

Website

<https://courses.engr.illinois.edu/bioe582>



University of Illinois BIOE 582 (Fall 2017) Statistics and Algorithms in Genomic Biology

[Home](#) [Schedule](#)

Instructors

Sergei Maslov: maslov@illinois.edu

Office: IGB 3406

Office hours: by appointment

Until about mid-November when he will be replaced with

Paul Jensen: pjens@illinois.edu



University of Illinois BIOE 582 (Fall 2017) Statistics and Algorithms in Genomic Biology

[Home](#) [Schedule](#)

Date	Topic	Slides	Matlab code	Homework/Exams
8/22	Why probability and statistics in Genomics?	Lecture 0		

Syllabus

Topical Outline:

- Discrete, continuous, and multivariate probability distributions in genomics
- Parameter estimation and confidence intervals
- Hypotheses testing for one or multiple samples
- Linear regression
- Techniques of gene expression analysis: clustering algs, co-expression networks
- Graph-theoretical approaches to biological networks
- Intro to genomes
- Intro to sequencing technologies
- Microbial genome assembly
- How to find and annotate genes (annotation RAST and NCBI annotation pipeline)
- Final Project: Sequencing and assembling clinical isolates of oral streptococci.

Grading: midterm (40%), final project (60%).

Final Project: Teams of students will complete a de novo assembly for an unsequenced bacteria of clinical relevance. Students will load and run a Oxford Nanopore sequencer and apply standard tools to assemble, quality control, and annotate a genome. The final sequence will be deposited in the NCBI genome repository.

Foundations of Probability

Random experiments

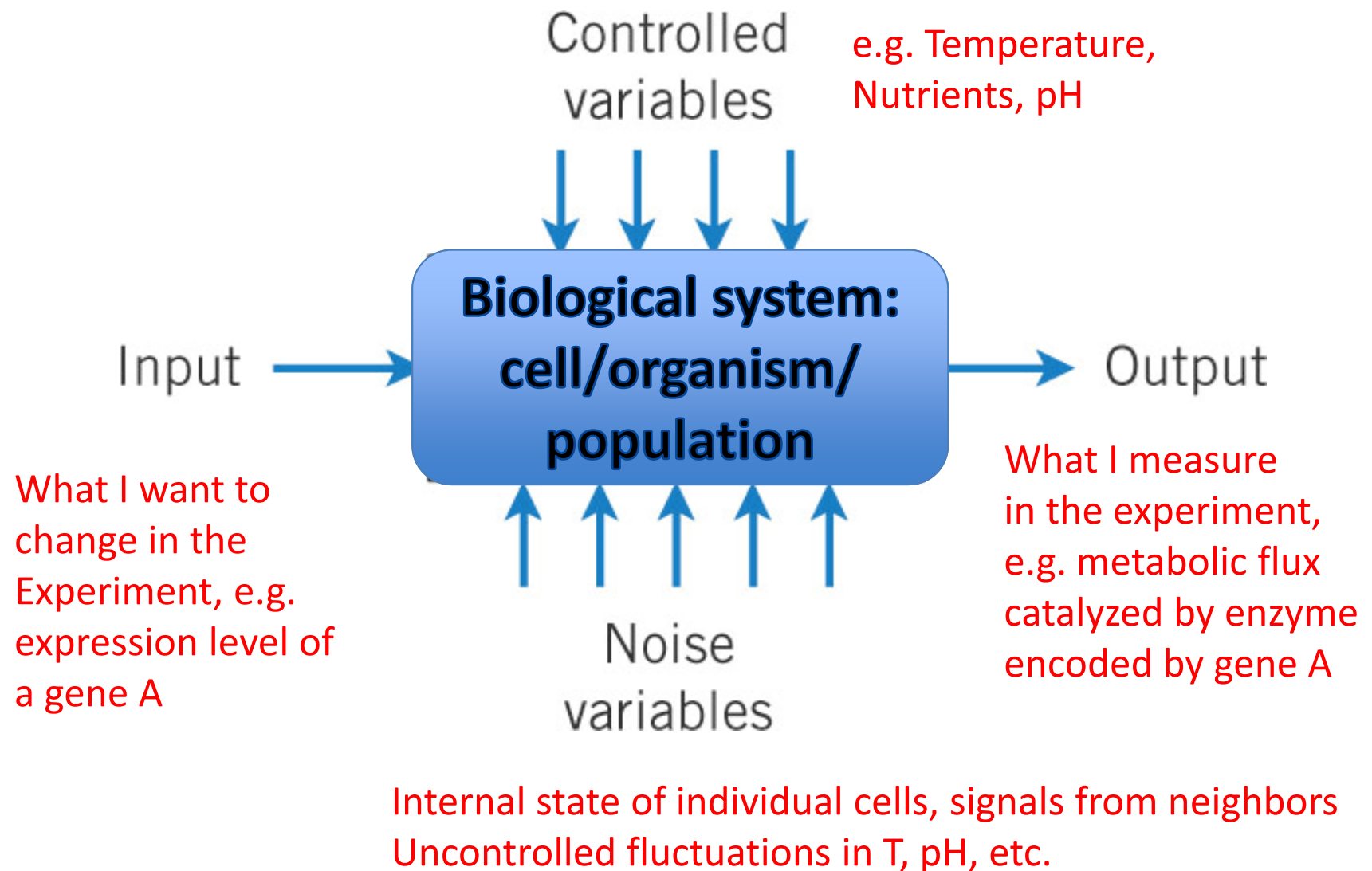
Sample spaces

Venn diagrams of
random events

Random Experiments

- An **experiment** is an operation or procedure, carried out under controlled conditions
 - Example: measure the metabolic flux through a reaction catalyzed by a given enzyme
- An experiment that can result in **different outcomes**, even if repeated in the same manner every time, is called a **random experiment**
 - Cell-to-cell variability due to history/genome variants
 - Noise in external parameters such as temperature, nutrients, pH, etc.
- **Evolution** offers ready-made random experiments
 - Genomes of different species
 - Genomes of different individuals within a species
 - Individual cancer cells

Changes in Input & Controlled parameters + Noise Produce Output Variation



Sample Spaces

- Random experiments have unique outcomes.
- The set of all possible outcomes of a random experiment is called the sample space, S .
- S is discrete if it consists of a finite or countable infinite set of outcomes.
- S is continuous if it contains an interval (either a finite or infinite width) or multiple intervals of real numbers.

Examples of Sample Spaces

- Experiment measuring the abundance of mRNA expressed from a single gene

$S = \{x | x > 0\}$: continuous.

- Bin it into four groups

$S = \{\textit{below 10}, \textit{10-30}, \textit{30-100}, \textit{above 100}\}$: discrete.

- Is gene “on” (say ≥ 30) or “off” (< 30)?

$S = \{\textit{true}, \textit{false}\}$: logical/Boolean.

- Time course of expression of 25,000 human genes measured every 1hr for 3 hours:

$S = \{x_1(1), x_2(1), \dots, x_{25000}(1),$
 $x_1(2), x_2(2), \dots, x_{25000}(2),$
 $x_1(3), x_2(3), \dots, x_{25000}(3) \mid \text{all } x_i(t) > 0\}$

Event

An event (E) is a **subset of the sample space** of a random experiment, i.e., **one or more** outcomes of the sample space.

- The **union** of two events is the event that consists of all outcomes that are contained in either of the two events. We denote the union as $E_1 \cup E_2$
- The **intersection** of two events is the event that consists of all outcomes that are contained in both of the two events. We denote the intersection as $E_1 \cap E_2$
- The **complement** of an event in a sample space is the set of outcomes in the sample space that are not in the event. We denote the complement of the event E as E' (sometimes E^c or \bar{E})

Examples

Discrete

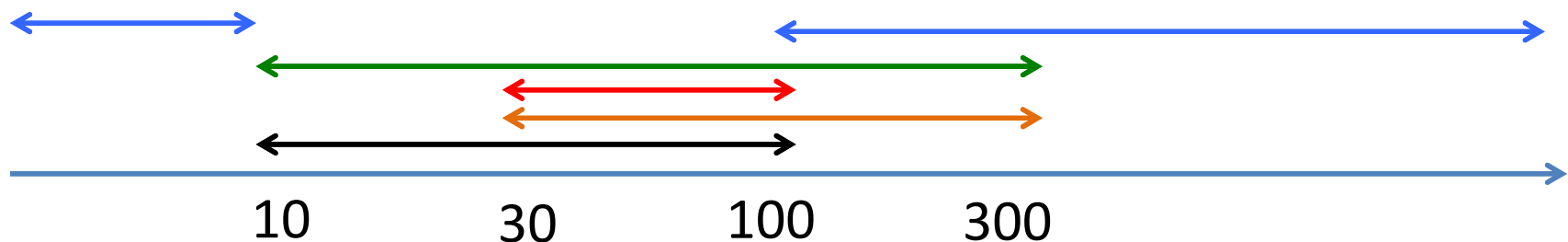
1. Assume you toss a coin once. The sample space is $S = \{H, T\}$, where H = head and T = tail and the event of a head is $\{H\}$.
2. Assume you toss a coin twice. The sample space is $S = \{(H, H), (H, T), (T, H), (T, T)\}$, and the event of obtaining exactly one head is $\{(H, T), (T, H)\}$.

Continuous

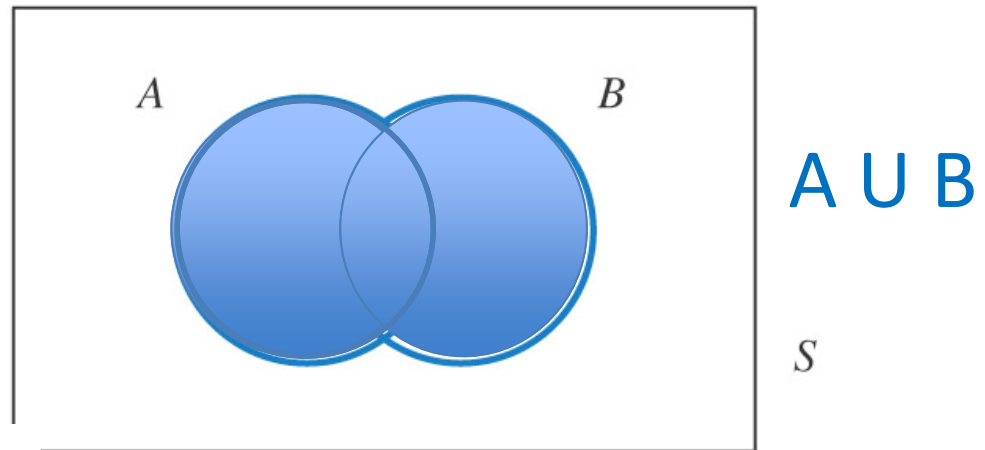
Sample space for the expression level of a gene: $S = \{x | x > 0\}$

Two events:

- $E1 = \{x | 10 < x < 100\}$
- $E2 = \{x | 30 < x < 300\}$
- $E1 \cap E2 = \{x | 30 < x < 100\}$
- $E1 \cup E2 = \{x | 10 < x < 300\}$
- $E1' = \{x | x \leq 10 \text{ or } x \geq 100\}$



Venn diagrams



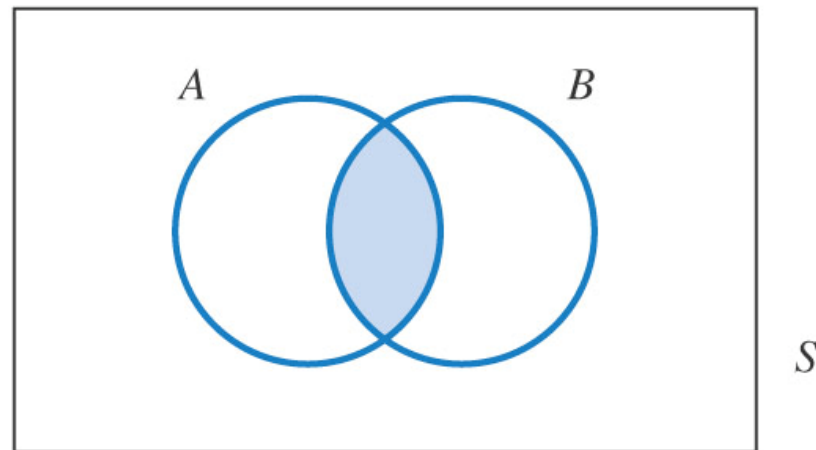
John Venn (1843-1923)
British logician

Find
5 differences
in beard and
hairstyle



John Venn (1990-)
Brooklyn hipster

Venn diagrams

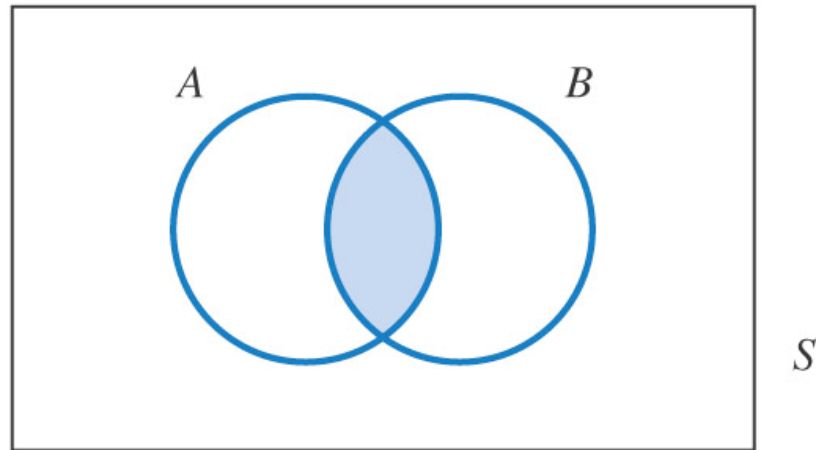


Which formula describes the blue region?

- A. $A \cup B$
- B. $A \cap B$
- C. A'
- D. B'

Get your i-clickers

Venn diagrams



Which formula describes the blue region?

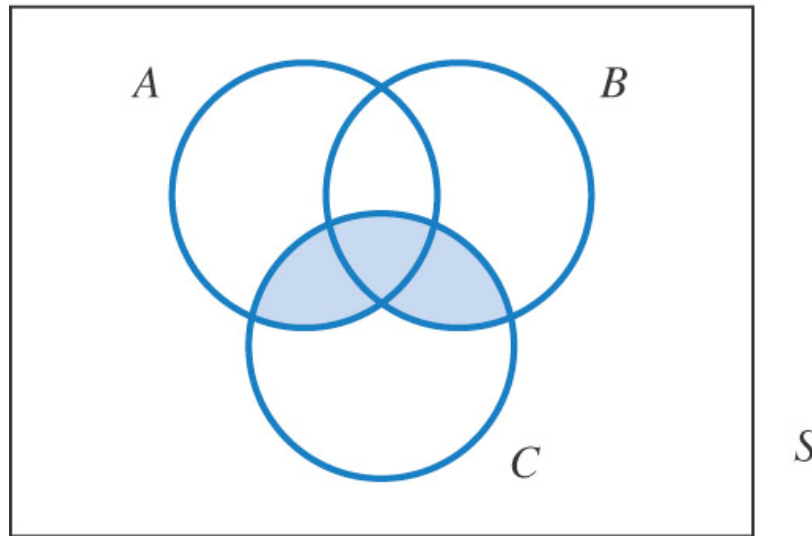
A. $A \cup B$

B. $A \cap B$

C. A'

D. B'

Venn diagrams

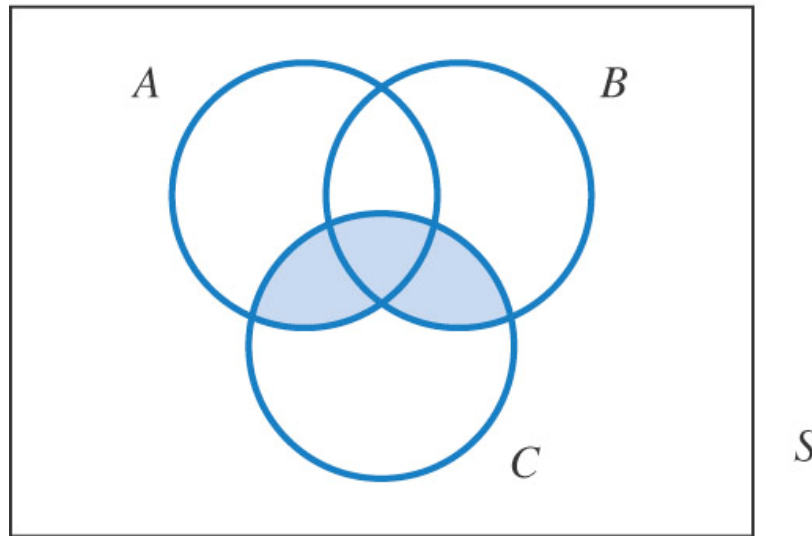


Which formula describes the blue region?

- A. $(A \cup B) \cap C$
- B. $(A \cap B) \cap C$
- C. $(A \cup B) \cup C$
- D. $(A \cap B) \cup C$

Get your i-clickers

Venn diagrams



Which formula describes the blue region?

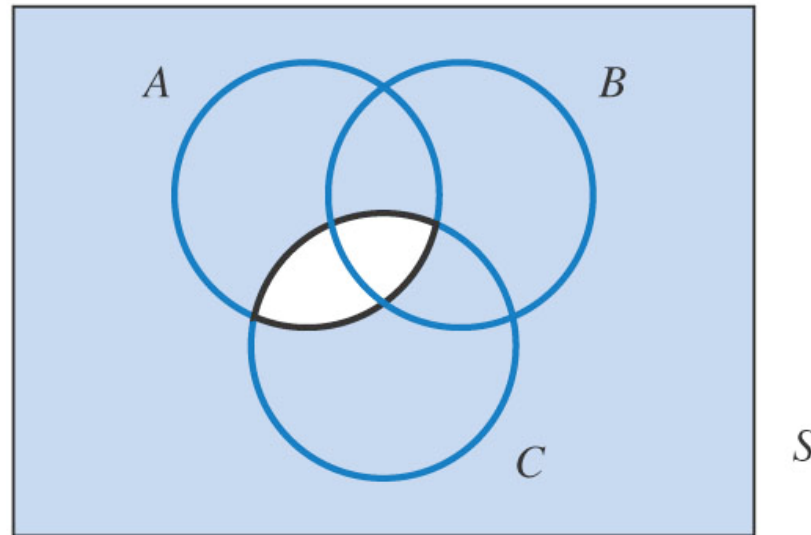
A. $(A \cup B) \cap C$

B. $(A \cap B) \cap C$

C. $(A \cup B) \cup C$

D. $(A \cap B) \cup C$

Venn diagrams

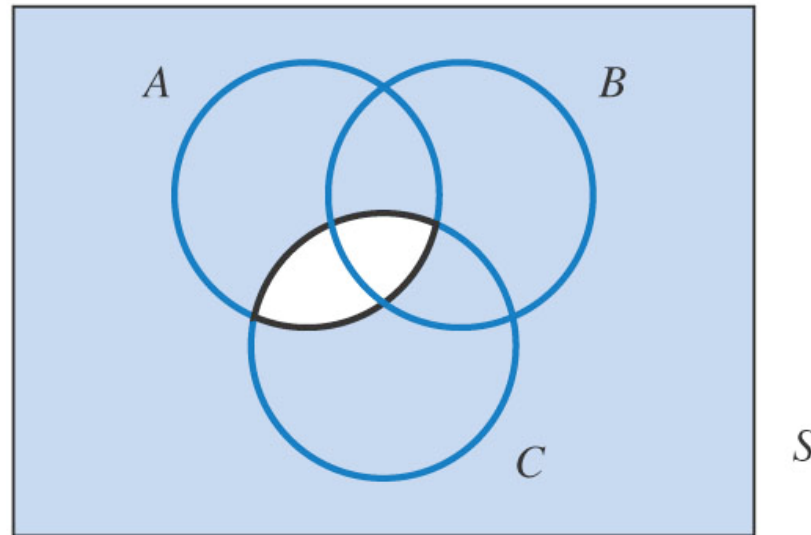


Which formula describes the blue region?

- A. $A \cap C$
- B. $A' \cup C'$
- C. $(A \cap B \cap C)'$
- D. $(A \cap B) \cap C$

Get your i-clickers

Venn diagrams



Which formula describes the blue region?

A. $A \cap C$

B. $A' \cup C'$

C. $(A \cap B \cap C)'$

D. $(A \cap B) \cap C$

Definitions of Probability

Two definitions of probability

- (1) **STATISTICAL PROBABILITY**: the relative frequency with which an event occurs in the long run
- (2) **INDUCTIVE PROBABILITY**: the degree of belief which it is reasonable to place in a proposition on given evidence

Statistical Probability

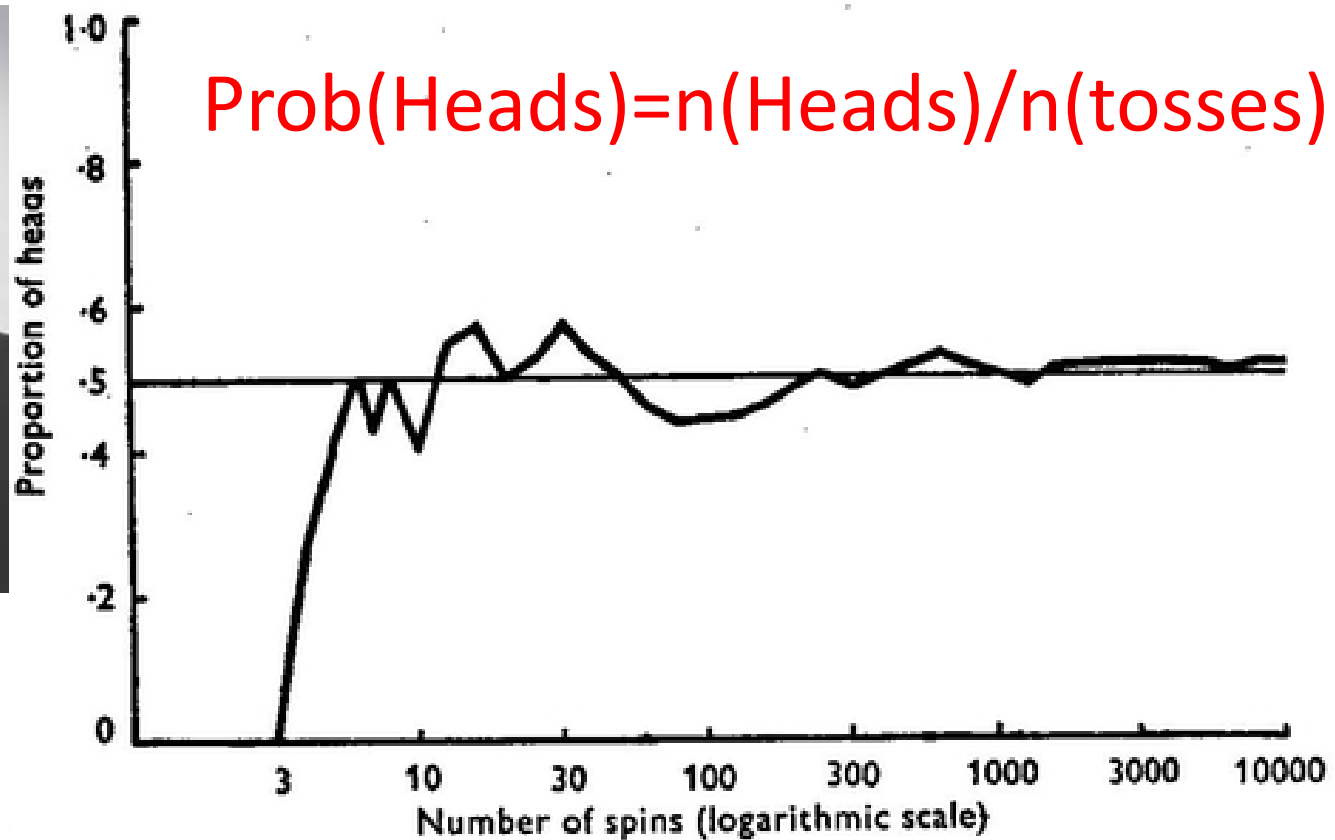
A **statistical probability** of an event is the **limiting value** of the **relative frequency** with it occurs in a **very large number** of **independent trials**

Empirical

Statistical Probability of Coin Toss



John Edmund Kerrich
(1903–1985)
British/South African
mathematician



Proportion of heads among 10,000 coin tosses (Kerrich 1946)

Matlab is easy to learn

- Matlab is the lingua franca of all of **engineering**
- We are working with CS department to reinsta teaching Matlab in CS 101
- Meanwhile, use online tutorials e.g.:
<https://www.youtube.com/watch?v=82TGgQApFIQ>
- **Matlab** is designed to work with **Matrices** → symbols ***** and **/** are understood as **matrix multiplication** and **division**
- Use **.*** and **./** for regular (non-matrix) multiplication
- Add **;** in the end of the line to avoid displaying the output on the screen
- **Loops**: `for i=1:100; f(i)=2.*rand; end;`
- **Conditional statements**: `if rand>0.5; count=count+1; end;`
- **Plotting**: `plot(x,y,'ko-')`; or `semilogx(x,y,'ko-')`; or `loglog(x,y,'ko-')`; .
To keep **adding plots onto the same axes** use: `hold on;`
To **create a new axes** use `figure;`
- **Generating matrices**: `rand(100)` – generates square matrix 100x100.
Confusing! Use `rand(100,1)` or `zeros(30,20)`, or `randn(1,40)`;
- If Matlab complains multiplying matrices **check sizes** using `whos` and if needed **use transpose** operation: `x=x'`;

A Matlab Cheat-sheet (MIT 18.06, Fall 2007)

Basics:

save 'file.mat' save variables to *file.mat*
 load 'file.mat' load variables from *file.mat*
 diary on record input/output to file *diary*
 diary off stop recording
 whos list all variables currently defined
 clear delete/undefine all variables
 help command quick help on a given *command*
 doc command extensive help on a given *command*

Defining/changing variables:

$x = 3$ define variable x to be 3
 $x = [1 \ 2 \ 3]$ set x to the 1×3 row-vector (1,2,3)
 $x = [1 \ 2 \ 3];$ same, but don't echo x to output
 $x = [1; 2; 3]$ set x to the 3×1 column-vector (1,2,3)
 $A = [1 \ 2 \ 3 \ 4; 5 \ 6 \ 7 \ 8; 9 \ 10 \ 11 \ 12];$
 set A to the 3×4 matrix with rows 1,2,3,4 etc.
 $x(2) = 7$ change x from (1,2,3) to (1,7,3)
 $A(2,1) = 0$ change $A_{2,1}$ from 5 to 0

Arithmetic and functions of numbers:

$3*4$, $7+4$, $2-6$ $8/3$ multiply, add, subtract, and divide numbers
 3^7 , $3^{(8+2i)}$ compute 3 to the 7th power, or 3 to the $8+2i$ power
 $\text{sqrt}(-5)$ compute the square root of -5
 $\text{exp}(12)$ compute e^{12}
 $\log(3)$, $\log_{10}(100)$ compute the natural log (ln) and base-10 log (\log_{10})
 $\text{abs}(-5)$ compute the absolute value $|-5|$
 $\sin(5\pi/3)$ compute the sine of $5\pi/3$
 $\text{besselj}(2,6)$ compute the Bessel function $J_2(6)$

Arithmetic and functions of vectors and matrices:

$x * 3$ multiply every element of x by 3
 $x + 2$ add 2 to every element of x
 $x + y$ element-wise addition of two vectors x and y
 $A * y$ product of a matrix A and a vector y
 $A * B$ product of two matrices A and B
 $x * y$ not allowed if x and y are two column vectors!
 $x .* y$ element-wise product of vectors x and y
 A^3 the square matrix A to the 3rd power
 x^3 not allowed if x is not a square matrix!
 $x.^3$ every element of x is taken to the 3rd power
 $\cos(x)$ the cosine of every element of x
 $\text{abs}(A)$ the absolute value of every element of A
 $\text{exp}(A)$ e to the power of every element of A
 $\text{sqrt}(A)$ the square root of every element of A
 $\text{expm}(A)$ the matrix exponential e^A
 $\text{sqrtm}(A)$ the matrix whose square is A

Transposes and dot products:

$x.', A.'$ the transposes of x and A
 x', A' the complex-conjugate of the transposes of x and A
 $x' * y$ the dot (inner) product of two column vectors x and y
 $\text{dot}(x, y)$, $\text{sum}(x.*y)$...two other ways to write the dot product
 $x * y'$ the outer product of two column vectors x and y

Constructing a few simple matrices:

$\text{rand}(12,4)$ a 12×4 matrix with uniform random numbers in $[0,1)$
 $\text{randn}(12,4)$ a 12×4 matrix with Gaussian random (center 0, variance 1)
 $\text{zeros}(12,4)$ a 12×4 matrix of zeros
 $\text{ones}(12,4)$ a 12×4 matrix of ones
 $\text{eye}(5)$ a 5×5 identity matrix I ("eye")
 $\text{eye}(12,4)$ a 12×4 matrix whose first 4 rows are the 4×4 identity
 $\text{linspace}(1.2, 4.7, 100)$
 row vector of 100 equally-spaced numbers from 1.2 to 4.7
 $7:15$ row vector of 7,8,9,...,14,15
 $\text{diag}(x)$ matrix whose diagonal is the entries of x (and other elements = 0)

Portions of matrices and vectors:

$x(2:12)$ the 2nd to the 12th elements of x
 $x(2:\text{end})$ the 2nd to the last elements of x
 $x(1:3:\text{end})$ every third element of x , from 1st to the last
 $x(:)$ all the elements of x
 $A(5,:)$ the row vector of every element in the 5th row of A
 $A(5,1:3)$ the row vector of the first 3 elements in the 5th row of A
 $A(:,2)$ the column vector of every element in the 2nd column of A
 $\text{diag}(A)$ column vector of the diagonal elements of A

Solving linear equations:

$A \setminus b$ for A a matrix and b a column vector, the solution x to $Ax=b$
 $\text{inv}(A)$ the inverse matrix A^{-1}
 $[L,U,P] = \text{lu}(A)$ the LU factorization $PA=LU$
 $\text{eig}(A)$ the eigenvalues of A
 $[V,D] = \text{eig}(A)$ the columns of V are the eigenvectors of A , and
 the diagonals $\text{diag}(D)$ are the eigenvalues of A

Plotting:

$\text{plot}(y)$ plot y as the y axis, with 1,2,3,... as the x axis
 $\text{plot}(x,y)$ plot y versus x (must have same length)
 $\text{plot}(x,A)$ plot columns of A versus x (must have same # rows)
 $\text{loglog}(x,y)$ plot y versus x on a log-log scale
 $\text{semilogx}(x,y)$ plot y versus x with x on a log scale
 $\text{semilogy}(x,y)$ plot y versus x with y on a log scale
 $\text{fplot}(@(\text{expression}), [a,b])$
 plot some expression in x from $x=a$ to $x=b$
 axis equal force the x and y axes of the current plot to be scaled equally
 $\text{title}('A \text{ Title}')$ add a title $A \text{ Title}$ at the top of the plot
 $\text{xlabel}('blah')$ label the x axis as *blah*
 $\text{ylabel}('blah')$ label the y axis as *blah*
 $\text{legend}('foo', 'bar')$ label 2 curves in the plot *foo* and *bar*
 grid include a grid in the plot
 figure open up a new figure window

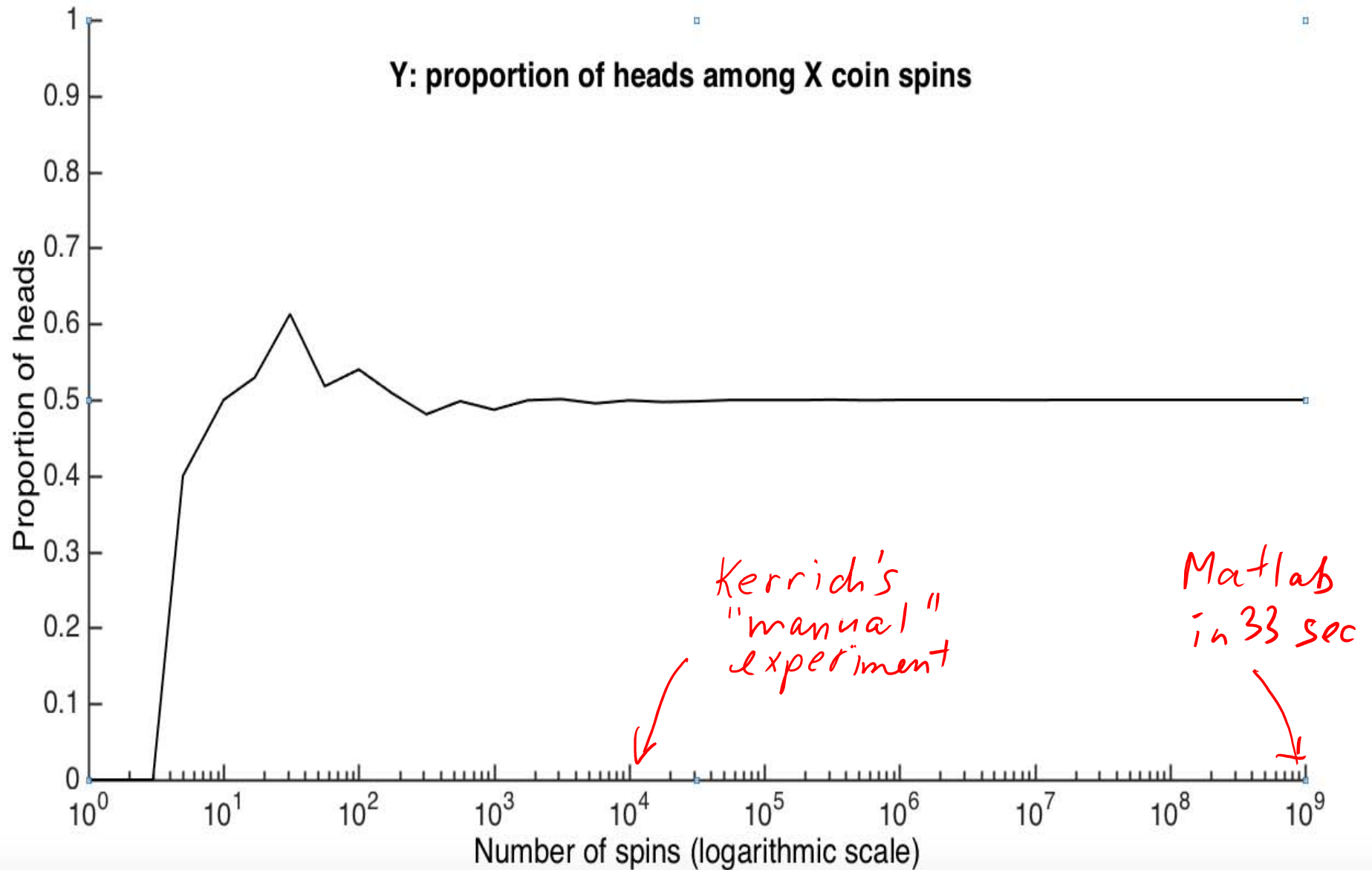
<http://web.mit.edu/18.06/www/Spring09/matlab-cheatsheet.pdf>

Matlab group exercise

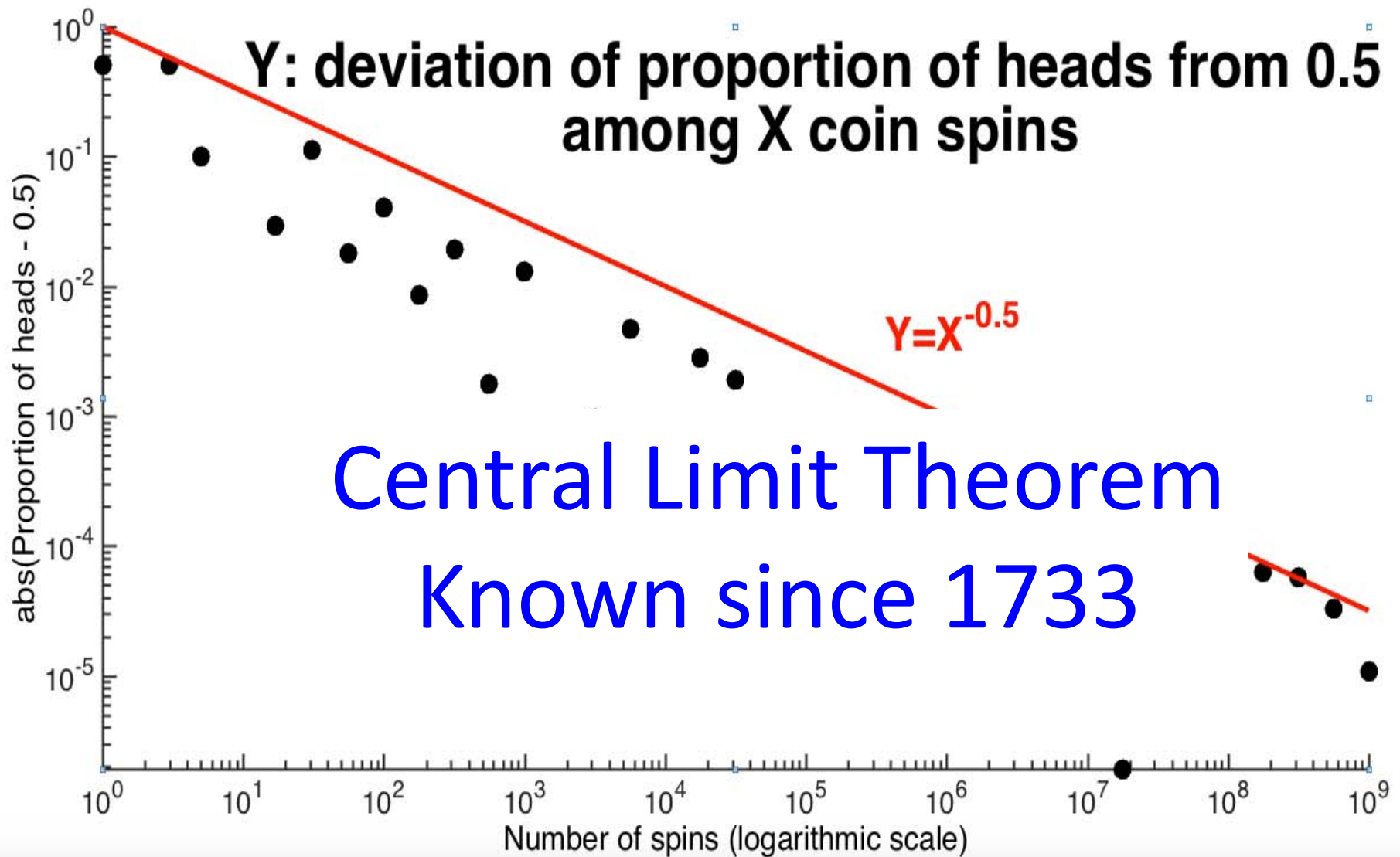
- I want each one of you to write a matlab script to:
 - Simulate a fair coin toss experiment. For this Use `floor(2.*rand)` to generate a fair coin toss:
1 – heads, 0 - tails
 - Calculate the fraction of heads `f_heads(t)` in
t=10;
1000;
100,000;
10,000,000 throws
 - Plot `f_heads(t)` vs t with a logarithmic t-axis
 - Plot `abs(f_heads(t)-0.5)` vs t on a log-log plot
(both axes are logarithmic)

How I did it

- Stats=1e7;
- r0=rand(Stats,1);
- r1=floor(2.*r0);
- n_heads(1)=r1(1);
- for t=2:Stats; n_heads(t)=n_heads(t-1)+r1(t); end;
- tp=[1, 10,100,1000, 10000, 100000, 1000000, 10000000]
- np=n_heads(tp)
- fp=np./tp
- figure; semilogx(tp,fp,'ko-');
- hold on; semilogx([1,10000000],[0.5,0.5],'r--');
- figure; loglog(tp,abs(fp-0.5),'ko-');
- hold on; loglog(tp,0.5./sqrt(tp),'r--');



Proportion of heads among 1,000,000,000 coin tosses
(10⁵ more than Kerrich) took me 33 seconds on my Surface Book



ABS(Proportion of heads-0.5)
among 100,000,000 coin tosses

