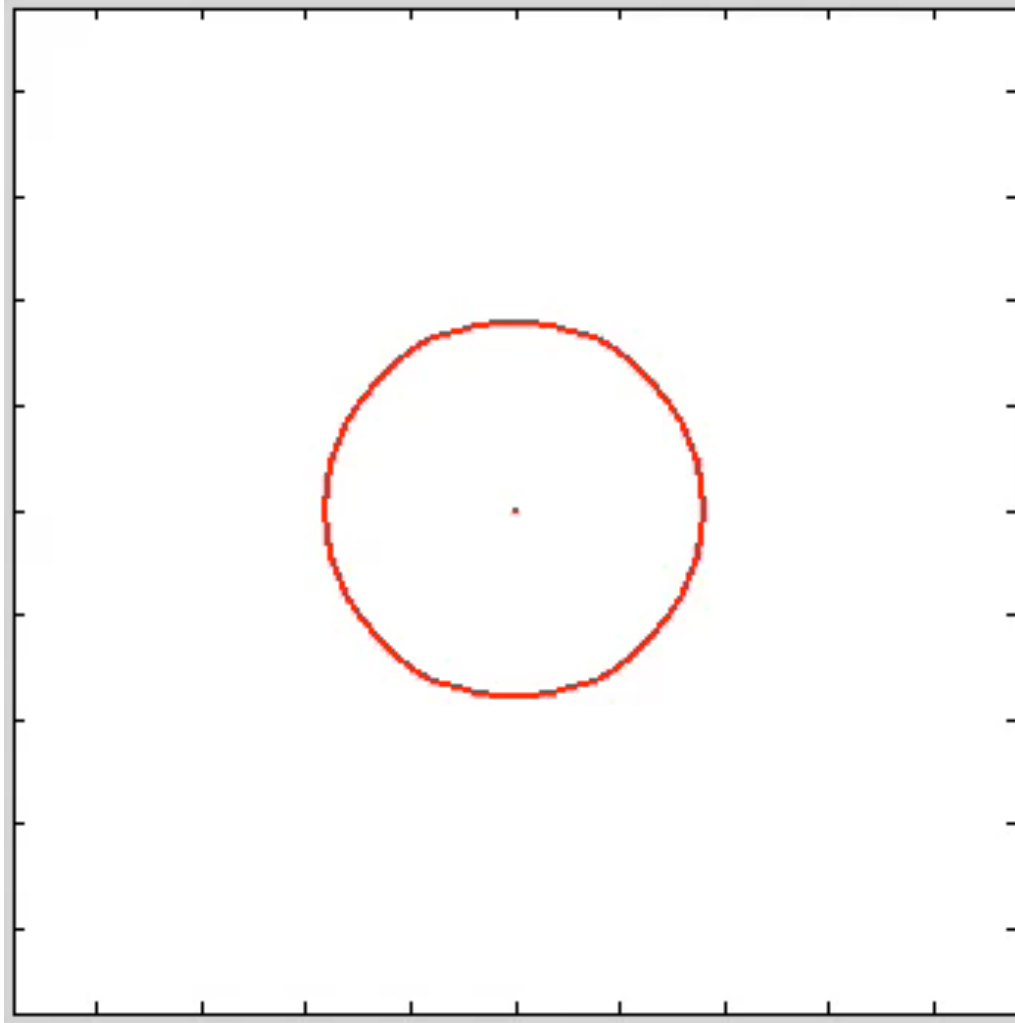
6

sensor model:
range-only sensor
w/ error $N(0,1)$

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

# SLAM revisited

- Uses a recursive version of parallel importance sampling: *particle filter*

  ‣ each sample (particle) = trajectory over time

  ‣ sampling extends trajectory by one step

  ‣ recursively update importance weights and renormalize

  ‣ resampling trick to avoid keeping lots of particles with low weights

# Particle filter example

# Monte-Carlo revisited

- Recall: wanted

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

- Would like to search for areas of high P(x)

- But searching could bias our estimates

# Markov-Chain Monte Carlo



- Randomized search procedure

- Produces sequence of RVs $X_1, X_2, \ldots$

  ‣ Markov chain: satisfies Markov property

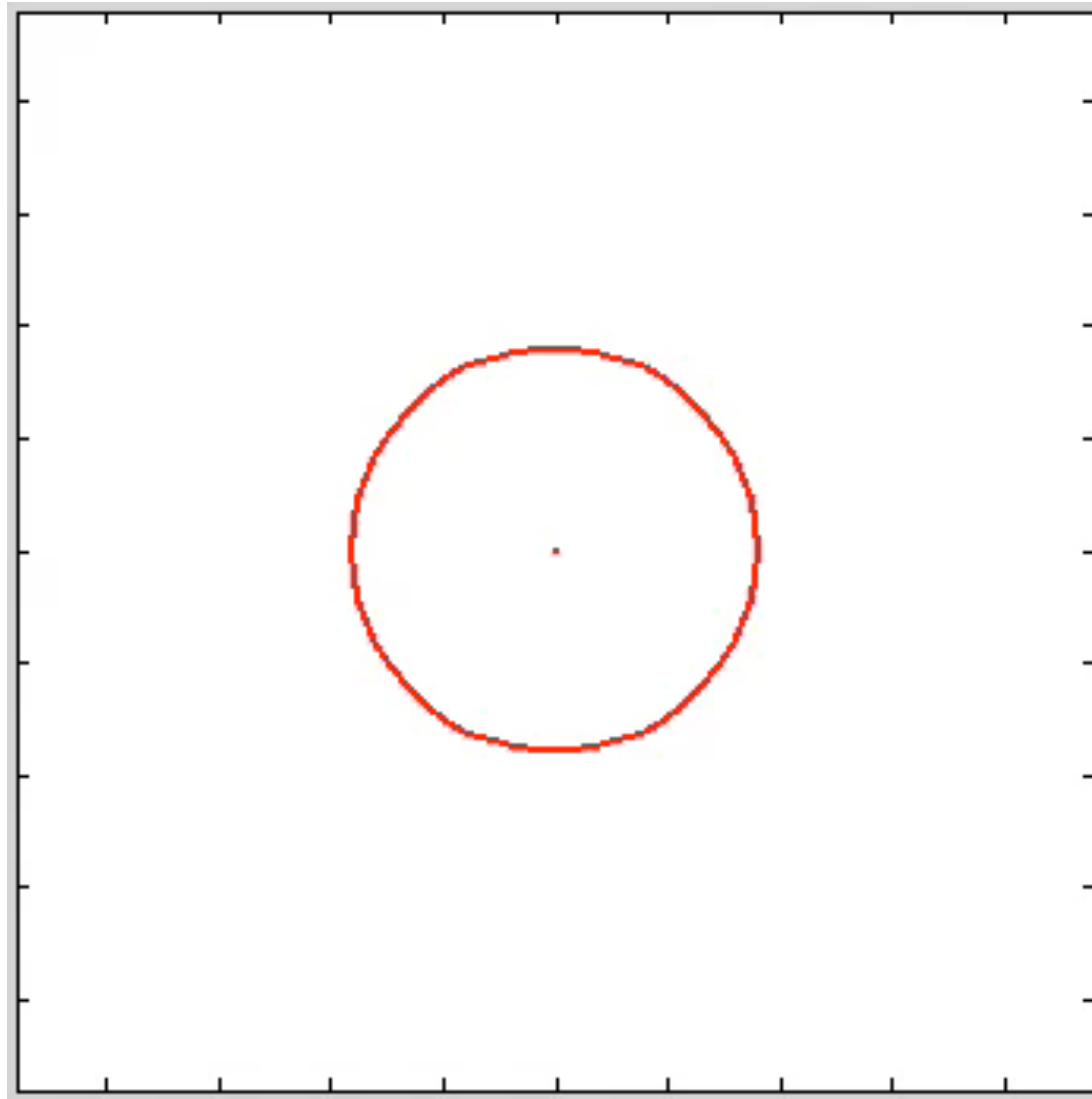  $$(X_1, X_2, \ldots X_{t-1}) \perp (X_{t+1} \ldots) \mid X_t$$

- If $P(X_t)$ small, $P(X_{t+1})$ tends to be larger

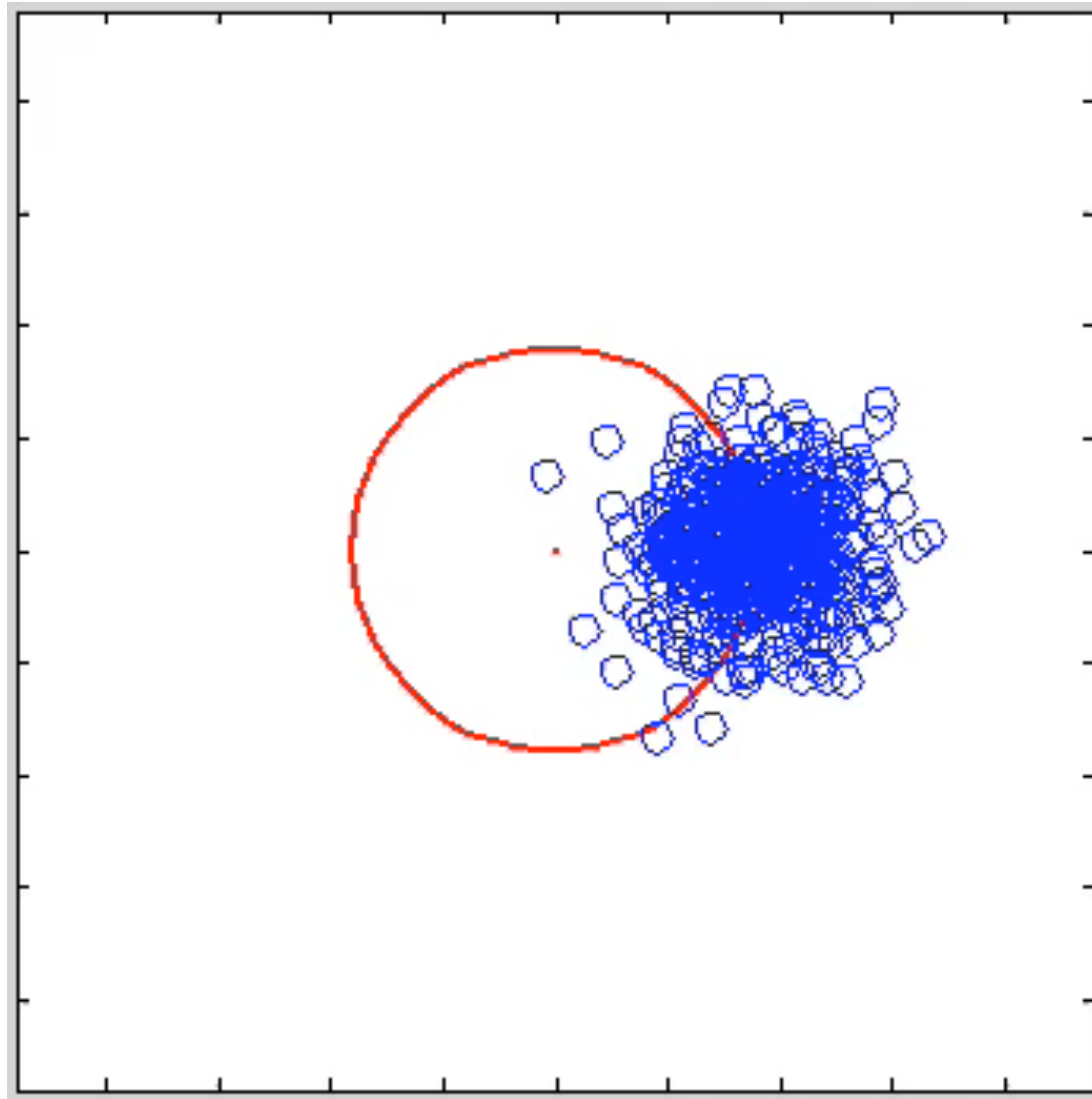- As $t \to \infty$, $X_t \sim P(X)$   "stationary" or limiting dist'n is $P$

- As $\Delta \to \infty$, $X_{t+\Delta} \perp X_t$

  forgetting

# Markov chain

# Stationary distribution

# Markov-Chain Monte Carlo

- As $t \to \infty$, $X_i \sim P(X)$; as $\Delta \to \infty$, $X_{t+\Delta} \perp X_t$

  *burn-in time*        *mixing time*

- For big enough $t$ and $\Delta$, an approximately i.i.d. sample from $P(X)$ is

  ‣ $\{ X_t, X_{t+\Delta}, X_{t+2\Delta}, X_{t+3\Delta}, \dots \}$

- Can use i.i.d. sample to estimate $E_P(g(X))$    *better if $g$ expensive*

$$\hat{G} = \frac{1}{N} \sum_{i=1}^{N} g(X_{t + (i-1)\Delta})$$

$$\hat{G}' = \frac{1}{N} \sum_{i=1}^{N} g(X_{t+1+(i-1)\Delta})$$

- Actually, don't need independence:

$$\hat{G} = \frac{1}{\Delta N} \sum_{i=1}^{\Delta N} g(X_{t+i-1})$$

*better if $g$ cheap*

# Metropolis-Hastings

*or* $ZP(X)$
↑
*i.e., don't need to know normaliz. constant*

- Way to design chain w/ stationary dist'n P(X)

- Basic strategy: start from arbitrary X

- Repeatedly tweak X to get X'

  ▸ If P(X') ≥ P(X), move to X'

  ▸ If P(X') ≪ P(X), stay at X

  ▸ In intermediate cases, randomize

*→key to getting stationary dist'n P*

# Proposal distribution

- Left open: what does "tweak" mean?

- Parameter of MH: Q(X' | X)

  *one-step proposal dist'n*

- Good proposals explore quickly, but remain in regions of high P(X)

- Optimal proposal?   $Q(x'|x) = P(x')$

# Simplest proposal

- **Random walk** MH: → for continuous $X$
  - ▸ $Q(X' \mid X) = N(X' \mid X, \sigma^2 I)$
  - ▸ big σ: move quickly
  - ▸ small σ: remain in regions of high $P$

- Not usually a great proposal, but sometimes the best we have

# MH algorithm

- Initialize $X_1$ arbitrarily $\rightarrow$ *must have* $P(X_1) > 0$

- For t = 1, 2, …:

  ▸ Sample X' ~ Q(X' | $X_t$)    $\rightarrow$ *current iterate*

  ▸ Compute p = $\dfrac{P(x')}{P(x_t)} \dfrac{Q(x_t | x')}{Q(x' | x_t)}$    *acceptance prob.*

  ▸ With probability min(1, p), set $X_{t+1}$ := X'   *accept*

    ▸ else $X_{t+1}$ := $X_t$   *reject*

- Note: sequence $X_1$, $X_2$, … will usually contain duplicates

*If $Q(x' | X_t) = P(x')$ then mix time = 1*

18

# Acceptance rate

- Want **acceptance rate** (avg of min(1,p)) to be large, so we don't get big runs of same X

- Want $Q(X' \mid X)$ to move long distances (to explore quickly)

- Tension between long moves, acceptance rate:

# Random walk MH revisited

- Suppose we always accepted. Then:

$$X_t \sim X_0 + \sum_{i=1}^{t} N(0, \sigma^2 I)$$

$$X_t \sim N(X_0, t\sigma^2 I)$$

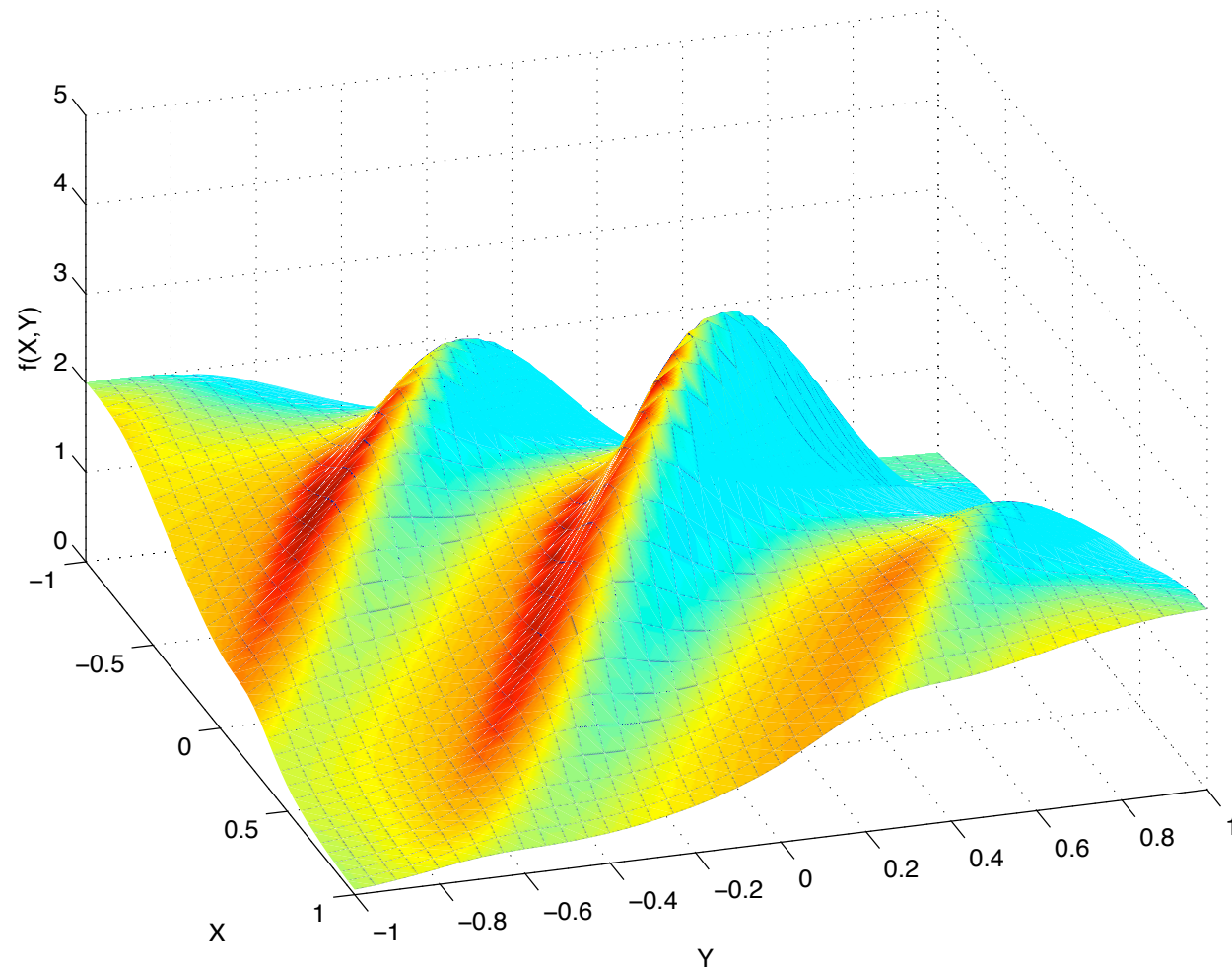$$\text{stddev}(X_t) = \sqrt{t}\, \sigma$$

of one component

Standard deviation of one (any) component of $X_t$

$$= \sqrt{t}\, \sigma$$

- Variance can only be smaller if we reject

# Mixing rate, mixing time

- If we pick a good proposal, we will move rapidly around domain of P(X)

- After a short time, won't be able to tell where we started

- This is short **mixing time** = # steps until we can't tell which starting point we used

- **Mixing rate** = 1 / (mixing time)

# MH example

# MH example