

Introduction to data science & artificial intelligence (INF7100)

Arthur Charpentier

#132 Uncertainty and Randomness

été 2020

Uncertainty



Data



Storing Data: Tally sticks, used starting in the Paleolithic area



A tally (or tally stick) was an ancient memory aid device used to record and document numbers, quantities, or even messages.

Data

Collecting Data: John Graunt conducted a statistical analysis to curb the spread of the plague, in Europe, in 1663

Data

Data Manipulation: Herman Hollerith created a Tabulating Machine that uses punch cards to reduce the workload of US Census, in 1881, see [1880 Census](#), $n = 50$ million Americans.

1	06/20/189	Frank Ford	W M 40		1	Salomon
2	—	Carolina	W F 30	Wife	1	Hopkinson
3	—	Gustavus	W D 5	Son 1	1	St. Steel
4	—	Anna	W F 4	daughter 1	1	
5	—	Richard	W M 2	Son 2	1	
6	70/115	Hannibal Peter	W D 35		1	Julia Evans
7	—	Constance	W F 25	Wife	1	Hopkinson
8	—	Emmale	W F 8	daughter 1	1	St. Steel
9	—	Johel	W m 6	Son 2	1	St. Steel
10	—	Emilia	W F 4	daughter 1	1	
11	—	Nina	W F 2	daughter 1	1	
12	—	Willy	W D 5	son 1	1	

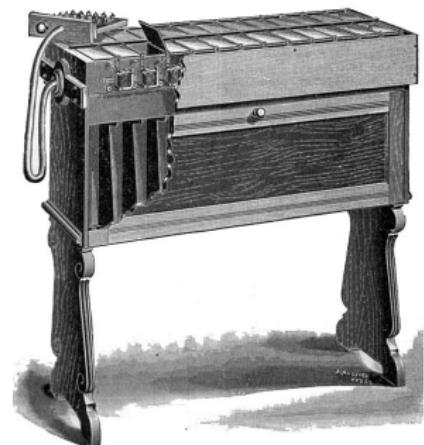


Fig. 3.—Sorting Machine.

Hollerith's Electric Sorting and Tabulating Machine.

Allen-Scott Report

Data Center Plan Called Privacy Invasion

By ROBERT S. ALLEN
and PAUL SCOTT

WASHINGTON — A special White House task force is recommending the creation of a federal data center which eventually could have a comprehensive file on every man, woman and child in the country.

Now under study in inner administration circles, the still-secret report advocates the gradual transfer of all governmental records and statistics to magnetic computer tape, which would be turned over to a newly-created agency that would function as a general data center.

The computerized information would be available, at the push of a button, to a wide range of government authorities.

Estimated cost of the pro-

cial Security, census data, medical, credit and criminal reports.

"Comprehensive information of this kind, centralized in one agency," says Gallagher, "could constitute a highly dangerous dossier bank. Such an agency would be a distinct departure from our American tradition."

Subcommittee investigators have ascertained that the task force's report states that a vast accumulation of government records already is on computer tape and could be turned over to the proposed general data center immediately. Listed as among these available files are:

Internal Revenue Service — 742 million personal and corporate tax returns.

Defense Department — 14

the most intimate information, the investigators learned, are freely passed around among agencies. Graphically illustrative of this practice and its harsh consequences are the following two instances:

A teenager visiting Washington stayed with an uncle, at his mother's suggestion. During the night the boy was sexually assaulted by the uncle. Years later, as a Phi Beta Kappa graduate from a leading Eastern university, the boy applied for a job with the National Security Agency. During a required lie detector test he told about the assault. His frank admission cost him the desired job.

But that wasn't all. This affair, in which he was an innocent victim, haunted him again

Data Center: The US Government plans the world's first data center to store 742 million tax returns and 175 million sets of fingerprints, in 1965.

Data

Literary Digest Poll based on 2.4 million readers, Özbeyp (2018)

A. Landon: 57% vs. F.D. Roosevelt: 43%

George Gallup sample of about 50,000 people

A. Landon: 44% vs. F.D. Roosevelt: 56%

Actual results

A. Landon: 38% vs. F.D. Roosevelt: 62%



The Literary Digest NEW YORK OCTOBER 31, 1936

Topics of the day

LANDON, 1,293,669; ROOSEVELT, 972,897

Final Returns in The Digest's Poll of Ten Million Voters

With the great battle of the bellies in the Poll of ten million voters, scattered throughout the forty-eight States of the Union, now finished, and in the table below we record the returns received up to the hour going to press.

These figures are exactly as received from more than one in every five voters polled in the poll, and have neither weighted, adjusted nor interpreted.

Never before in an experience covering

more than a quarter of a century in taking polls have we received such a varied variety of criticism—praise from many; condemnation from many others—and yet it has been just the same type that has come to us time after time as the Poll has been taken in all these years.

A telegram from a newspaper in California asks: "Is it true that Mr. Roosevelt has got more votes than Franklin D. Roosevelt?" A telephone message only the day before these lines were written: "Has the Repub-

ican National Committee purchased Ten Literary Digests?" And all types and varieties, including: "Have the Jews purchased Ten Literary Digests?" "Is the Pope of Rome a member of the Literary Digest?" And so it goes—all equally absurd and amusing. We could add more to this list, and yet all these questions in their day and age have been asked by us who have been experiencing all down the years from the very first Poll.

Problem. Now, are the figures in this Poll correct? Is there any question we could simply refer to a telegram we sent to a young man in Massachusetts the other day in answer to his challenge to us to wager our last ducat on the accuracy of our Poll. We wired him as follows:

"For nearly a quarter century, we have been taking Polls of the voters in the forty-eight States of the Union for almost half a century, and we have always merely mailed the ballots, counted and recorded those

returned and let the people of the Nation draw their conclusions as to our accuracy. So far, we have been right in every Poll. Will we be right in the current Poll? That, however, is a question which the President's reelection, is in the 'lap of the gods.'

"We never make any claims before election but we respectfully refer you to the record of the accuracy of our Literary Digest." And so it goes—all equally absurd and amusing.

"We could add more to this list, and yet all these questions in their day and age have been asked by us who have been experiencing all down the years from the very first Poll.

Problem. Now, are the figures in this Poll correct? Is there any question we could simply refer to a telegram we sent to a young man in Massachusetts the other day in answer to his challenge to us to wager our last ducat on the accuracy of our Poll. We wired him as follows:

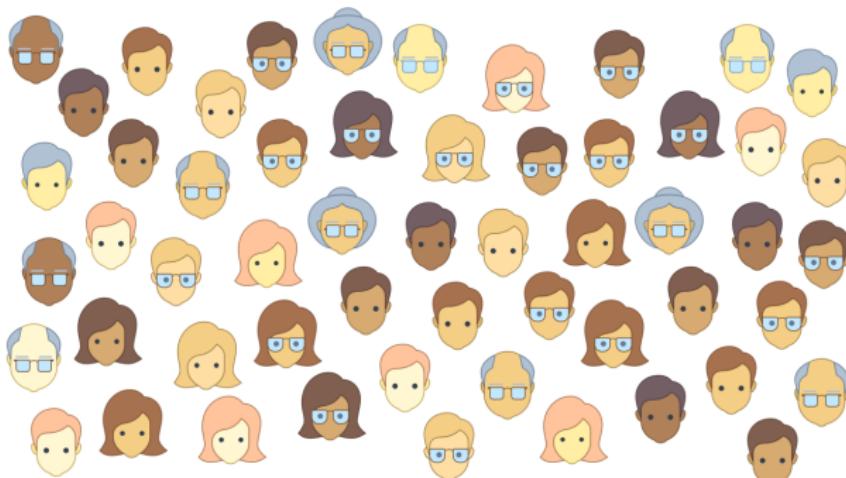
"The statistics used in this article in this article are property of Frank J. Wagner Company and have been copyrighted by it; neither the original author nor the copyright owner may be reproduced or published without the special permission of the copyright owner."

In studying the table of the voters from

The statistics used in this article in this article are property of Frank J. Wagner Company and have been copyrighted by it; neither the original author nor the copyright owner may be reproduced or published without the special permission of the copyright owner."

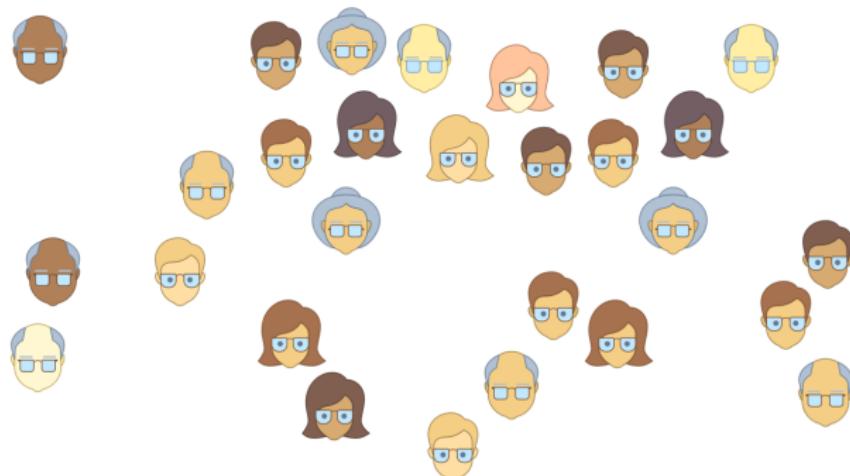
1936 Election & Sampling

A badly chosen big sample is much worse than a well-chosen small sample, see [Why the 1936 Literary Digest Poll Failed](#) or [How biases can arise in sampling](#)

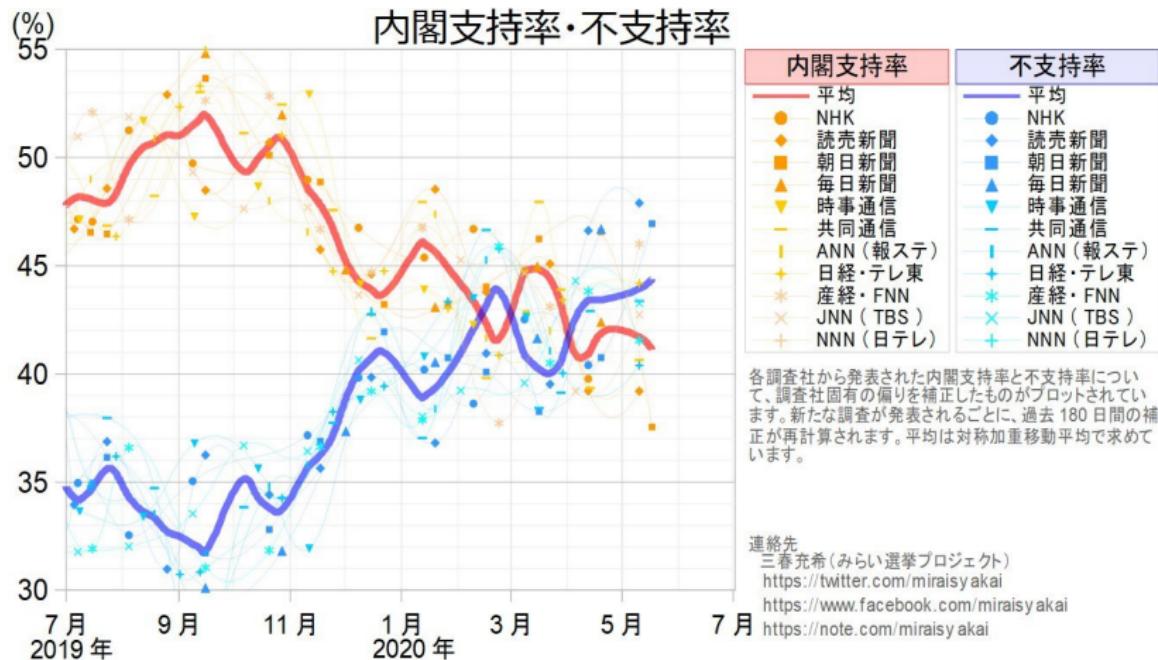


1936 Election & Sampling

A badly chosen big sample is much worse than a well-chosen small sample, see [Why the 1936 Literary Digest Poll Failed](#) or [How biases can arise in sampling](#)



Sampling Uncertainty



“Dans toute statistique, l'inexactitude du nombre est compensée par la précision des décimales” Alfred Sauvy (see #331)

Sampling Error and Uncertainty

Sample $\mathcal{X} = \{x_1, x_2, x_3, \dots, x_{n-1}, x_n\}$

Consider some statistics $t(\mathcal{X})$, e.g. $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

What if we had had another sample \mathcal{X}' ?

How different $t(\mathcal{X}')$ could have been...?

Under the assumption that x_i 's are observed from some $\mathcal{N}(\mu, \sigma^2)$ distribution, we know (mathematical theory) that $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$

- ▶ we need to know the distribution (here Gaussian)
- ▶ we need to know the parameters (here σ^2)

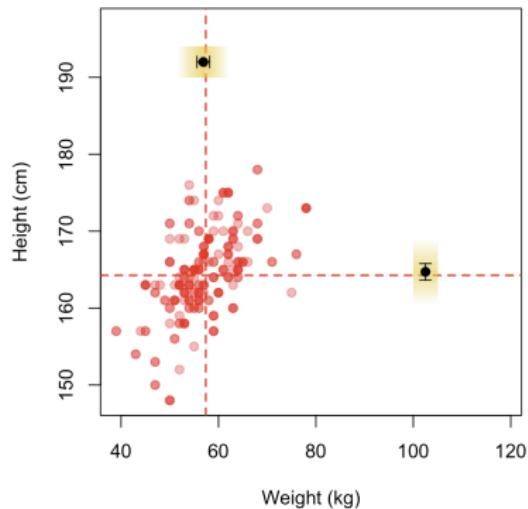
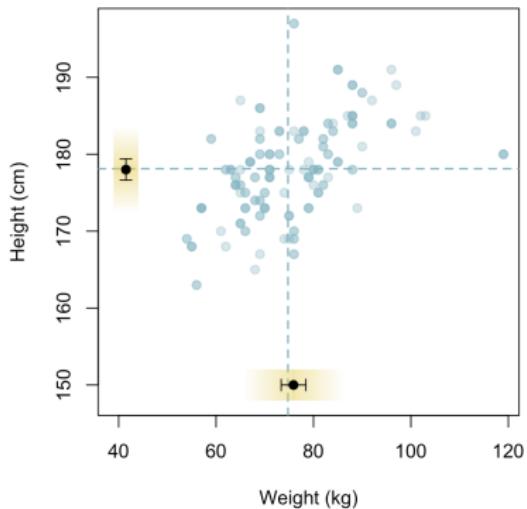
Sampling Error and Uncertainty

We can create some fake sample by **resampling**...

1	19	10	1	8	12	12	5	9	20	6	10	18	8	20	13	16	16	13	16	11	11	7	14	17	1	19	9	10	6	1	11	8	19	8
2	4	16	6	17	20	9	5	13	16	7	17	12	10	12	3	17	12	12	11	7	11	9	20	14	8	16	3	2	5	19	11	3	10	18
3	18	1	1	2	5	14	8	5	18	18	20	3	4	7	10	16	17	8	10	7	3	8	3	1	12	3	3	10	18	3	12	11	7	2
4	9	3	16	1	11	20	3	20	16	15	18	10	15	13	8	5	4	13	20	16	14	8	14	4	2	13	15	8	18	7	14	9	2	20
5	3	16	12	18	15	8	17	8	15	10	2	5	6	13	4	1	16	10	5	13	13	18	18	6	11	2	12	2	15	4	9	20	20	1
6	17	10	15	14	8	5	9	15	19	4	5	6	6	7	6	9	17	14	9	20	6	2	19	13	5	6	11	18	4	14	14	1	8	16
7	18	8	10	19	5	18	11	1	17	11	4	20	4	7	7	6	15	4	7	11	4	17	16	18	14	1	8	8	10	10	12	10	16	6
8	11	1	5	3	11	1	20	15	5	3	19	9	16	15	4	20	19	4	11	7	3	17	6	19	10	12	4	15	16	9	1	14	3	8
9	11	2	2	9	20	5	14	7	6	7	4	7	10	4	13	10	9	16	4	8	20	14	8	17	2	10	2	11	16	19	6	16	20	4
10	9	10	11	4	9	19	16	15	18	5	16	7	17	12	15	10	4	4	1	3	19	8	19	9	18	18	18	7	17	4	2	16	7	
11	9	13	16	13	14	19	2	4	17	4	12	6	7	5	8	16	15	8	17	17	5	18	2	4	6	20	9	5	15	3	6	11	15	18
12	7	13	18	19	3	16	4	6	14	6	8	19	13	5	20	10	2	15	9	17	15	16	20	1	13	8	5	17	13	2	9	12	4	18
13	12	17	12	14	11	8	15	10	3	11	9	16	18	6	13	9	15	13	2	1	8	19	9	1	18	16	4	12	16	8	20	12	13	8
14	9	15	13	14	15	6	14	8	12	11	4	20	13	9	14	9	8	2	20	16	7	4	14	11	5	17	18	6	19	2	11	16	17	9
15	19	20	1	2	11	14	18	15	16	3	12	9	14	3	7	13	10	6	16	3	19	2	1	3	13	9	6	20	20	13	17	18	20	
16	7	4	8	18	4	14	3	17	11	3	12	20	7	9	2	11	12	2	16	2	19	10	9	6	13	1	13	7	1	7	14	6	4	20
17	2	16	7	1	6	14	20	6	19	1	8	12	12	14	7	5	9	17	11	9	7	7	14	6	11	5	15	5	5	11	16	13	5	9
18	6	20	1	14	3	15	13	11	5	12	18	5	7	16	18	8	17	1	18	4	4	19	8	20	12	13	7	7	6	2	10	10	16	
19	9	17	14	3	10	7	8	11	19	12	10	8	4	9	1	3	2	18	3	17	20	13	9	9	13	9	1	7	11	20	20	1	17	3
20	8	18	4	9	19	3	10	16	7	5	10	1	14	16	9	17	9	7	2	13	6	5	1	19	19	1	1	17	7	17	13	16	12	19

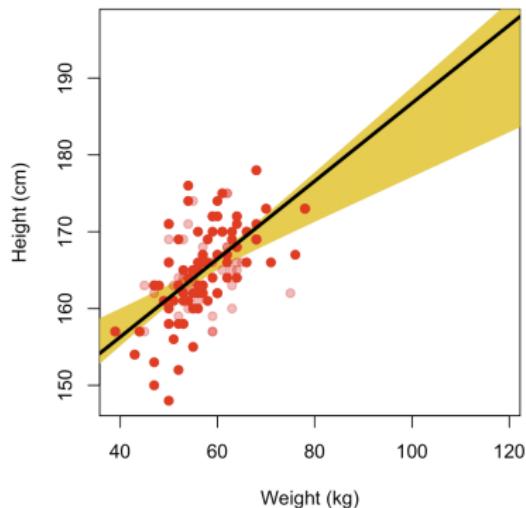
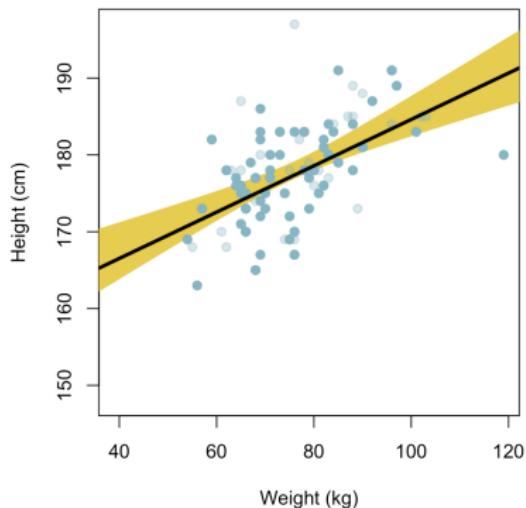
Resampling & Bootstrapping

Students height and weight, of boys (left) and girls (right),



Resampling & Bootstrapping

Regression line of height against weight



Bootstrap & Bootstrapping

Algorithm 1: Bootstrapping

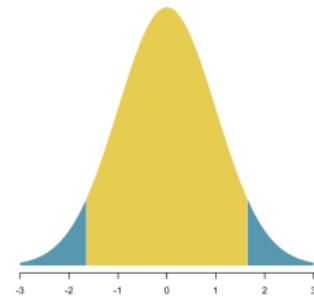
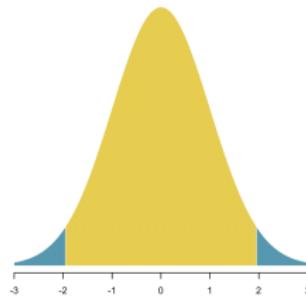
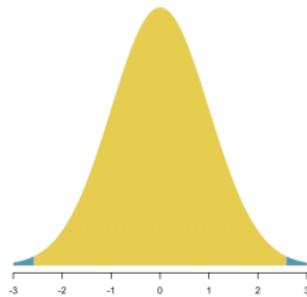
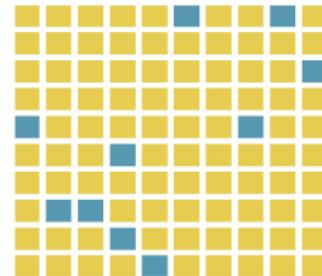
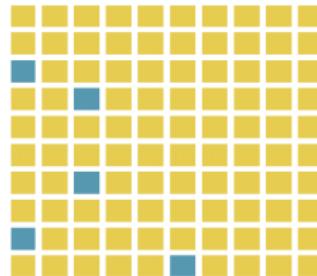
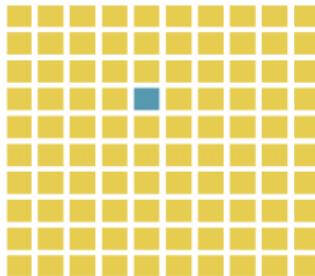
```
1 initialization :  $x = \{x_1, \dots, x_n\}$ ;  
2 for  $b = 1, 2, \dots, B$  do  
3    $x_b \leftarrow$  resample from  $x$  (with replacement, same size);  
4    $t_b = T(x_b)$ 
```

In python and R

```
1 > from sklearn.utils import resample  
2 > x = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6]  
3 > xb = resample(x, replace=True)  
4 > xb  
5 [0.5, 0.3, 0.4, 0.5, 0.6, 0.5]
```

```
1 > x = c(0.1, 0.2, 0.3, 0.4, 0.5, 0.6)  
2 > xb = sample(x, replace=TRUE)  
3 > xb  
4 [1] 0.4 0.4 0.5 0.1 0.5 0.6
```

1%, 5% or 10% chance



(Parametric) Probability Distribution

Is that the only way...?

What if we know the distribution of X_i 's ?

We have here [for the Poisson distribution] a special case of the remarkable fact that there exist a few distributions of great universality which occur in a surprisingly great variety of problems. The three principal distributions, with ramifications throughout probability theory, are the binomial distribution, the normal distribution...and the Poisson distribution, William Feller, mentioned in Computational Models, Paul Humphreys, (see 431, ... 434)



The Monte Carlo method is a statistical sampling technique that over the years has been applied to a wide variety of scientific problems. Although the name suggests that it must involve random sampling, it has given ever more sophisticated descriptions of the method as captured in some original papers. One quote is in 1951 about schools:

"The first thoughts and attempts I made at Monte Carlo sampling were suggested by a question which occurred to me in 1944 while I was consulting the literature and planning solutions. The question was what are the odds of getting a California license plate laid out with 52 cards will come out randomly? I thought it would be a good time trying to estimate them by pure continental calculations. I wondered whether a more practical method than "abstract fiddling" might not be to lay out 52 cards and then simply observe and count the number of times the sequence of 52 cards is possible to emerge with the beginning and the end of the new set of 52 cards. This was the first time I had heard of problems of neutron diffusion and other applications of Monte Carlo and other generalizations of the Monte Carlo method. I began to think about how to change processes described by the 52 cards into an equivalent form interpretable in terms of the 52 cards. I then read a paper by Ulam and von Neumann, and we began to work on the problem."

"The first thoughts and attempts I made at Monte Carlo sampling were suggested by a question which occurred to me in 1944 while I was consulting the literature and planning solutions. The question was what are the odds of getting a California license plate laid out with 52 cards will come out randomly? I thought it would be a good time trying to estimate them by pure continental calculations. I wondered whether a more practical method than "abstract fiddling" might not be to lay out 52 cards and then simply observe and count the number of times the sequence of 52 cards is possible to emerge with the beginning and the end of the new set of 52 cards. This was the first time I had heard of problems of neutron diffusion and other generalizations of the Monte Carlo method. I began to think about how to change processes described by the 52 cards into an equivalent form interpretable in terms of the 52 cards. I then read a paper by Ulam and von Neumann, and we began to work on the problem."

Von Neumann was intrigued. Statistical sampling was already well known.

10

Generating Random Numbers

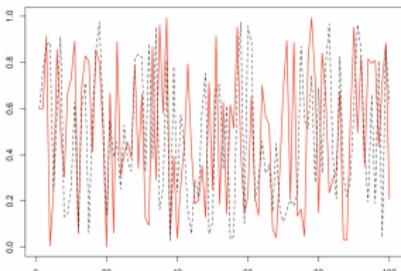
predicting without understanding
or
understanding without predicting

Quelle responsabilité pour les algorithmes ?

Sequences generated with Sedgewick algorithm, $u_n = x_n/m$ where

$$x_n = (ax_{n-1} + c) \text{ modulo } m$$

$a = 1013904223$, $c = 1664525$ and $m = 2^{32}$
with $u_1 = 0.6$ and $u_1 = 0.60001$.



See [Histoire du hasard et de la simulation](#) for further discussion

Generating Random Numbers

A Million Random Digits with 100,000 Normal Deviates, RAND, 1955

★★★★★ Errors Throughout

Reviewed in the United States on December 30, 2013

They sure don't come up with random numbers like they used to. If you look closely, you will note that every tenth digit or so is just a repeat of the last digit and every hundredth or so is a just the same digit repeated three times. How sloppy!

A sampling of this "work":

Page 36 - Line 6 - 15 characters in should be 5, not 4.

Page 99 - Line 18 - first three characters should be "453" not "345".

Page 145 - Line 2 - 7th and 19th characters transposed.

Page 190 - Whole line of numbers omitted between 6th and 7th lines.

Pages 210 and 211 - Two sections appear quasi-randomized, instead of randomized.

★★★★★ Predictable

Reviewed in the United States on October 30, 2013

It seemed like about 10% of the time I was able to predict which number was next. It was still better than Life of Pi which, aside from being irrational, included no estimations of Pi at all.

★★★★★ Not really random

Reviewed in the United States on September 26, 2012

I bought two copies of this book. I find that the first copy perfectly predicts what the numbers will be in the second copy. I feel cheated.

★★★★★ great yet lacking...

Reviewed in the United States on May 1, 2013

So far, it's great, but what I REALLY want is an audio version that I can listen to at the beach or the gym. Maybe narrated by Morgan Freeman???

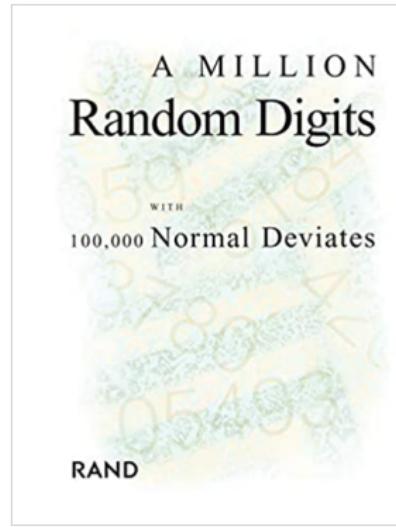


TABLE OF RANDOM DIGITS

00000	10097	32533	76320	13386	34673	54876	08502	09117	39292	74945
00001	37541	04905	64894	74296	24805	24037	20343	10402	00822	91665
00002	00000	00000	00000	00000	00000	00000	00000	00000	00000	00000
00003	99019	03229	05276	70713	38311	31160	88676	74397	04436	27659
00004	12807	89997	00157	36147	64032	36453	98951	16877	12171	76833
00005	66068	74717	34072	78930	36697	36170	65812	29868	11199	29170
00006	85269	77802	02051	65692	88665	74818	72623	85247	18622	88732
00007	85269	77802	02051	65692	88665	74818	72623	85247	18622	88732
00008	63573	32135	05325	47048	30553	57548	28468	28700	83491	25624
00009	00000	00000	00000	00000	00000	00000	00000	00000	00000	00000
00010	98530	17767	14803	68007	32109	40538	66979	93423	50000	73986
00011	11860	05431	39808	27732	30725	58248	28402	32775	47981	
00012	83452	99634	06288	99083	13746	70679	18475	40610	68713	77817
00013	00000	00000	00000	00000	00000	00000	00000	00000	00000	00000
00014	99594	67344	87317	64649	91830	88528	93783	61346	23478	34113
00015	65481	17874	17468	00580	58047	76974	73239	57186	40218	16544
00016	80124	30853	17272	08151	45312	23274	31115	78252	14385	53763
00017	00000	00000	00000	00000	00000	00000	00000	00000	00000	00000
00018	69916	25803	66232	25148	34932	87203	70621	13990	94400	56418
00019	00893	20303	14225	68014	46427	56788	96397	78822	54328	14598
00020	91499	14923	68647	27086	46162	83054	94750	88932	37069	20048
00021	00000	00000	00000	00000	00000	00000	00000	00000	00000	00000
00022	44104	81949	00157	47954	32979	28275	57600	400340	42050	82341
00023	13564	73742	11100	02940	12860	74097	96644	89439	28707	25815
00024	63608	46629	14005	34844	40219	52603	30351	77082	20707	31797

Generating Random Numbers

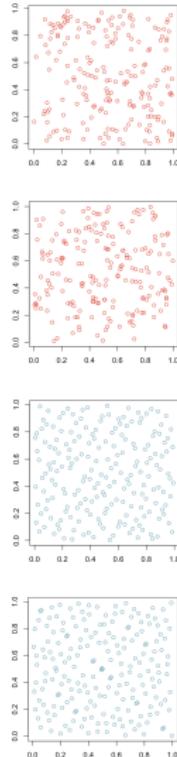
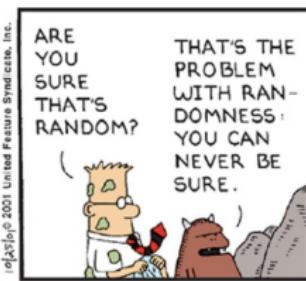
One can easily draw 'random' numbers from any distribution
most popular : uniform on $[0, 1]$ and Gaussian $\mathcal{N}(0, 1)$
In python and R

```
1 > import random
2 > random.uniform(0, 1)
3 0.15473258779393262
4 > import numpy as np
5 > x = np.random.normal(0, 1, 10)
6 > x
7 array([ 1.3136,   0.0829,   1.1754,  -0.7242,   0.5681,
8          0.0701,  -0.0567,   0.3344,  -1.1355,  1.1510])
```

```
1 > runif(10)
2 [1] 0.1118923 0.1714409 0.3878056 0.4111353 0.3819566
3 [6] 0.0145239 0.2615349 0.8561546 0.9593199 0.2472850
4 > rnorm(10)
5 [1] -0.234491 -0.891424 -1.357470 -0.5840280  1.66809
6 [6]  1.085743  0.557528 -1.636832 -0.271538  1.091798
```

Generating Random Numbers

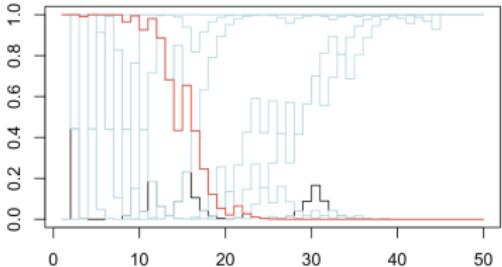
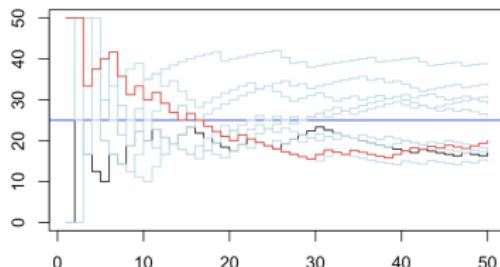
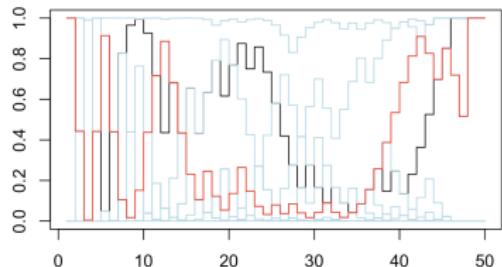
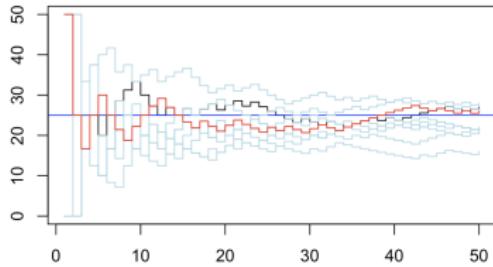
- ▶ draw $n = \text{one million values}$ and plot the histogram to check
- ▶ values should be “independent”
(between n and $n + 1$)



via **Dilbert**

Forecasting Elections

Evolution of $n \mapsto \hat{p}_n = \bar{x}_n$ (left)
and $n \mapsto \mathbb{P}[\bar{X}_N > 1/2 | x_1, \dots, x_n]$ (right)



Forecasting Mortality

Important to add **confidence bands** on predictions

See recently [@Alex_Blanchet](#)'s tweet on mortality in Québec

