

Random Number Generation

PSEUDO-RANDOM NUMBER GENERATION

Some function when queried gives a number in $(0, 1)$

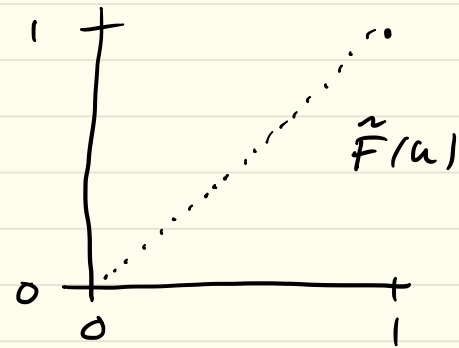
Desired Properties

1. UNIFORM

- many many many samples should look like coming from $U(0, 1)$

e.g. the "empirical CDF"

$$\tilde{F}(u) = \# \text{ samples } \leq u$$



INDEPENDENCE

- ? How would you test?

- more to come


PSEUDO-RANDOM NUMBER GENERATION (RNG)

IMPORTANT CONSIDERATIONS

- FAST
- PORTABLE BETWEEN COMPUTERS
- REPLICABLE
- CLOSELY APPROXIMATE IDEAL STATISTICAL PROPERTIES OF UNIFORMITY AND INDEPENDENCE
- "LONG CYCLE"

RNGs based on functions produce a sequence of numbers

$x_0 \ x_1 \ x_2 \ \dots \ x_j \ \dots \ x_k \ \dots \ x_{m-1}$


period of cycle

◦ "cryptographically secure"

- Given long sequence of past samples,
very hard to guess next sample

Common RNGs

Linear Congruential Method (LCM)

produces sequence x_0, x_1, x_2, \dots by

$$x_{i+1} = (ax_i + c) \bmod m \quad i=0, 1, 2, \dots$$

multiplier

increment

modulus

Values chosen for a , c , m , and x_0 have ENORMOUS impact on statistical properties and cycle length

Random numbers are in $[0, m-1]$, converted by

$$R_i = \frac{x_i}{m}$$

Tausworthe Generator

- an option in Möebius

operates on sequence of binary digits
 $b_0 b_1 \dots b_k \dots$

numbers formed by grouping bits, e.g. every 64
to get 64 bit unsigned integers

$$b_n = c_{q-1} \cdot b_{n-1} \oplus c_{q-2} \cdot b_{n-2} \oplus \dots \oplus c_0 \cdot b_{n-q}$$

c_i are binary coefficients
· is logical AND

Need good choices for c_i

COMBINED LINEAR CONGRUENTIAL GENERATOR

- technique to create longer period generator

- Use k different LCM

j^{th} generator has m_j prime, a_j multiplier

- $X_{i,j}$ is i^{th} output from generator j

$$X_{i,j} \sim U[0, m_j - 1]$$

$$Y_i = \left(\sum_{j=1}^k (-1)^{j-1} X_{i,j} \right) \pmod{m_1 - 1} \quad R_i = \begin{cases} X_i / m_1, & X_i > 0 \\ (m_1 - 1) / m_1, & X_i = 0 \end{cases}$$

maximum period

$$P = \frac{(m_1 - 1)(m_2 - 1) \cdots (m_k - 1)}{2^{k-1}}$$

CRYPTOGRAPHICALLY SECURE RNG

CRNG should pass the "next bit" test

Given the 1st k bits of a random sequence, no polynomial time algorithm can predict the next bit with probability of success $> 50\%$.

CRNG should pass "state compromise extensions"

- state of LCM is last number generated
- state of Mersenne Twister is 624
64-bit integers

IN EVENT THAT SOME OR ALL OF STATE IS COMPROMISED,
SHOULD BE IMPOSSIBLE TO RECONSTRUCT STREAM
OF RANDOM NUMBERS PRIOR TO REVOLATION

Random Number Streams

seed : starting point in sequence

$x_0 \ x_1 \ x_2 \ \dots \ x_j \ \dots \ x_{m-1} \ x_0 \ x_1 \ \dots$



start here and $R_1 = \frac{x_j}{m}$, $R_2 = \frac{x_{j+1}}{m}$ etc.

YOU CAN run multiple streams from the same mathematical sequence, separated by b

$x_0 \ x_1 \ \dots \ x_{b-1} \ x_b \ x_{b+1} \ \dots \ x_{2b-1} \ x_{2b} \ x_{2b+1} \ \dots$



$\frac{x_0}{m}, \frac{x_1}{m}, \dots$



$\frac{x_b}{m}, \frac{x_{b+1}}{m}, \dots$



$\frac{x_{2b}}{m}, \frac{x_{2b+1}}{m}, \dots$

? For really large period b , how do you FIND these starting points...
points???

TESTS FOR RANDOM NUMBERS

Reminder of how statistical testing works...

$S = \{S_1, S_2, \dots, S_n\}$ a set of DATA SAMPLES, from some (unknown) probability distribution

We wish to look for statistical support for some **hypothesis** about the data

H_0 :

usually called **null hypothesis**

H_A :

usually called **alternative hypothesis**

- Construct some statistic with known distribution using S
- ASK "What's the probability of seeing this statistic if H_0 is true?"

"CONFIDENCE LEVEL" specifies decision

IF probability of observing statistic $< \alpha$
WE **REJECT THE NULL HYPOTHESIS**
 α typically .1, .05, .01

LOGIC

" IF H_0 is true, then there is only a
(10%, 5%, 1%) chance of data set S
yielding up that statistic. So probably H_0
is not true"

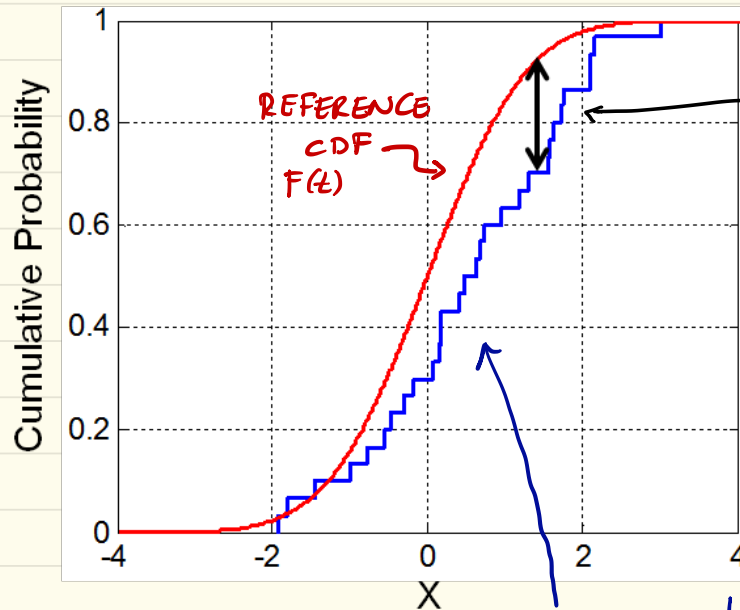
called $1-\alpha$ confidence level

" WITH 99% confidence we reject H_0 "

TEST FOR STATISTICAL UNIFORMITY

One test is "one sample" Kolmogorov-Smirnov Test

- compare empirical CDF with that of the "reference distribution"



$$D = \max_t |B(t) - F(t)|$$

$$B(t) = \frac{|\{s_i \in S \mid s_i \leq t\}|}{|S|}$$

NOTICE — IS LARGEST DEVIATION

TEST FOR STATISTICAL UNIFORMITY

USE K-S TABLE (look up) example

$n \backslash \alpha$	0.01	0.05	0.1
2	0.929	0.842	0.776
4	0.733	0.624	0.564
10	0.490	0.410	0.368
20	0.356	0.294	0.264
50	0.230	0.190	0.170
OVER 50	$\frac{1.63}{\sqrt{n}}$	$\frac{1.36}{\sqrt{n}}$	$\frac{1.22}{\sqrt{n}}$

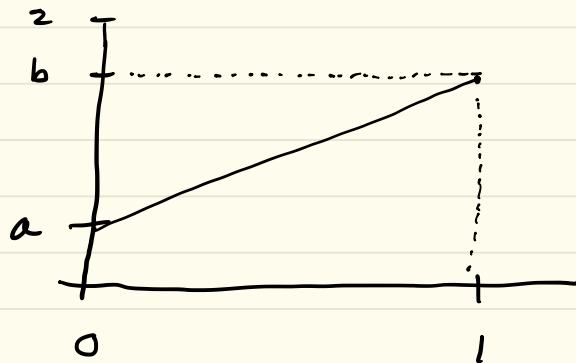
n is # samples

How to use :

1. Choose α for $1-\alpha$ confidence, e.g. 95%
2. Compute empirical distribution function, measure D
3. Find statistic threshold $d(n, \alpha)$ given n & α
4. Reject H_0 if $D > d(n, \alpha)$

Example

consider family of distributions "similar" to U



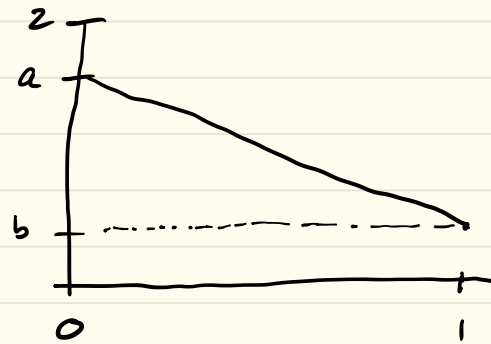
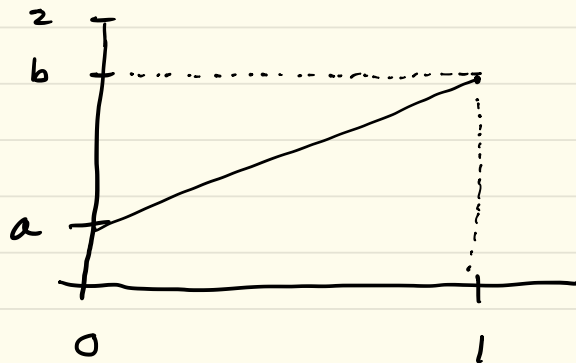
$$f(t) = \frac{a + (b-a) \cdot t}{N(a,b)}$$

$$N(a,b) = \text{area under curve} = \min\{a,b\} + |b-a|/2.0$$

Notice that $b = 2.0 - a$, and that $U(0,1)$ has $a = b = 1.0$

Example

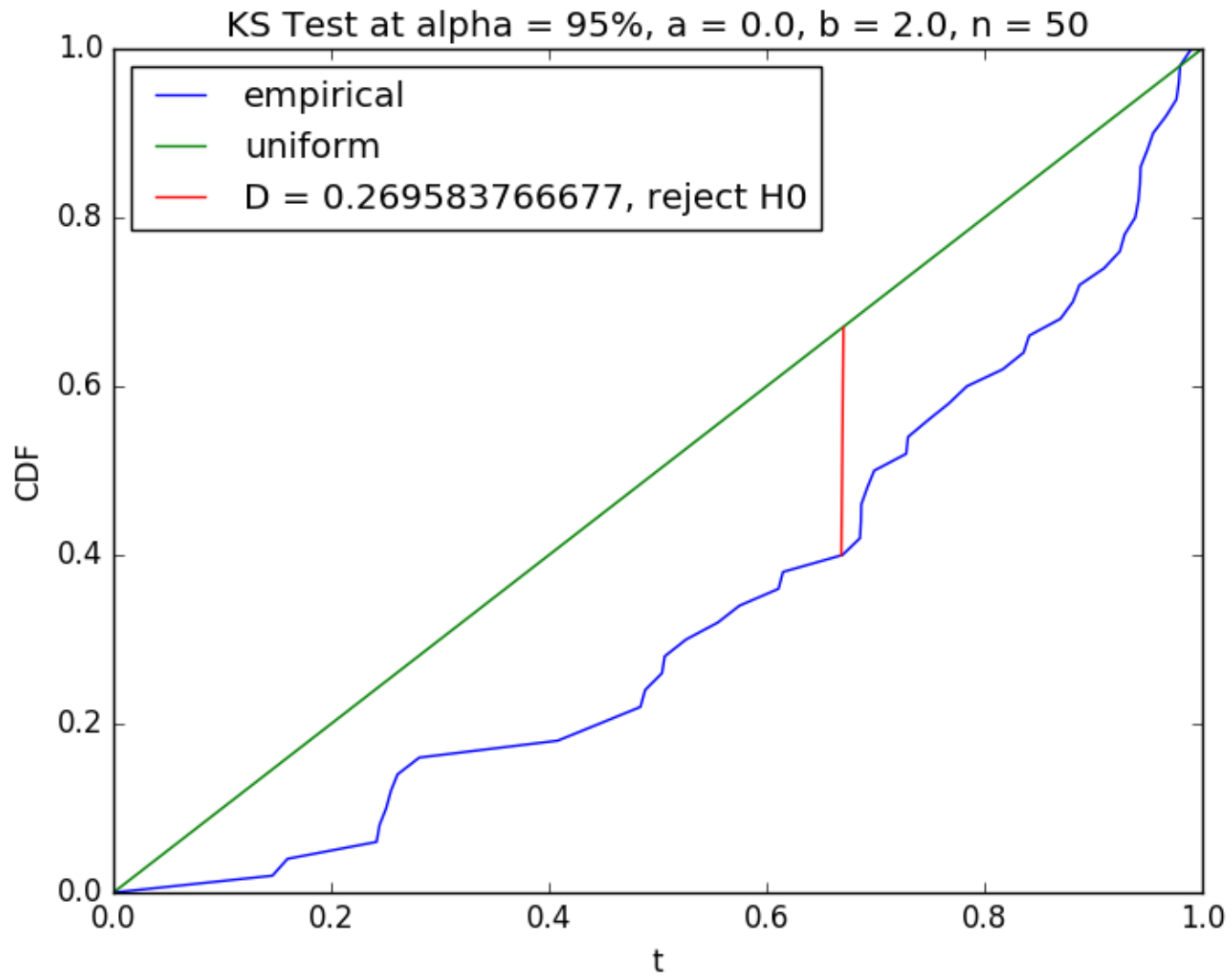
consider family of distributions "similar" to U

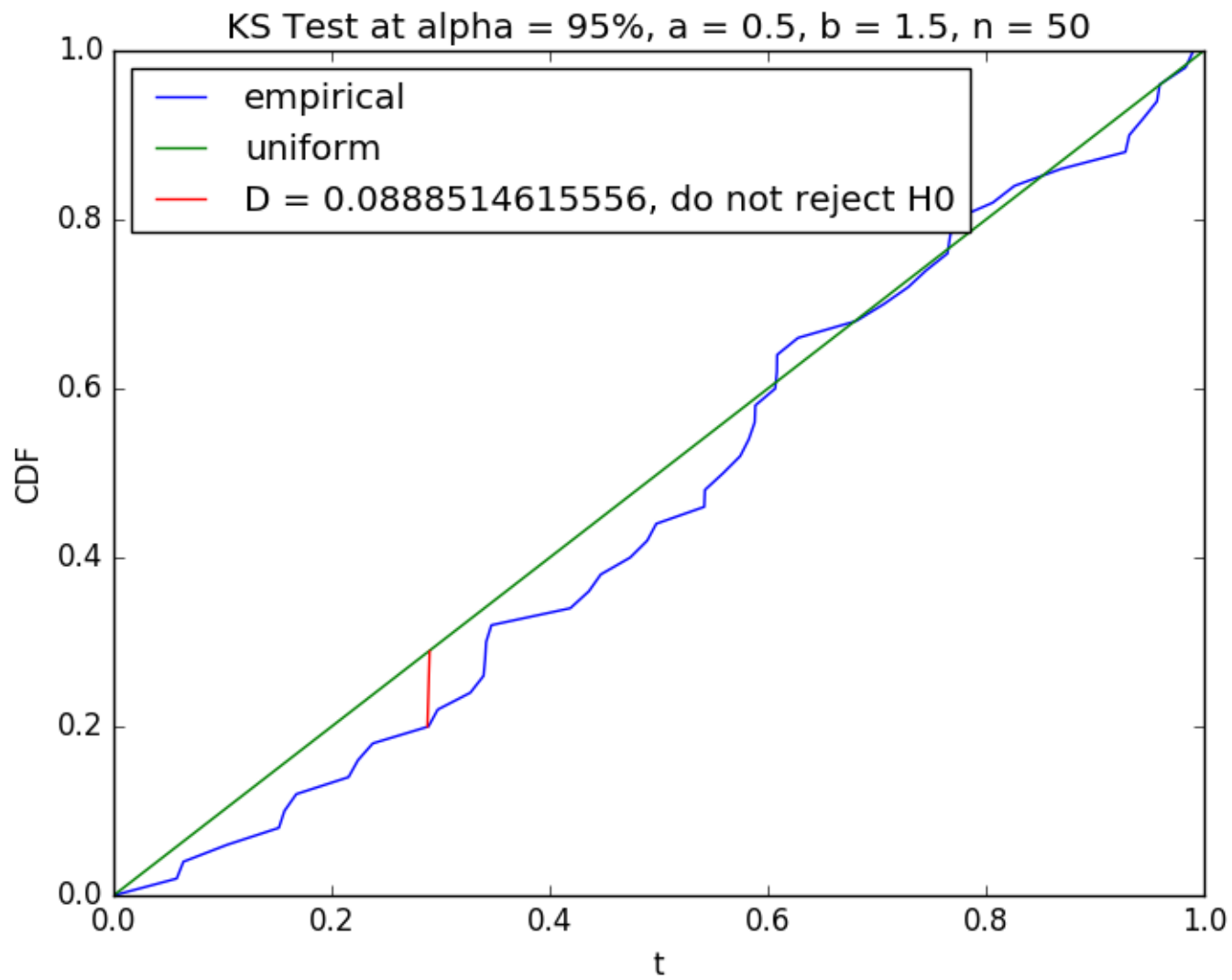


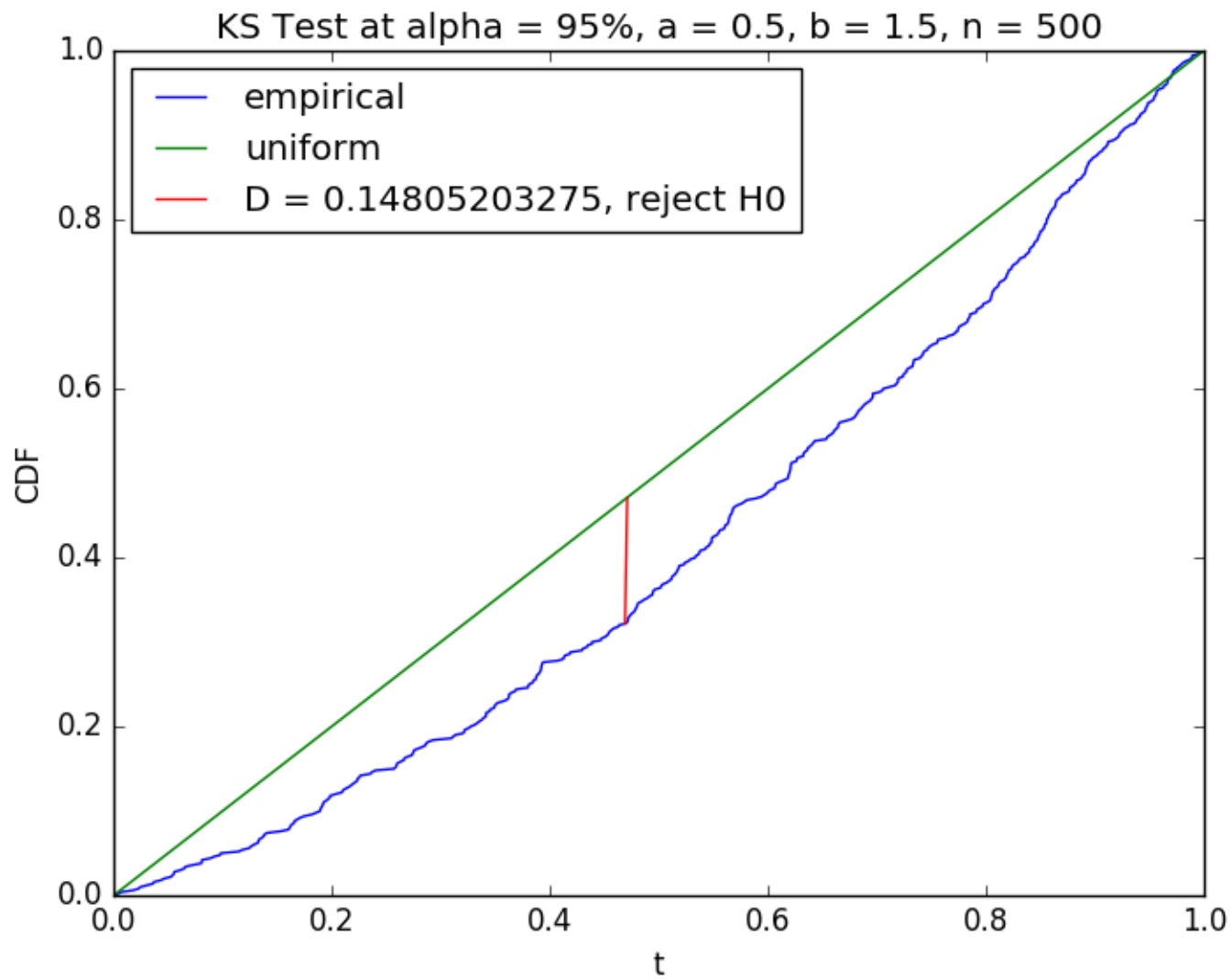
$$f(t) = \frac{a + (b-a) \cdot t}{N(a,b)}$$

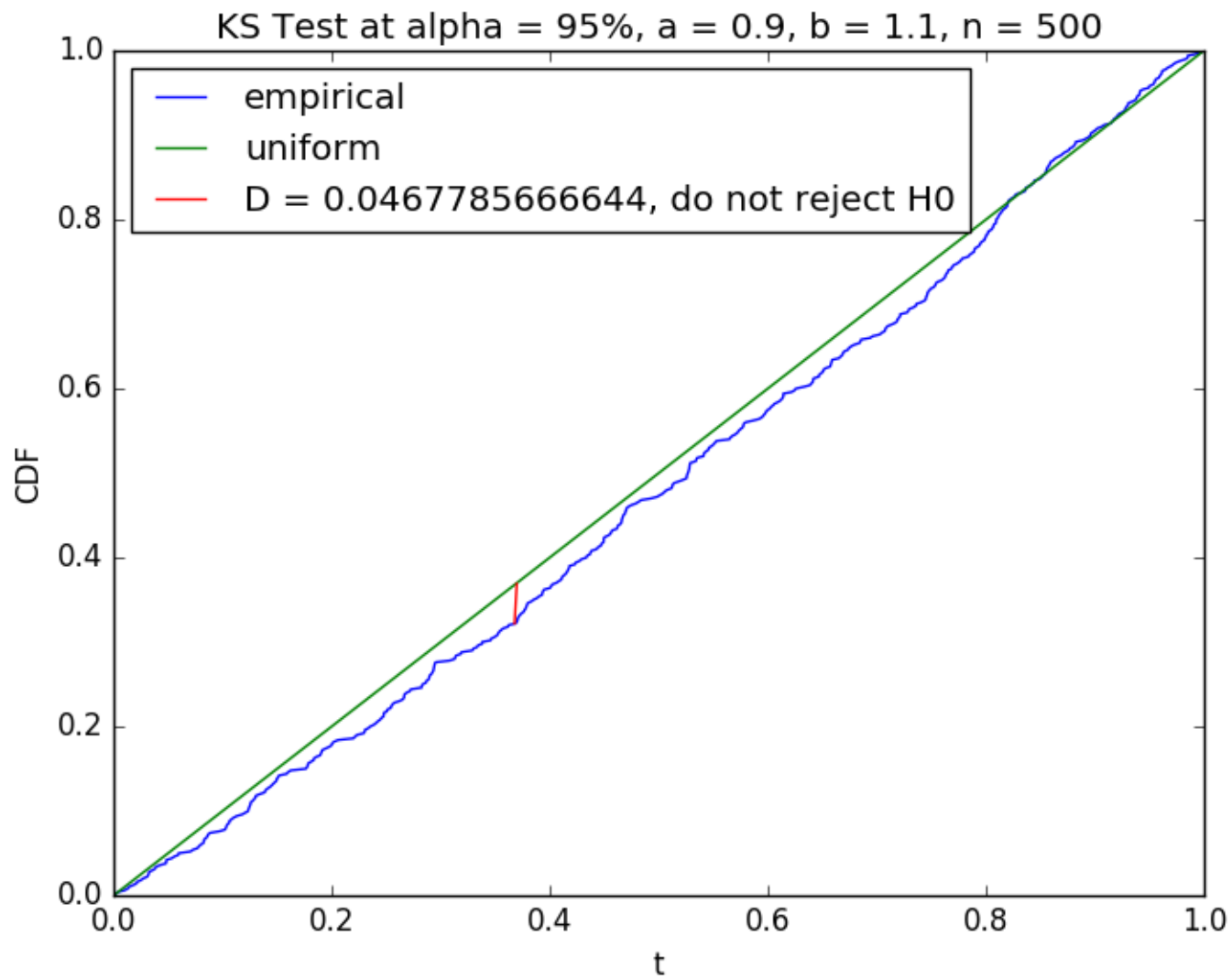
$$N(a,b) = \text{area under curve} = \min\{a,b\} + |b-a|/2.0$$

Notice that $b = 2.0 - a$, and that $U(0,1)$ has $a = b = 1.0$









Chi-square test

- create n equi-length partitions

$$\left[0, \frac{1}{n}\right), \left[\frac{1}{n}, \frac{2}{n}\right) \dots \left[\frac{n-1}{n}, 1\right)$$

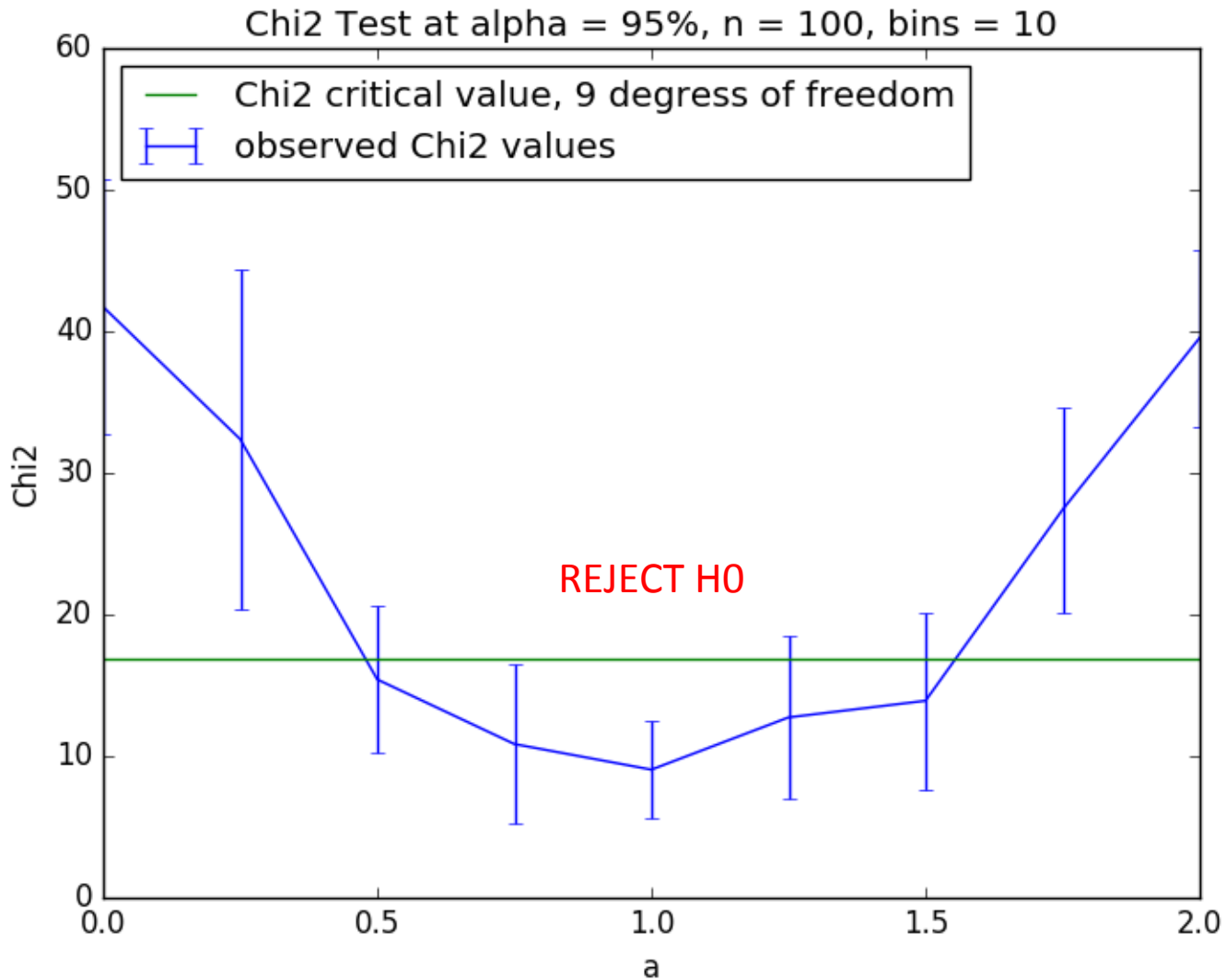
- Bin each $s_i \in \{s_1, s_2, \dots, s_N\}$
- Under the assumption of uniformity, bins will get approximately the same number of samples, N/n
- Let O_i be the OBSERVED number of samples in

$$\left[\frac{i-1}{n}, \frac{i}{n}\right)$$

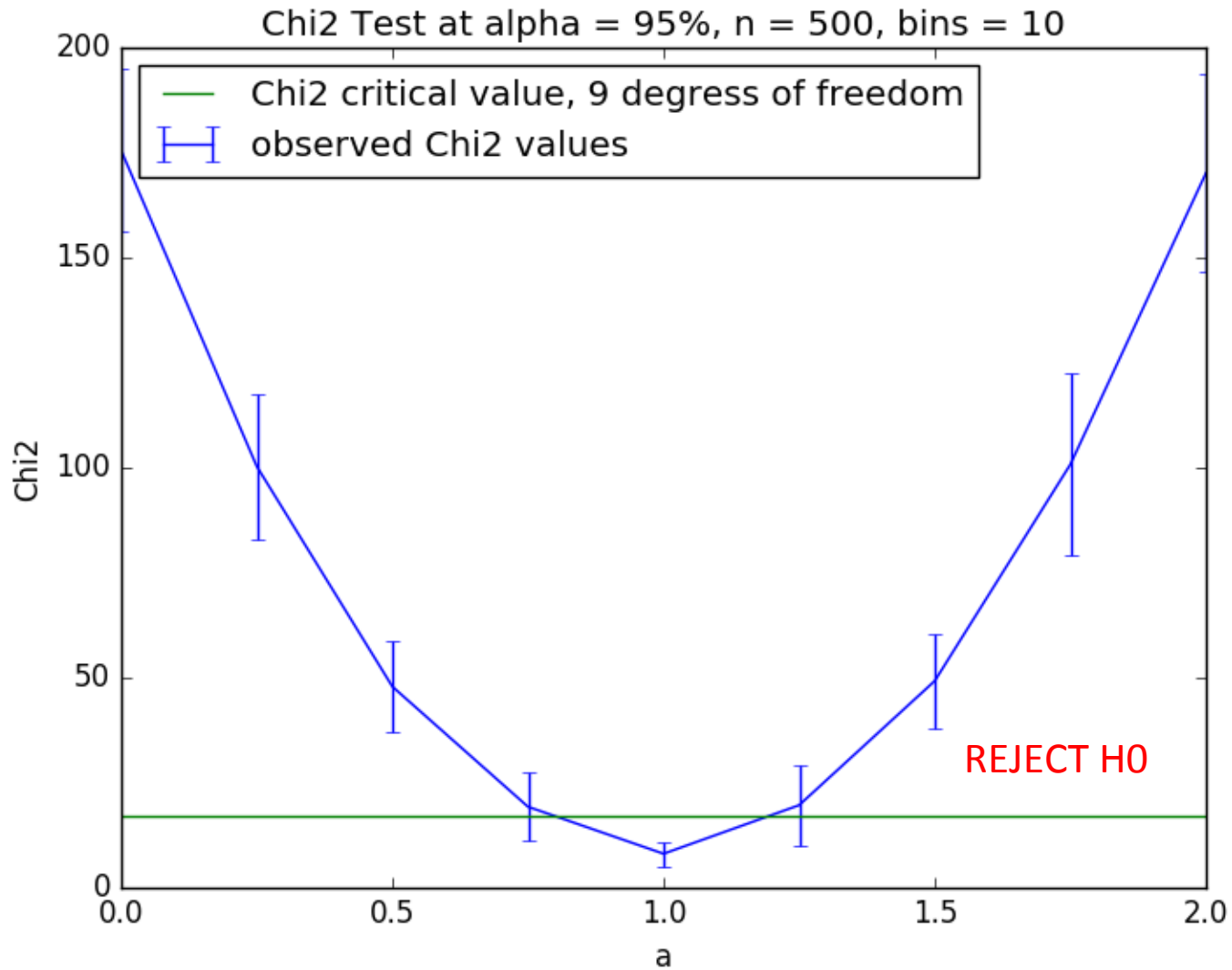
- compute
$$\chi_o^2 = \sum_{i=1}^n \frac{(O_i - N/n)^2}{N/n} \quad N \geq 50$$

- said to have $n-1$ *degrees of freedom*

- compare against critical value of χ_o^2 distribution
- "TOO LARGE" taken to mean uniformity is unlikely



a is parameter from distribution studied using KS test
a = 1 is perfectly uniform

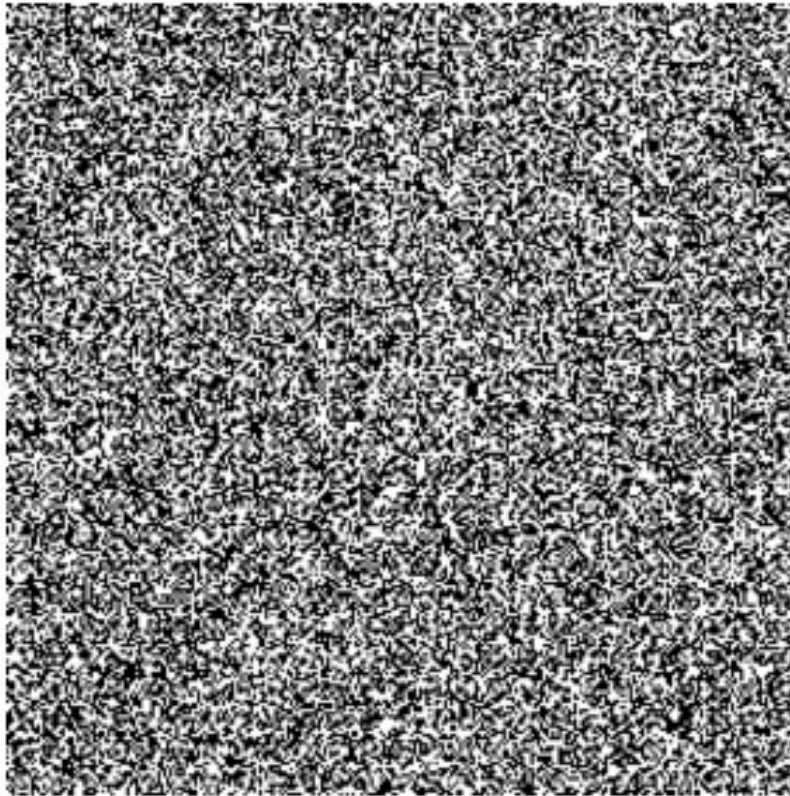


a is parameter from distribution studied using KS test
a = 1 is perfectly uniform

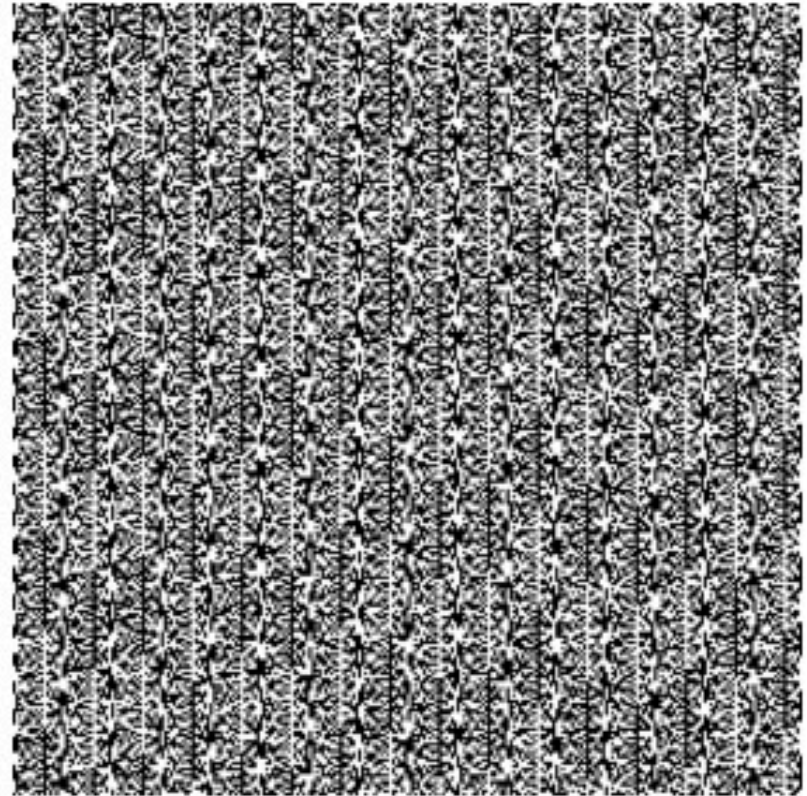
How can you tell whether your RNG produces independent numbers?

DILBERT By SCOTT ADAMS





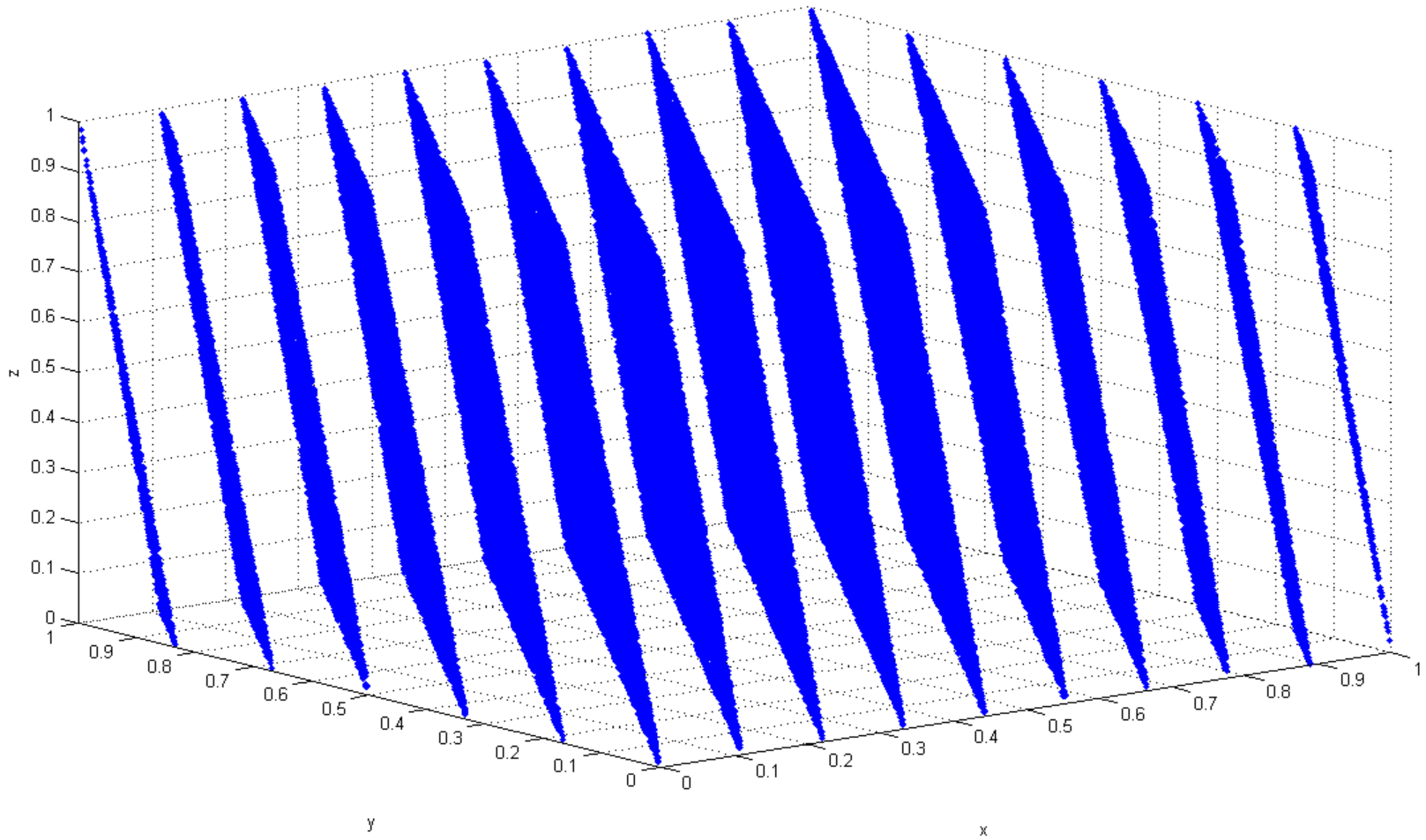
RANDOM.ORG



PHP rand() on Microsoft Windows

Test: plot (u_n, u_{n+1})

Patterns suggest lack of independence



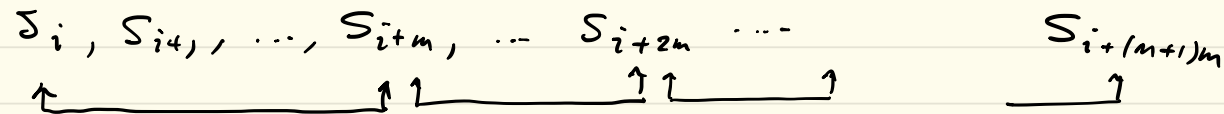
Plot of (u_{n-1}, u_n, u_{n+1}) using (old) rand() generation from early Unix

TEST FOR INDEPENDENCE

Given sequence of generated random numbers

$$S_1, S_2, \dots, S_k, \dots, S_N \quad \text{with } \bar{S} = \frac{1}{N} \sum_{i=1}^N S_i$$

We can test for **autocorrelation** between every m numbers

$$S_i, S_{i+m}, \dots, S_{i+2m}, \dots, S_{i+(m+1)m}$$


k -lag autocorrelation

$$\tilde{\gamma}_k = \frac{\sum_{i=1}^{N-k} (S_i - \bar{S})(S_{i+k} - \bar{S})}{\sum_{i=1}^N (S_i - \bar{S})^2}$$

When $\{S_i, S_{i+m}, S_{i+2m}, \dots\}$ are independent and set size is large, then

$\tilde{\gamma}_k$ is approximately normal, mean = 0, $\sigma = 1$

TESTING FOR RANDOMNESS, TYPICALLY ONLY LAG-1 TEST

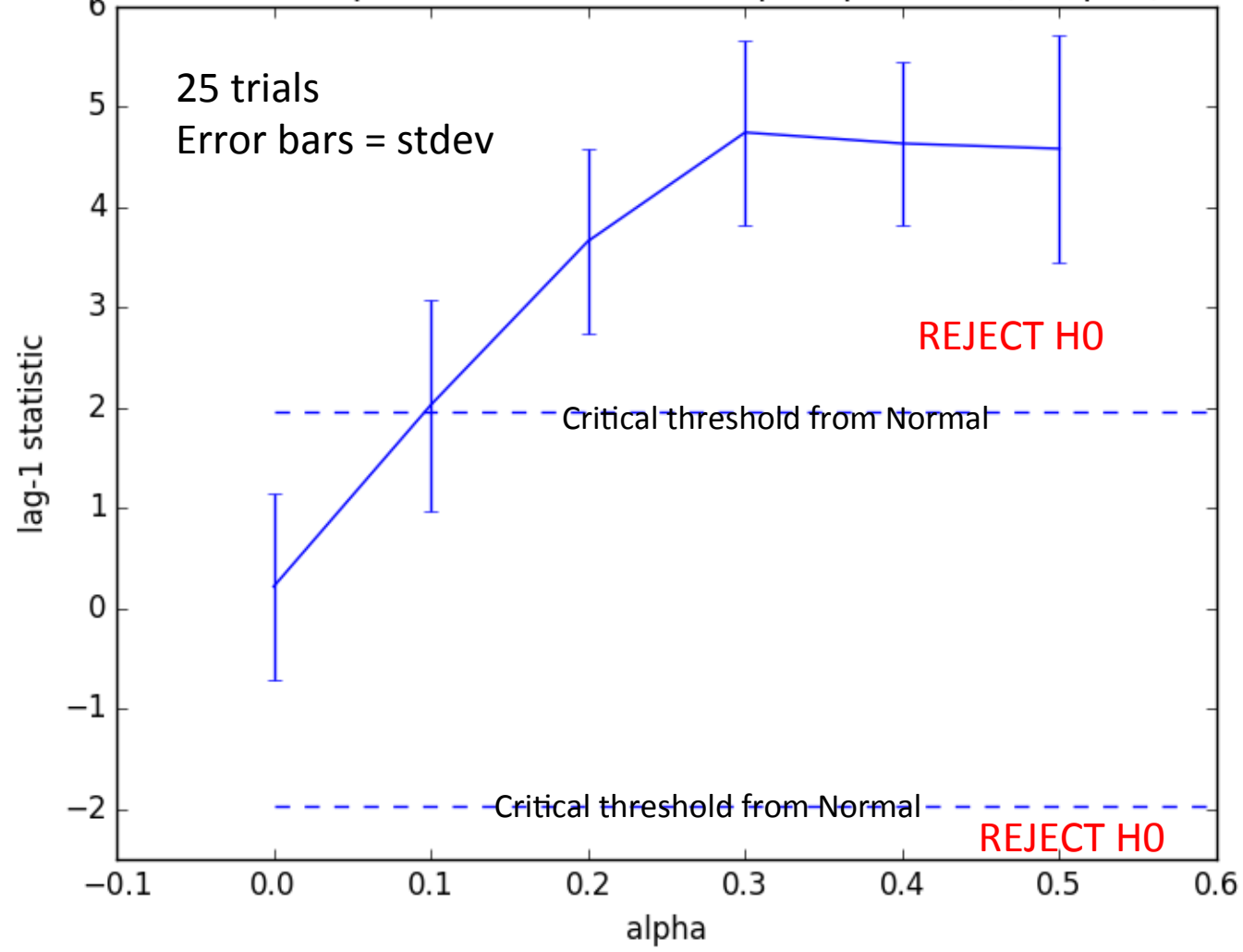
Subject $\tilde{\delta}_1$ to a "two tail" test



a positive value above $t_{\alpha/2}$ or a negative value below $-t_{\alpha/2}$ is evidence of non-independence

look up $t_{\alpha/2}$ values from standard normal table.

1 Correlation Test at alpha = 95%, newU = alpha*prevU + (1-alpha)*U, 10K sar



$$U_n = \alpha * U_{n-1} + (1-\alpha) * U$$