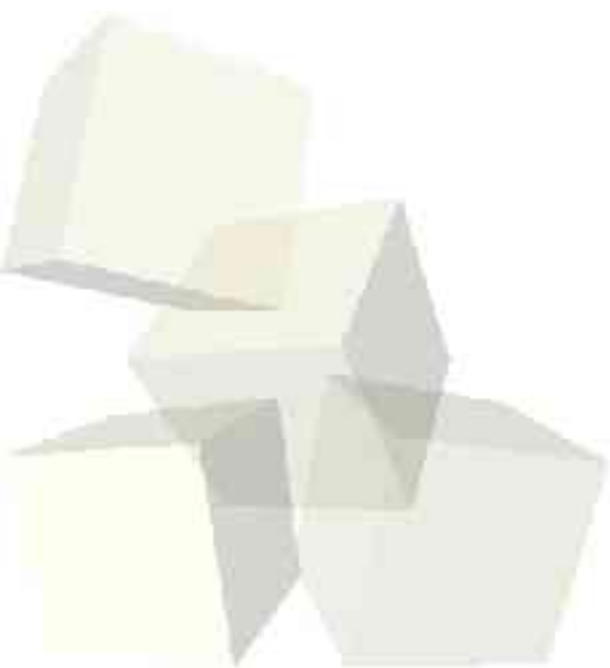# Numerical Integration Monte Carlo style

Mark Huber
Dept. of Mathematics and Inst. of Statistics and Decision Sciences
Duke University
mhuber@math.duke.edu
www.math.duke.edu/~mhuber

Nature laughs at the difficulties of integration.

Pierre-Simon de Laplace

# Darwin visited the Galapagos in 1835

# Darwin noted 14 species of finches



(these 11 photographed by Dr. Robert Rothman)

## Not all finches on all islands!

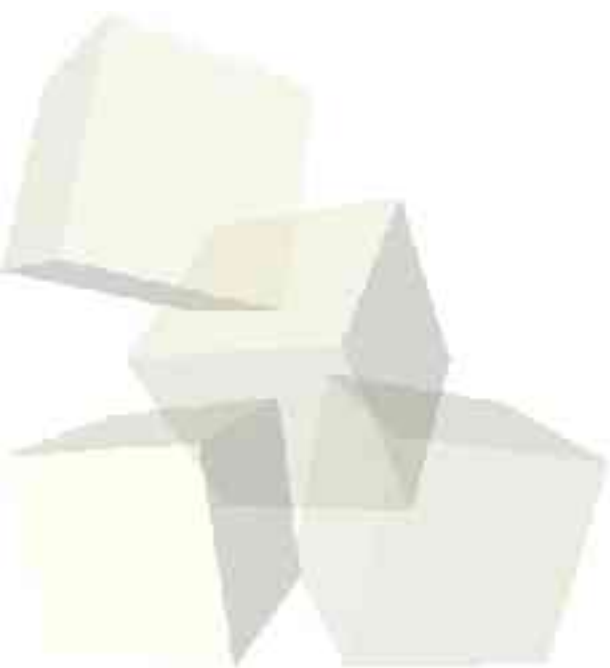|  | *A* | *B* | *C* | *D* | *E* | *...* | *Sums* |
|---|---|---|---|---|---|---|---|
| large ground | 0 | 0 | 1 | 1 | 1 |  | 14 |
| medium ground | 1 | 1 | 1 | 1 | 1 |  | 13 |
| small ground | 1 | 1 | 1 | 1 | 1 |  | 14 |
| sharp-beaked | 0 | 0 | 1 | 1 | 1 |  | 10 |
| ... |  |  |  |  |  |  |  |
| sums | 4 | 4 | 11 | 10 | 8 |  |  |

## 14 types of finches, 17 islands

Is this data random?

Or is it evidence of evolution?

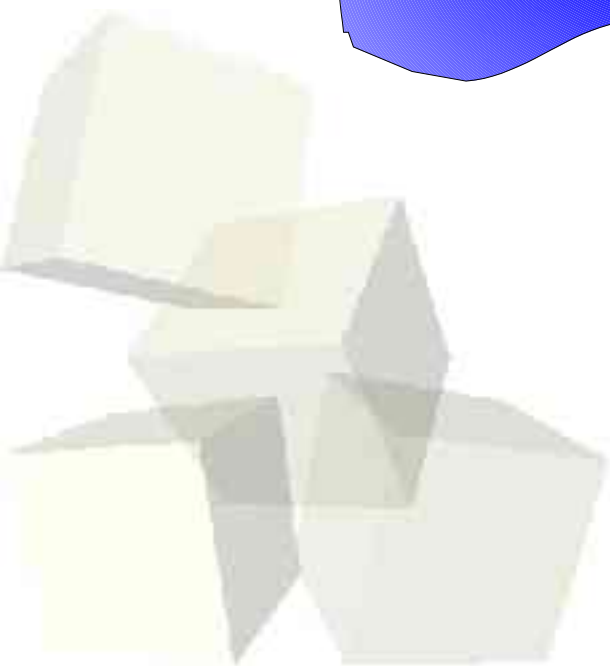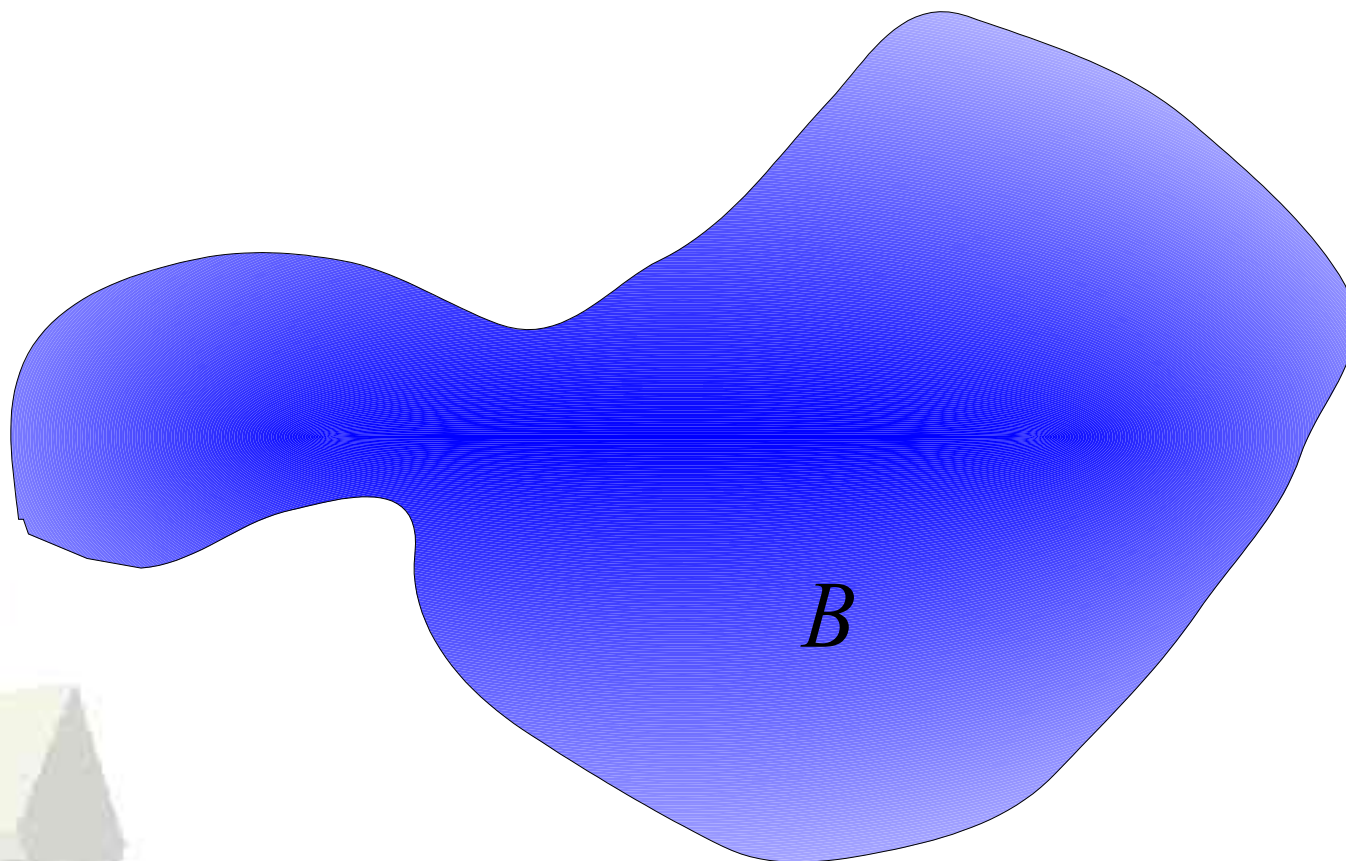To answer deterministically, sum over all tables with same row and column sums

$$2.2 \times 10^{16} \text{ tables!}$$
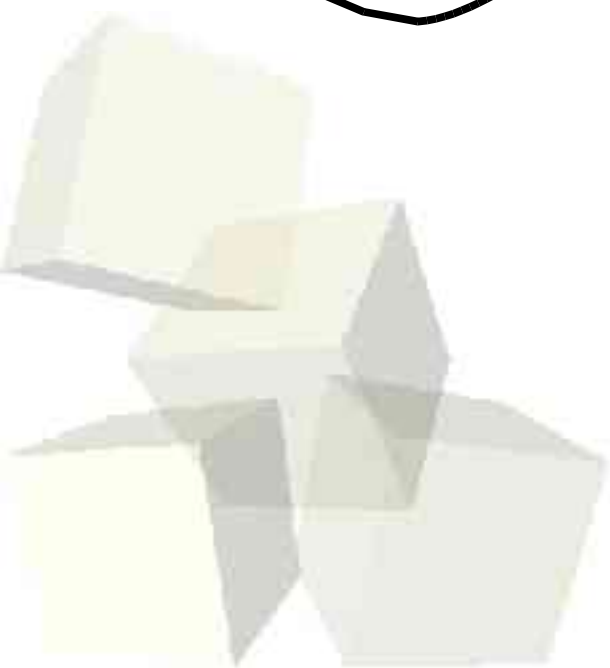
*What is the area of $B$ ?*

$B$

*How many integer points in* $B$ *?*



$B$

These problems have very high dimension

Examples
   Statistical problems
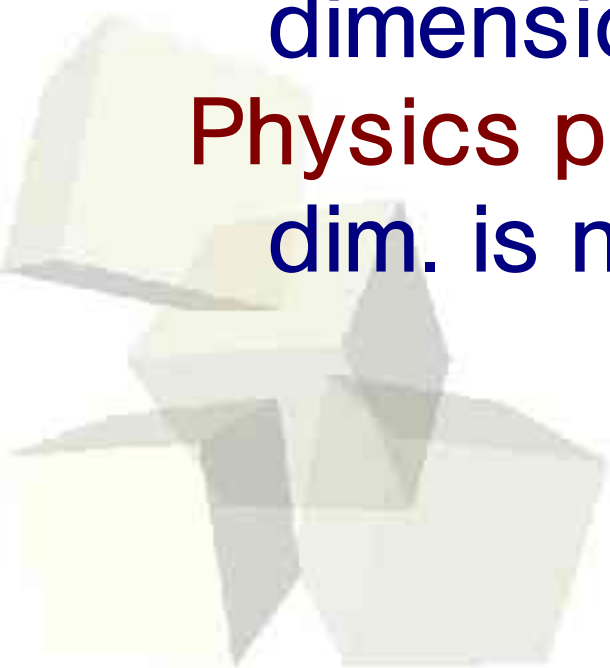      dimension is number of data points
   Network (graph) problems
      dimension is number of nodes
   Physics problems
      dim. is number of interacting entities

Deterministic methods exist
>> Directly count the integer points
>> Running time grows exponential with dim.
>> Trapezoidal Rule, Simpson's Rule, etcetera
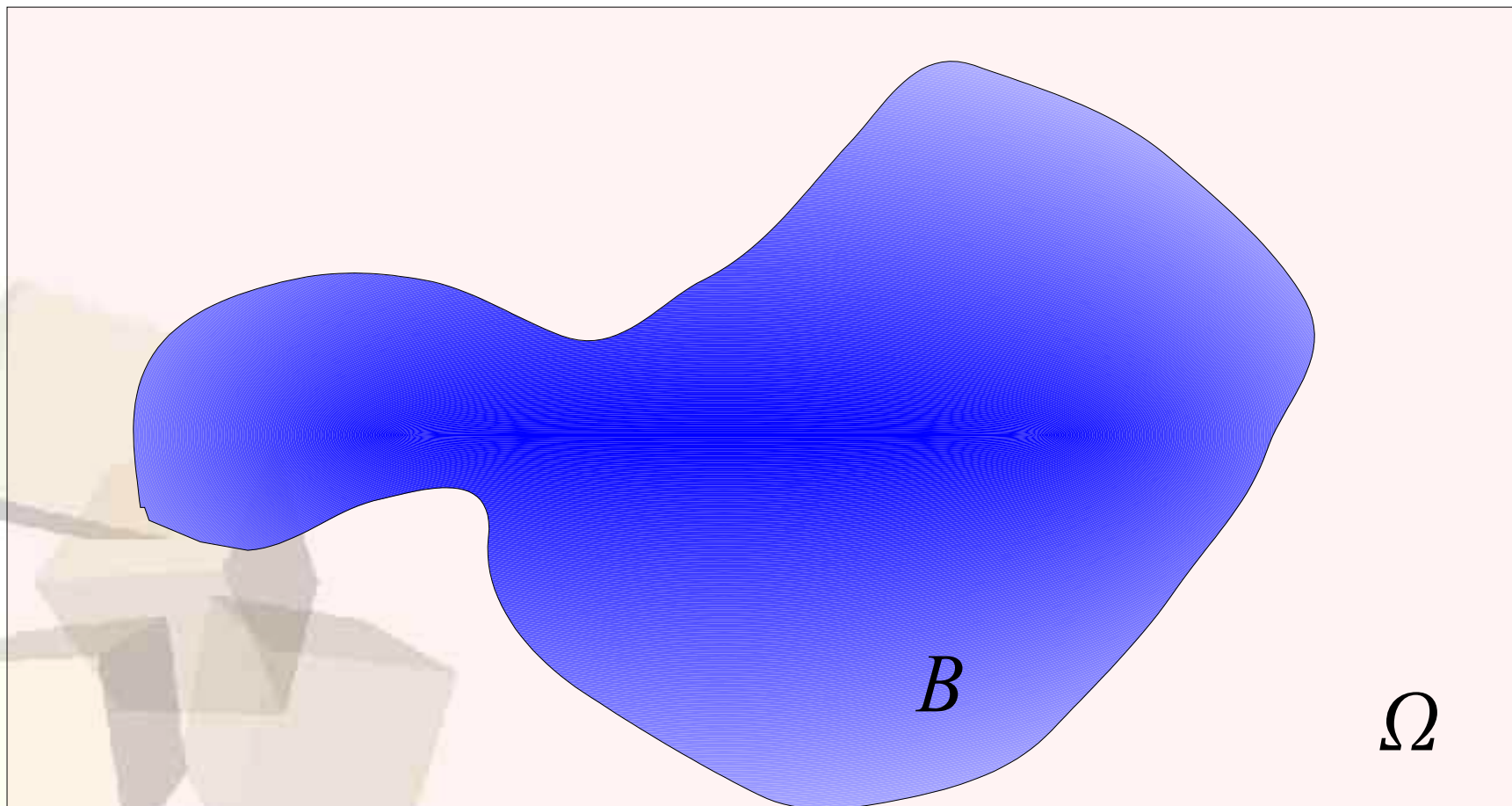>>> Effectively reduce dimension by 1

#P hard
>> Counting the proper colorings of a graph
>> Counting Hamiltonian cycles in a graph

## Acceptance/Rejection

1) Generate samples from bounding region
2) Find percentage lie in $B$
3) Multiply by area of bounding region



$B$
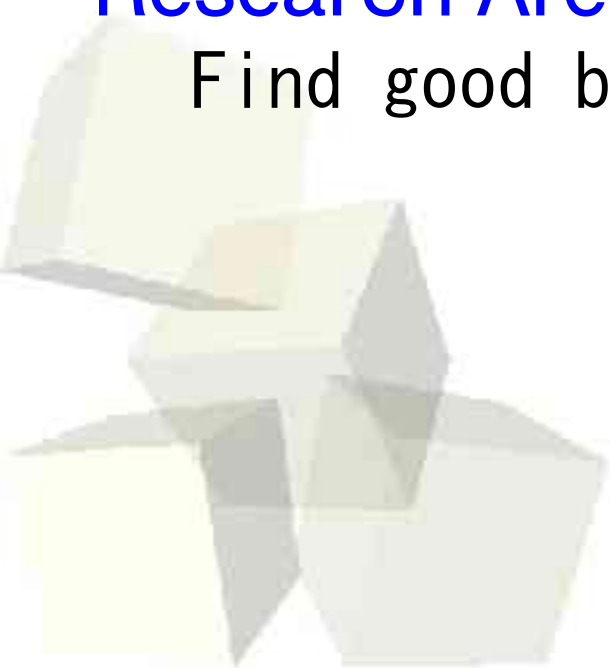
$\Omega$

## The Problem

Need "tight" bounding box
Otherwise need lots of samples for good estimate
Difficult to get in high dimensions

## Research Area #1

```
Find good bounding boxes for actual high
              dimensional problems of interest.
```
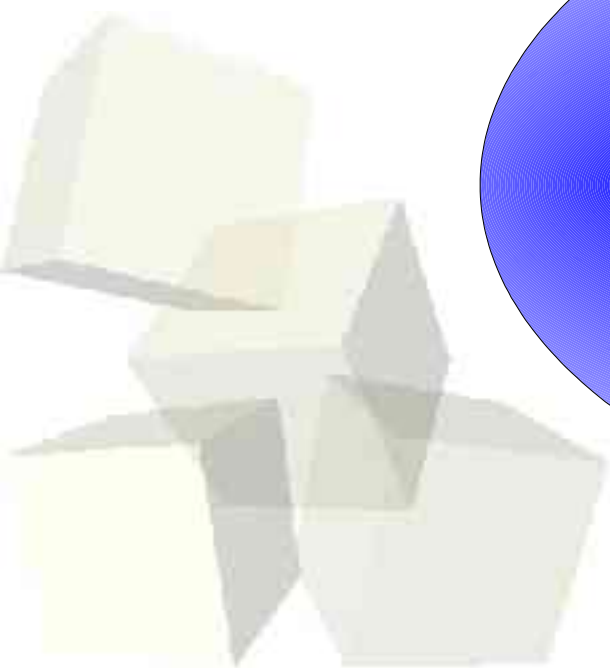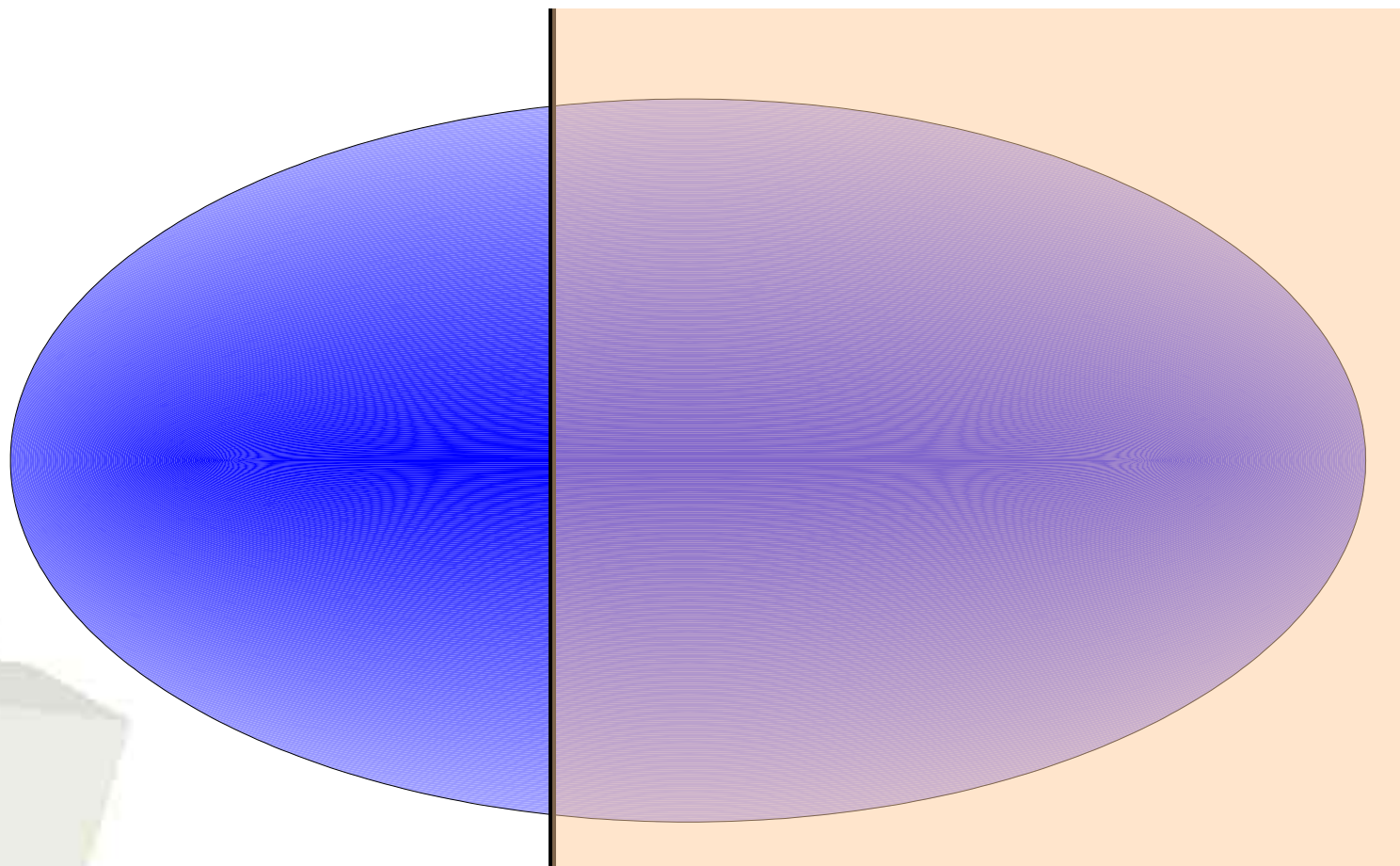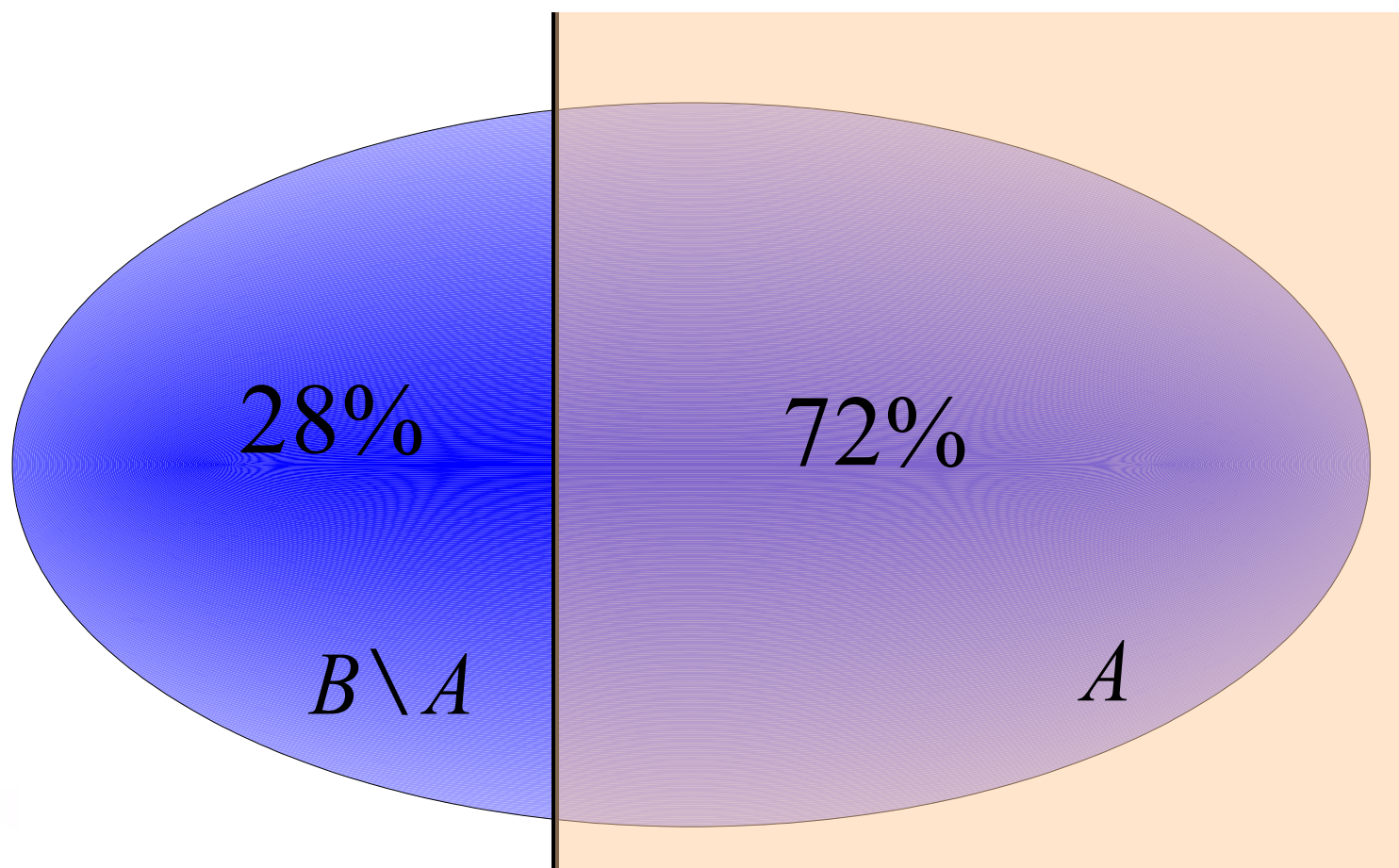
Many times, problem reducible
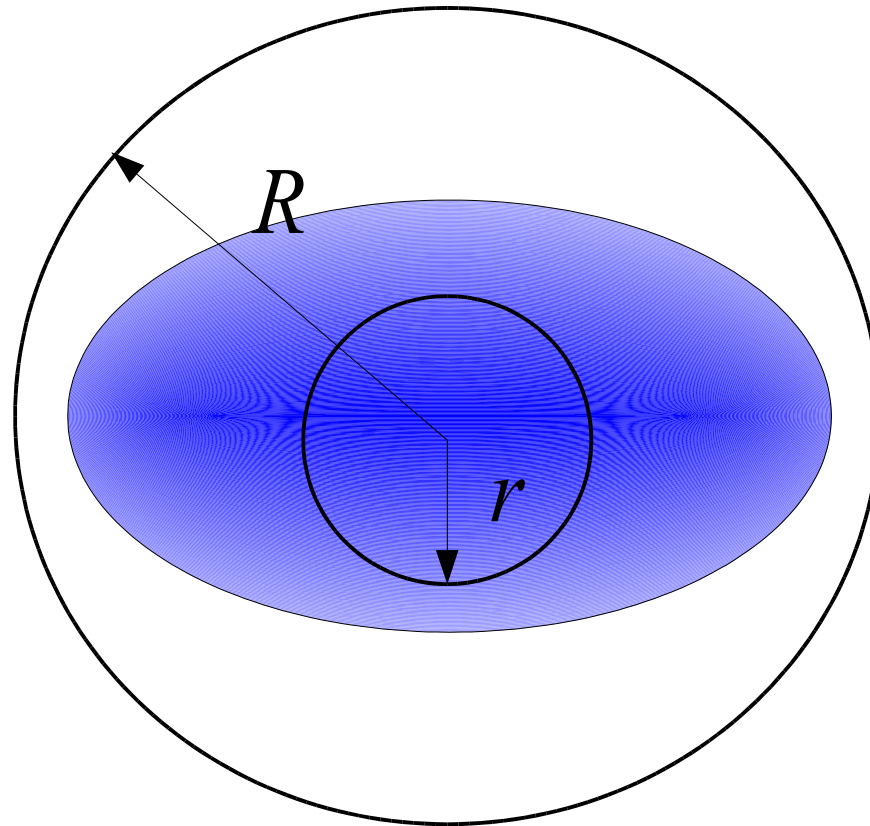Jerrum, Valiant, Vazirani, 1986
Example: convex regions

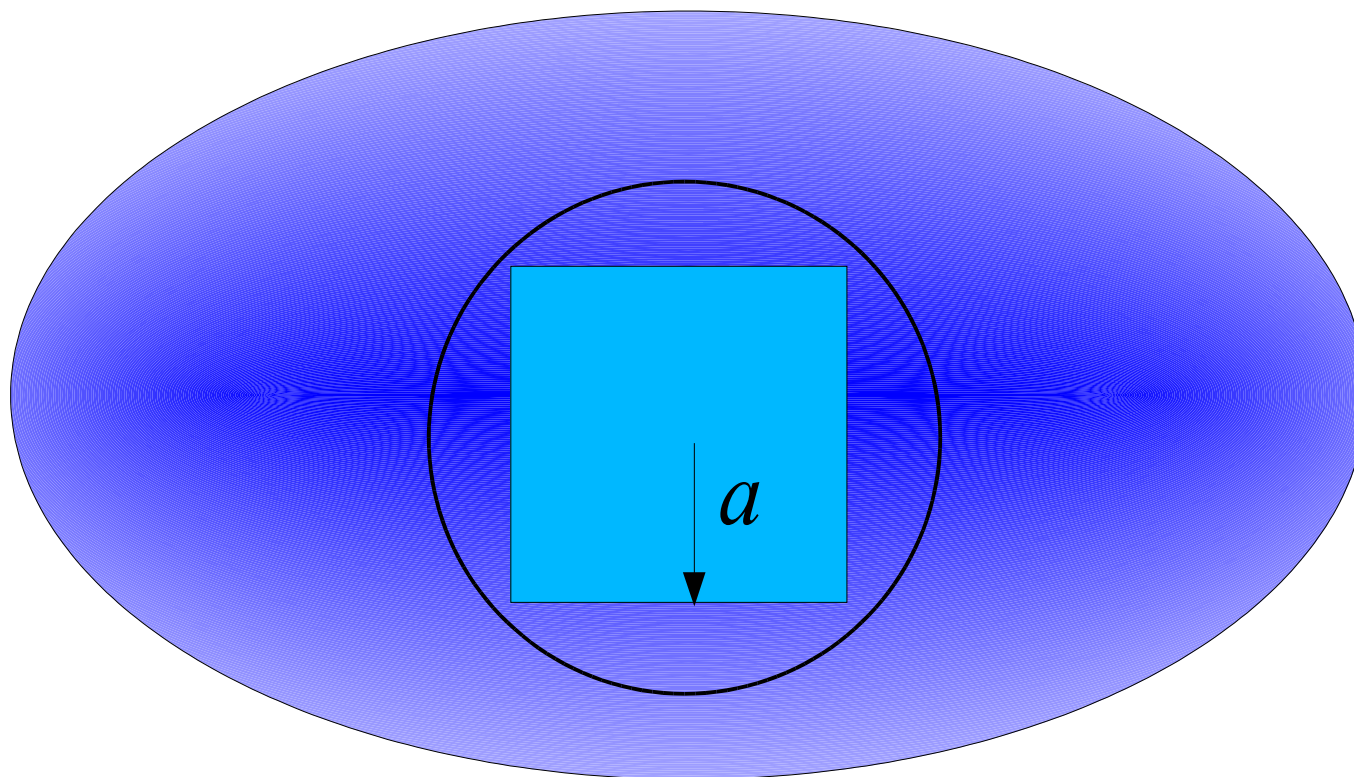$$\text{vol}(B) = \text{vol}(A) \times \frac{\text{vol}(B)}{\text{vol}(A)}$$

Estimate $\text{vol}(B)/\text{vol}(A)$

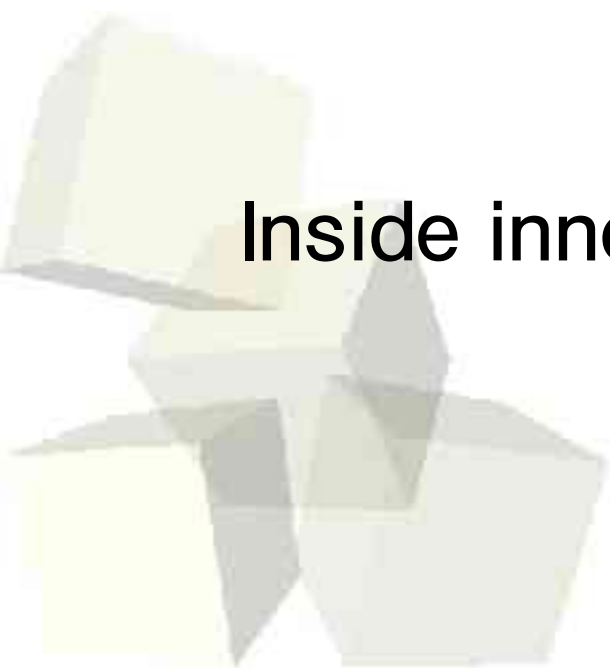(even with this help, can't come within factor of 2 efficiently with deterministic methods [Elekes 86])
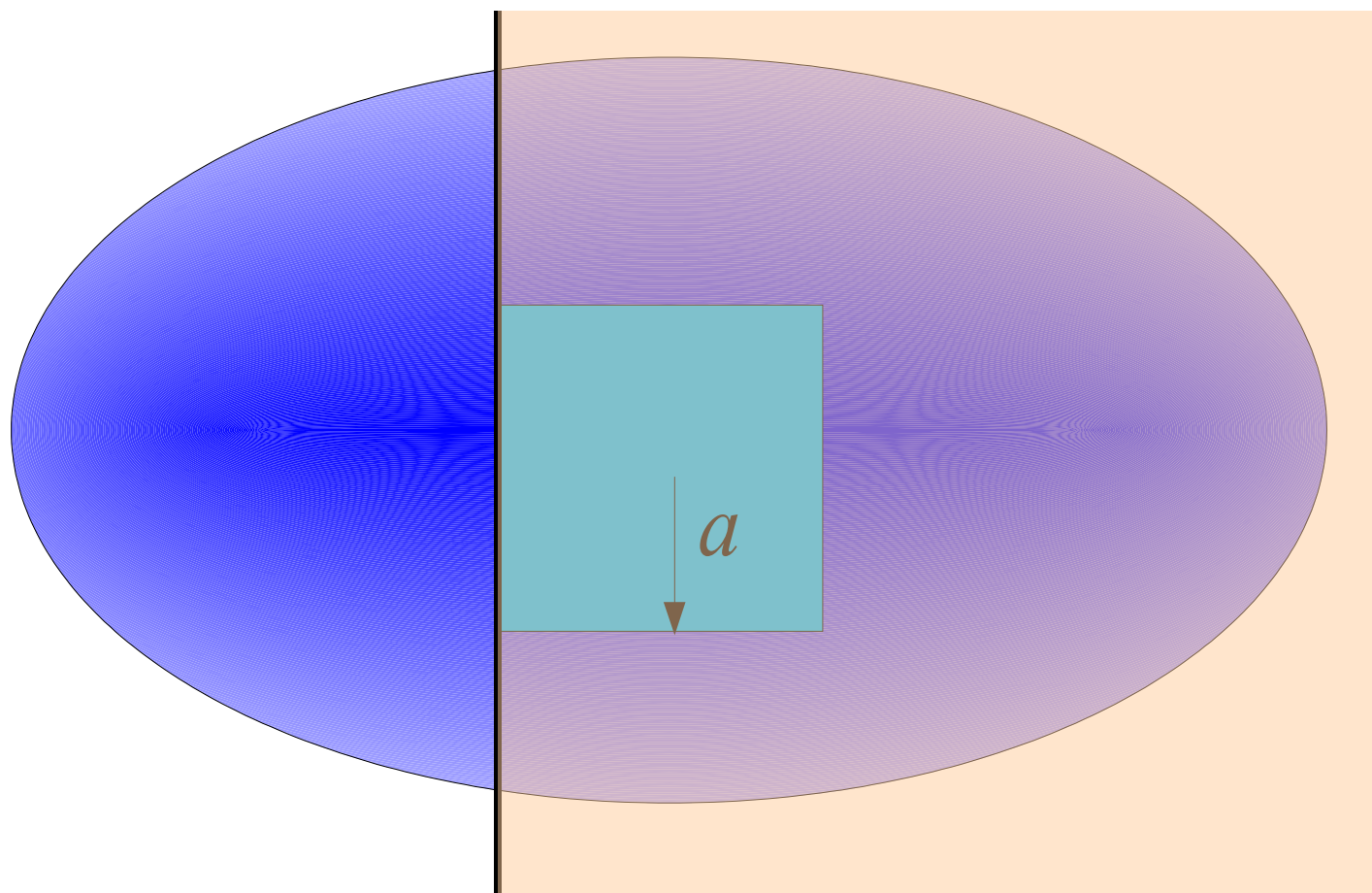


$$\rho = \frac{r}{R} \text{ large}$$

Inside inner ball, box half edge length

$$a = r/\sqrt{(\dim)}$$
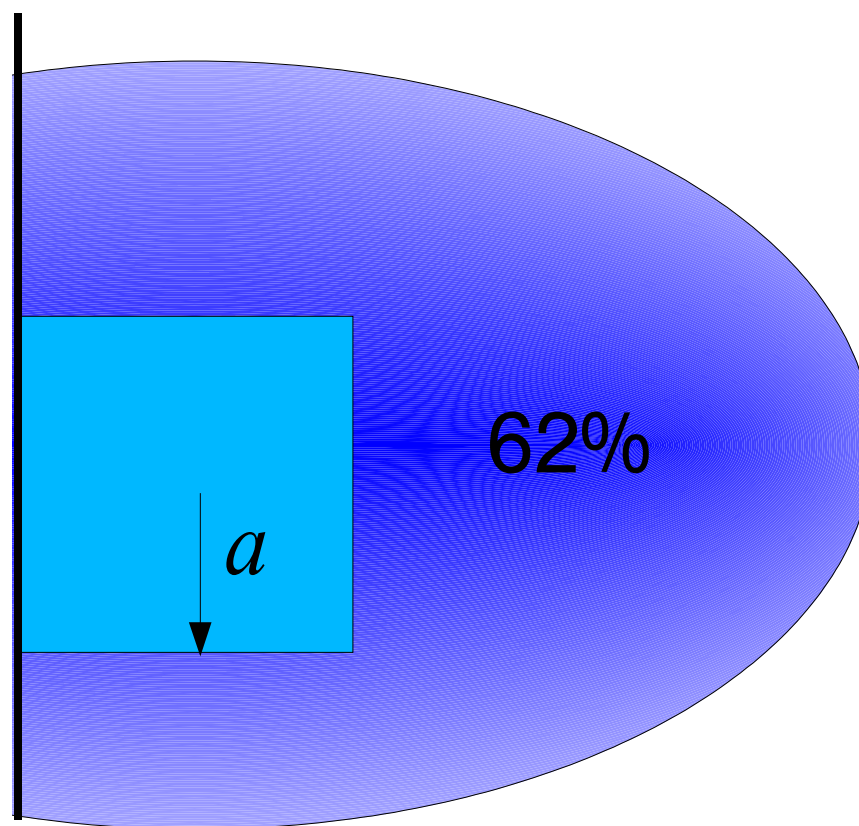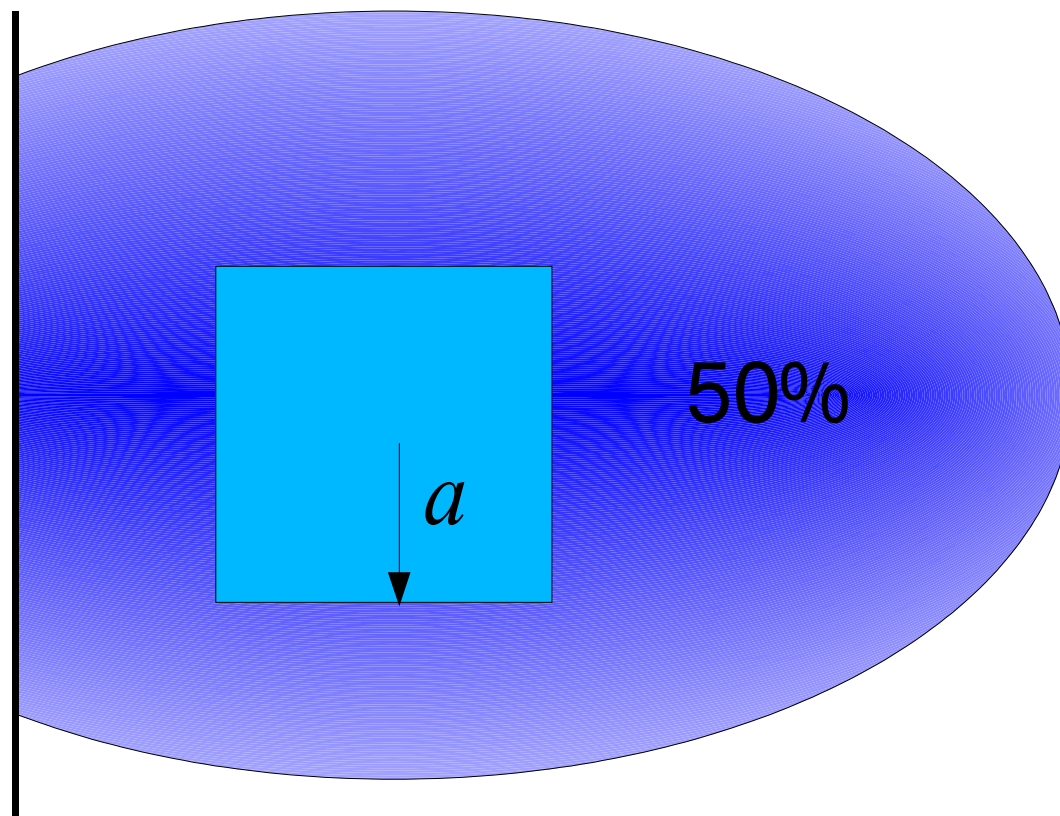
Slice off region to right of box
  Generate lots of random samples
  Estimate percent of area in sliced region

62%

$a$

If region with box at least 50%
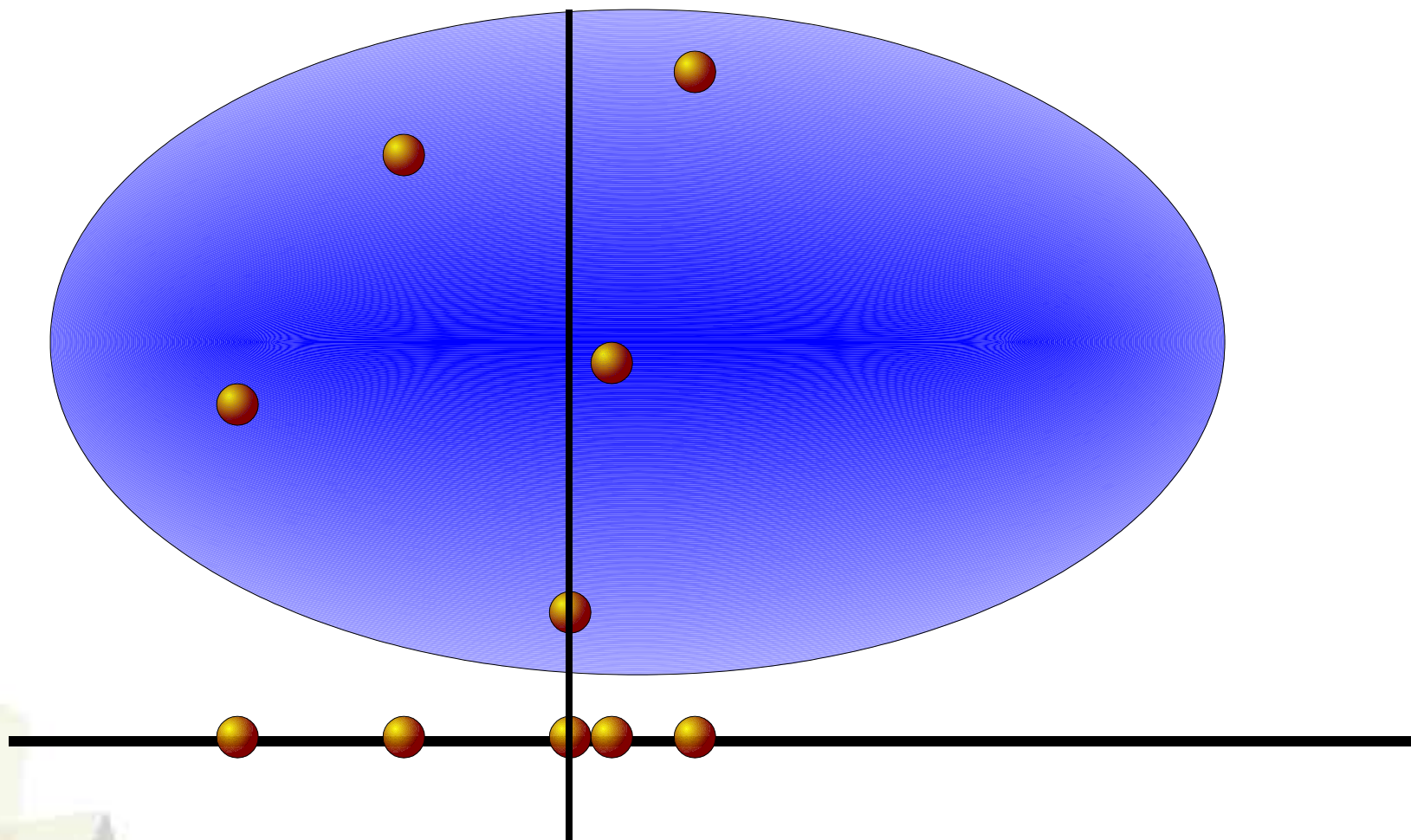use as reduced problem

50%

$a$

Else
find median, use that instead

1) grab samples from body
2) project onto one dimension
3) take median of projections

$50\%$

$a$

Either
1) Match one facet of box or
2) Volume of body reduced by 1/2

Note

$$(2R)^{\dim} \geq \mathrm{vol}(\text{original } B)$$

Volume of body after many steps

$$(2R)^{\dim}(1/2)^n \geq \mathrm{vol}(B \text{ after } n \text{ steps})$$

For center box

$$\mathrm{vol}(\text{center box}) = (2a)^{\dim} \geq [2R\rho/\sqrt{(\dim)}]^{\dim}$$

So most steps that can be taken

$$M := 2d + d(\log(d/\rho))/\log(2)$$

To get median need [Cohen 97][Huber 98]

$$O\left(\log\left(1/\delta\right)/\epsilon^2\right) \text{ samples}$$

To get within $\epsilon$ of answer with probability $1-\delta$

Overall, if $M$ steps taken need $\epsilon' = \epsilon/M$

$$O\left(M^3 \log\left(M/\delta\right)\right) \text{ total samples}$$

$$O\left(\dim^3 \log^2\left(\dim/\delta\right)\right) \text{ total samples}$$

Polynomial in the dimension!

## Most used method:  Markov chains



Pick a direction uniformly at random
Move to a uniform point staying inside body

$$O\left(\dim^7\right) \text{time}$$    [Kannen, et. al. 94]
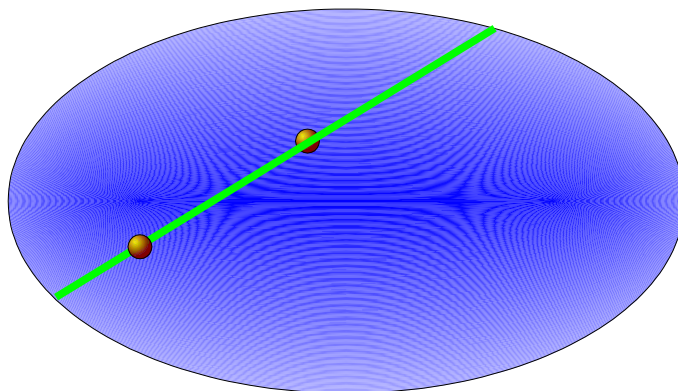
**Research Area #2**

Can bound for Markov chain be improved?

Originally $O\left(\dim^{27}\right)$ steps

**Research Area #3**

Can perfect sampling methods be used for this problem?

Some of my current research questions:

Data from unknown mixtures of distributions
(ex:  responders versus nonresponders to drugs)
Perfect matchings in a graph
(ex:  astronomical data is doubly truncated)
Multinormal distribution on positive orthant
Contingency tables with extra constraints
(ex:  perhaps columns represent age)
The many worlds version of the Ising model
Self organizing lists
(because who has time to organize their own lists?)

Monte Carlo methods are the only known way to handle high dimensional numerical integration

Many interesting questions remain:

Better envelopes for acceptance/rejection
Better Markov chains
Perfect sampling algorithms instead of MC