# Practical AI & Machine Learning Projects and Datasets



GAN started with seed = 103 (X axis = epoch, Y axis = loss function)

GAN started with seed = 102 (X axis = epoch, Y axis = loss function)

# Preface

This book is intended to participants in the AI and machine learning certification program organized by my AI/ML research lab MLtechniques.com. It is also an invaluable resource to instructors and professors teaching related material, and to their students. If you want to add enterprise-grade projects to your curriculum, with deep technical dive on modern topics, you are welcome to re-use my projects in your classroom. I provide my own solution to each of them.

This book is also useful to prepare for hiring interviews. And for hiring managers, there is plenty of original questions, encouraging candidates to think outside the box, with applications on real data. The amount of Python code accompanying the solutions is considerable, using a vast array of libraries as well as home-made implementations showing the inner workings and improving existing black-box algorithms. By itself, this book constitutes a solid introduction to Python and scientific programming. The code is also on my GitHub repository.

The topics cover generative AI, synthetic data, machine learning optimization, scientific computing with Python, experimental math, synthetic data and functions, data visualizations and animations, time series and spatial processes, NLP and large language models, as well as graph applications and more. It also includes significant advances on some of the most challenging mathematical conjectures, obtained thanks to modern computations. In particular, intriguing new results regarding the Generalized Riemann Hypothesis, and a conjecture regarding record run lengths in the binary digits of $\sqrt{2}$. Most projects are based on real life data, offered with solutions and Python code. Your own solutions would be a great addition to your GitHub portfolio, bringing your career to the next level. Hiring managers, professors, and instructors can use the projects, each one broken down in a number of steps, to differentiate themselves from competitors. Most offer off-the-beaten path material. They may be used as novel exercises, job interview or exam questions, and even research topics for master or PhD theses.

To see how the certification program works, check out our FAQ posted here, or click on the "certification" tab on our website MLtechniques.com. Certifications can easily be displayed on your LinkedIn profile page in the credentials section. Unlike many other programs, there is no exam or meaningless quizzes. Emphasis is on projects with real-life data, enterprise-grade code, efficient methods, and modern applications to build a strong portfolio and grow your career in little time. The guidance to succeed is provided by the founder of the company, one of the top experts in the field, Dr. Vincent Granville. Jargon and unnecessary math are avoided, and simplicity is favored whenever possible. Nevertheless, the material is described as advanced by everyone who looked at it.

The related teaching and technical material (textbooks) can be purchased at MLtechniques.com/shop/. MLtechniques.com, the company offering the certifications, is a private, self-funded AI/ML research lab developing state-of-the-art open source technologies related to synthetic data, generative AI, cybersecurity, geospatial modeling, stochastic processes, chaos modeling, and AI-related statistical optimization.

## About the author

Vincent Granville is a pioneering data scientist and machine learning expert, co-founder of Data Science Central (acquired by TechTarget), founder of MLTechniques.com, former VC-funded executive, author and patent owner.



Vincent's past corporate experience includes Visa, Wells Fargo, eBay, NBC, Microsoft, and CNET. Vincent is also a former post-doc at Cambridge University, and the National Institute of Statistical Sciences (NISS). He published in *Journal of Number Theory*, *Journal of the Royal Statistical Society* (Series B), and *IEEE Transactions on Pattern Analysis and Machine Intelligence*. He is also the author of multiple books, available here. He lives in Washington state, and enjoys doing research on stochastic processes, dynamical systems, experimental math and probabilistic number theory.

# Contents

```
print(data_A)

data_C = data_B.drop(['src_port','size'], axis=1)
data_C = data_C.groupby(data_C.columns.tolist(), as_index=False).size()
data_C1 = data_C[(data_C['bidirectional_mean_ps'] == 60) |
            (data_C['bidirectional_mean_ps'] == 1078) |
            (data_C['size'] > 1)]
data_C2 = data_C[(data_C['bidirectional_mean_ps'] != 60) &
            (data_C['bidirectional_mean_ps'] != 1078) &
            (data_C['size'] == 1)]
print(data_C)
data_C1.to_csv('iot_C1.csv')
data_C2.to_csv('iot_C2.csv')

data_B_full = data_B.join(data.set_index(features), on=features, how='inner')
features.remove('src_port')
data_C1_full = data_C1.merge(data_B_full, how='left', on=features)
data_C2_full = data_C2.merge(data_B_full, how='left', on=features)
data_C1_full.to_csv('iot_C1_full.csv')
data_C2_full.to_csv('iot_C2_full.csv')

map_C1 = data_C1_full.groupby('src_port')['src_port'].count()
map_C2 = data_C2_full.groupby('src_port')['src_port'].count()
map_C1.to_csv('iot_C1_map.csv')
map_C2.to_csv('iot_C2_map.csv')

data_C1 = data_C1_full.drop(['src_port','size_x', 'size_y'], axis=1)
data_C1 = data_C1.groupby(data_C1.columns.tolist(), as_index=False).size()
data_C2 = data_C2_full.drop(['src_port','size_x', 'size_y'], axis=1)
data_C2 = data_C2.groupby(data_C2.columns.tolist(), as_index=False).size()
data_C1.to_csv('iot_C1.csv')
data_C2.to_csv('iot_C2.csv')
```

## 2.3   Good GenAI evaluation, fast LLM search, and real randomness

In this section, I cover several topics in detail. First, I introduce one of the bests random number generators (PRNG) with infinite period. Then, I show how to evaluate the synthesized numbers using the best metrics. Finally, I illustrate how it applies to other contexts, such as large language models (LLM). In particular, the system is based on words with letters from an arbitrary alphabet, and it can be adapted to any prespecified multivariate distribution, not just uniform: the joint ECDF (empirical distribution) attached to a training set in GenAI systems, for instance. At each step, the focus is both on quality and speed, revisiting old methods or inventing new ones, to get solutions performing significantly better and requiring much less computing time. The three components of this system are:

**New powerful random number system**

In its simplest form, the random numbers are defined as the binary digits $d_n = x_n \bmod 2$, from the sequence $x_{n+1} = 3 \cdot (x_n // 2)$, where the double slash is the integer division [Wiki]. It is an improvement over binary digits of quadratic irrationals used previously (see section 4.4 in [12]) in the sense that $x_n$ grows only by a factor $3/2$ at each iteration, rather than 2. All sequences $(x_n)$ that do not grow indefinitely necessarily result in periodic numbers. This is the case for all PRNGs on the market.

In addition, despite having very long periods, these random generators with finite periods exhibit subtle patterns in rather low dimensions: in short, lack of randomness. They can be quite sensitive to the seed and may require many warm-up iterations before reaching higher randomness. See here how you can crack the Mersenne twister used in the Numpy random function. The question is this: how slowly can $x_n$ grow while preserving perfect randomness, fast implementation, and an infinite period? Read on to see how I managed to reduce the aforementioned exponential growth down to linear, while keeping an infinite period.

**Ultrafast, robust evaluation metrics**

The first step is to define what a strongly random sequence is, when it consists of deterministic digits. Details are again in chapter 4 in [12]. The takeaway: you need a metric that captures just that, when testing your system. This is true for all GenAI systems. Indeed, here I am re-using the full multivariate Kolmogorov-Smirnov distance (KS) specifically implemented in the context of synthetic data generation: see section 6.4.2 in [7] for

details. There, I showed how poorly implemented metrics used by vendors fail to capture subtle departures from the target distribution.

In this section, I present a very fast implementation of KS. I also include a few other tests. Very large test batteries exist, for instance Diehard [Wiki]. However, most rely on old statistical practice, offering a large number of disparate, weak tests, rather than a centralized approach to dealing with the problem. You can do a lot better with much fewer tests. This is one of the goals of this project, also with a focus on hard-to-detect patterns.

Also note that the KS distance relies on the CDF rather than the PDF (probability density function). The latter, used in many tests such as Chi-squared, does not work when you have billions of cross-feature buckets in high dimensions, each with very few observations. As in many GenAI systems, this is what we face. To give you an idea, think about counting occurrences of billions of "words" such as

$$3210232010310223034123103323103003110023102$$

in a sequence of trillions of digits in base 4 (in this case, the alphabet has 4 letters). Most counts will be zero. Likewise, the base (that is, the size of the alphabet) may be a very large integer. The KS distance handles this problem transparently by looking at closest strings found in the digit sequences, themselves having only one occurrence most of the time. Also, it easily takes care of conditional probabilities when needed.

My previous KS implementation involved thousands of Pandas SQL queries spanning across many features. The new version discussed here is based on the radix numeration system [Wiki], turning long strings in big integers (called blocks), allowing for fast retrieval with simple binary search in a list of big numbers. In this context, a block can have many digits: the $k$-th feature is the $k$-th digit, although blocks may have a varying number of digits. I implicitly rely on the Python Bignum library [Wiki] to deal with the computations. Finally, the binary search is further improved and called weighted binary search, accelerating the computations by a factor 3 or 4 in the examples tested. So far, I did not compare with other methods such as vector search based on KNN and ANN (approximate nearest neighbors). But these methods are very relevant in this context.

**Connection to LLM**

The above paragraphs establish the connection to large language models. The problem is strikingly similar to DNA sequence synthetization discussed in section 7.1, where the alphabet has four letters (A, C, G, T) and the words consist of DNA subsequences. The main difference is that DNA sequences are far from random. Yet, the methodology presented here can easily be adapted to arbitrary target distributions. In particular to empirical distributions like those associated to DNA sequencing, or keyword distributions in ordinary text.

Then, as illustrated in the DNA sequencing problem, predictive analytics for GenAI may rely on conditional probabilities such as $P(B_1|B_2)$, where $B_1, B_2$ are consecutive blocks. Transitioning from KS and the multivariate CDF to conditional probabilities is straightforward with the formula $P(B_1|B_2) = P(B_1, B2)/P(B_2)$.

## 2.3.1   Project description

We are dealing with sequences of positive integers defined by a recursion $x_{n+1} = f(x_n) \mod \tau_n$, where $x_0$ is the initial condition, referred to as the main seed. The choice of $f$ leads to $x_n$ randomly going up and down as $n$ increases. On average, the trend is up: eventually, $x_n$ gets bigger and bigger on average, albeit rather slowly to allow for fast computations. The **digits** – what the generated random numbers are – are simply defined as $d_n = x_n \mod b$, where $b$ is called the **base**.

The big contrast with classic random generators is that $\tau_n$ depends on $n$ rather than being fixed. Here $\tau_n$ strictly increases at each iteration, allowing $x_n$ to grow on average. It leads to more random digits, and an infinite period. In the current implementation, $\tau_{n+1} = s + \tau_n$ where $s$ is a moderately large integer. Thus, we have two additional seeds: the step $s$, and $\tau_0$. The function $f$ depends on two integer parameters $p, q$. More specifically,

$$x_{n+1} = p \cdot (x_n//q) \mod \tau_n, \text{ with } p > q > 1,\ x_0, \tau_0 \geq 0,\text{ and } \tau_{n+1} = s + \tau_n. \tag{2.1}$$

Again, the double slash stands for the integer division: $x_n//q = \lfloor x_n/q \rfloor$ where the brackets denote the floor function, or `int` in Python. I use the integer division in Python as it works when $x_n$ becomes extremely large, unlike the standard division combined with `int`.

I implemented three tests of randomness. Before going into the details, let me introduce the concept of **block**. A block consists of a fixed number of digits, used for testing purposes. They are encoded as big integers. For now, let's pretend that a block $B$ has $m$ digits, each with a variable base: the $k$-th digit is in base $b_k$, thus with $0 \leq d_k < b_k$. Then $B$ can be uniquely encoded with the one-to-one mapping

$$B = d_1 + d_2 \times b_1 + d_3 \times b_1 b_2 + d_4 \times b_1 b_2 b_3 + \cdots + d_m \times b_1 b_2 \cdots b_{m-1} \tag{2.2}$$

In this case, $B$ is represented by an integer strictly smaller than $b_1 b_2 \cdots b_m$. In the context of real data, each digit represents a feature; the $k$-th feature can take on $b_k$ different values after approximation or truncation, and a block corresponds to a row in a tabular data set. Of course, the $b_k$' can be quite large and different, depending on the feature, especially for numerical features. At the other extreme, for a 2-category feature, $b_k = 2$. The numeration system based on (2.2) is called the radix system. It is a generalization of the numeration system in fixed base $b$. When a new observation is generated, it is first encoded as a block, then the closest match to any row in the training set can be found with a binary search on all long integers $B$ between 0 and $(b_1 \cdots b_m) - 1$, each one representing a potential row in the training set. Note that in general, the existing $B$'s cover a small subset of all potential values: we are dealing with sparsity. We will use the same search algorithm here, with a fixed base $b$. It makes sense to call it radix search.

Now I can provide a quick overview of the three tests of randomness. The most involved is the block test: it is based on the multivariate KS distance, which in turn relies on radix search. The tests are:

- Run test. The max run test is implemented in section 4.3 for binary digits of quadratic irrationals. The solution presented here is more comprehensive, looking at runs or arbitrary length for all possible digits. These run lengths have a geometric distribution if digits show up randomly. When the base $b$ is very large, there are too many values to fit in a table, so I also compute the following statistics:

$$R(L) = \sum_{d=0}^{b-1} \left( \frac{\rho(L, d) - \rho_0(L, d)}{\sigma(L, d)} \right)^2, \quad L = 1, 2, \ldots \tag{2.3}$$

  where $\rho_0(L, d)$ is the number of runs of length $L$ found in the sequence for digit $d$, and $\rho(L, d)$ is the expected count if the sequence is random; then, it does not depend on $d$. Finally, $\sigma^2(L, d)$ is the theoretical variance, not depending on $d$, to normalize $R(L)$ so that it has a Chi-squared distribution with $b$ degrees of freedom. In short, an unexpectedly large $R(L)$ indicates lack of randomness.

- Spectral test. This test looks at autocorrelations of lag 1, 2, and so on in the digit sequence, to check out whether their behavior is compatible or not with randomness. Because the sequences can be extremely long, the computations are done on the flight, updated one new digit at a time using a buffer, without having to store the whole sequence anywhere.

- Block test. Create a sample of blocks $B$ of size $m$, for instance equally spaced, and independently of the blocks found in the digit sequence. Here $m$ is the number of digits per block, with each digit in base $b$. Hence, each block should occur in the infinite digit sequence with frequency $b^{-m}$, assuming randomness. Compute the CDF value $F(B)$ attached to the underlying theoretical (uniform) distribution, for each of these blocks. Then compute the ECDF or empirical CDF value $F_0(B)$ based on block occurrences: mostly nearest neighbors to $B$, found in the digit sequence. Finally, KS $= \sup |F(B) - F_0(B)|$, and the supremum is over all blocks $B$ in your sample. To magnify any subtle departure from randomness, plot the following: $\delta(B) = F(B) - F_0(B)$ on the Y-axis, and $B$ (in its integer representation) on the X-axis, for all $B$ in your sample.

There is a considerable amount of new material to master, if you were to implement the whole system on your own. The goal is not to create your own code from scratch, but rather to understand and use mine. You are welcome to improve it and create Python objects for the various components. One of the easy tasks is to use the code to compare different random generators and assess the impact of parameters. Another goal is to generalize the code to other contexts such as synthesizing DNA sequences, where the target distribution is not uniform, but comes as an ECDF (empirical CDF) computed on the training set. This also involves conditional probabilities: predicting the next block given the previous ones, when blocks are auto-correlated.

The project consists of the following steps:

**Step 1**: Read and understand the code in section 2.3.3. Identify the different components: the three tests, the base, the digits, the generation of the digits, the blocks and block encoding, the seeds and parameters, and the binary search to locate nearby blocks in the sequence, given a specific block. Run the code (it is on GitHub), see and analyze the results. Understand the computations for the autocorrelation coefficients, as it is done on the fly, in a non-traditional way using a buffer.

**Step 2**: What are the values of $\rho(L, d)$ and $\sigma(L, d)$ in Formula 2.3? How would you optimize a binary search so that it requires fewer steps? To help answer this question, look at the variable `trials`, which counts the combined number of steps used in all the binary searches. Finally, save the generated digits in a file if the sequence is not too long.

**Step 3**: To improve block tests, work with blocks of various sizes. Also, compute the autocorrelations in the block sequence, in the same way it is done in the digit sequence. Given a small-size block $B$ consisting

of two sub-blocks $B_1, B_2$, estimate $P(B_2|B_1)$ on the data, and confirm the independence between $B_1$ and $B_2$, that is, $P(B_2|B_1) = P(B_2)$.

- **Step 4**: Try different seeds $x_0, \tau_0$ and $s$ to check out how the random generators is sensitive to the seeds. Does it change the quality of the randomness? In particular, try to find the lowest possible seed values that still lead to good random numbers. Then, try different combinations of $p, q$ and $b$. Some work, some don't. Can you identify necessary conditions on $p, q, b$ to guarantee good random numbers? Finally, can you find a good combination that makes the sequence $(x_n)$ grow as slowly as possible, while still generating great random numbers. Obviously, you need $p > q$, for instance $p = q + 1$, and $s$ as small as possible. A large base $b$ also helps to extract as much as possible from the sequence $(x_n)$, yet some combinations $(p, q, b)$ don't work, and usually $b > q$ causes problems, making $b = q$ ideal.

### 2.3.2 Solution

Since much of the solution is my code in section 2.3.3, I cover only specific items in this section, with a focus on explaining and interpreting the output tables and visualizations. The main parameters are initialized in section [3] in the code, see lines $278 - 288$. In particular, the current values of $x_n, \tau_n$ are stored in x and tau respectively, while $s$ in Formula (2.1) is denoted as step in the code.

Figure 2.3 shows the function $\delta(B) = F(B) - F_0(B)$, that is, the difference between a perfectly uniform CDF and the approximate ECDF computed on the generated digits, using 500 equally-spaced test blocks, each with 6 digits. The number of blocks (500 here) is specified by n_nodes in the code. The higher the number, the better the approximation. The total number of digits if 640,000. I tried four sets of parameters labeled $HM_1$ to $HM_4$. The parameter sets are specific combinations of $p, q, b$. See lines 294, 316, 337 and 358 for the respective values in the code. I also added Numpy.random for comparison purposes. Note that HM stands for "Home-Made", by contrast to Numpy.
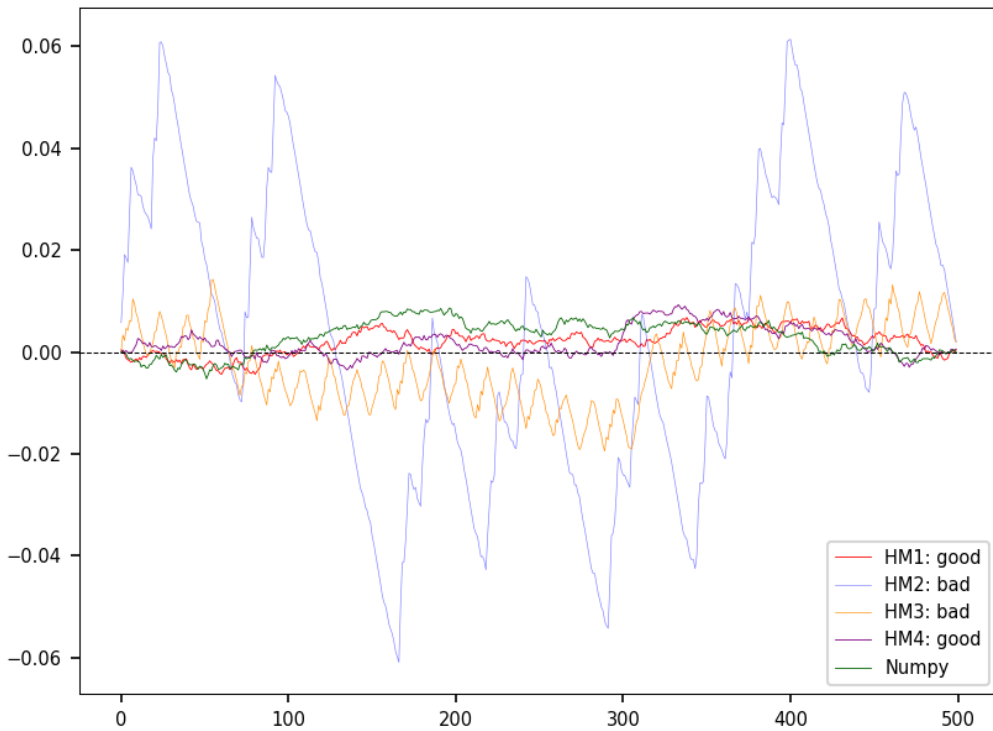


Figure 2.3: Four versions of my generator (HMx) compared to numpy.random

Clearly, $HM_2$ and $HM_3$ generate non-random numbers. As for $HM_1$, $HM_4$ and Numpy, they pass this test. It does not mean that they generate random digits. More tests are needed for verification, in particular with larger block sizes (the parameter block_size) and more nodes (the parameter n_nodes).

Table 2.1 focuses on $HM_1$ only. It shows the number of occurrences for runs of length $L$, for each of the 4 digits ($b = 4$) and various values of $L$. The sequence has 640,000 digits. The 4 rightmost columns are summary statistics: Exp and Avg are respectively the expected and observed counts, while Norm indicates whether or not all the counts, for a specific $L$, are compatible with perfect randomness. If the digits are truly random,

then `Norm` approximately follows a standard normal distribution. In particular, $\texttt{Norm} = (R(L) - b)/(2b)$ where $R(L)$ is defined by (2.3). Clearly, $\text{HM}_1$ passes this test.

| $L$ | $d=0$ | $d=1$ | $d=2$ | $d=3$ | Exp | Avg | Chi2 | Norm |
|---|---|---|---|---|---|---|---|---|
| 1 | 8968 | 8932 | 9074 | 8999 | 9000 | 8993 | 1.44 | $-0.3202$ |
| 2 | 2296 | 2227 | 2294 | 2149 | 2250 | 2242 | 6.81 | 0.3511 |
| 3 | 548 | 532 | 620 | 587 | 563 | 572 | 9.05 | 0.6315 |
| 4 | 146 | 132 | 143 | 131 | 141 | 138 | 1.44 | $-0.3204$ |
| 5 | 36 | 39 | 38 | 36 | 35 | 37 | 0.69 | $-0.4136$ |
| 6 | 3 | 11 | 8 | 10 | 9 | 8 | 4.61 | 0.0759 |
| 7 | 2 | 3 | 3 | 2 | 2 | 3 | 0.62 | $-0.4223$ |
| 8 | 0 | 2 | 0 | 0 | 1 | 1 | 3.83 | $-0.0211$ |

Table 2.1: Runs of length $L$ per digit for $\text{HM}_1$, with summary stats

Table 2.2 shows summary statistics for the 4 home-made generators (HM) and Numpy. One of them ($\text{HM}_3$) has $b = 256$. The normal approximation to `Norm` is especially good when $b > 10$, and bad when the counts are small, that is, when $L$ is large. All the values not compatible with the randomness assumption are highlighted in red. Thus, $\text{HM}_2$ and $\text{HM}_3$ do not produce random digits, confirming the findings from Figure 2.3. Here, AC stands for autocorrelations within digit sequences

| Metric | Level | $\text{HM}_1$ | $\text{HM}_2$ | $\text{HM}_3$ | $\text{HM}_4$ | Numpy |
|---|---|---|---|---|---|---|
| Norm | $L=1$ | $-0.3202$ | 1530.7 | 136.38 | 0.0885 | $-0.0471$ |
| | $L=2$ | 0.3511 | 249.3 | 95.80 | $-0.1944$ | 0.3141 |
| | $L=3$ | 0.6315 | 80.1 | 41.37 | 0.5242 | 0.8902 |
| AC | Lag 1 | 0.0010 | $-0.0046$ | $-0.0731$ | $-0.0002$ | $-0.0007$ |
| | Lag 2 | $-0.0011$ | 0.0022 | 0.0048 | 0.0022 | $-0.0014$ |
| KS | | 0.00676 | 0.06136 | 0.01932 | 0.00929 | 0.00859 |

Table 2.2: High-level comparison of HM and Numpy generators, with red flags

Besides listing the parameters used in the simulation, Table 2.3 features two interesting metrics: last $x_n$, and `Trials`. The former indicates how fast the sequence $x_n$ grows. Numpy is not based on growing sequences, resulting in digits that keep repeating themselves past the finite period. And for random binary digits based on quadratic irrationals, the last $x_n$ would be of the order $2^N \approx 10^{193,000}$, where $N = 640,000$ is the number of digits in the sequence. By contrast, for the HM generators explored here, it is around $10^{10}$.

| Feature | $\text{HM}_1$ | $\text{HM}_2$ | $\text{HM}_3$ | $\text{HM}_4$ | Numpy |
|---|---|---|---|---|---|
| Last $x_n$ | $10^{10}$ | $10^9$ | $10^{10}$ | $10^9$ | n/a |
| Trials | 178 | 2197 | 3497 | 2293 | 150 |
| Base | 4 | 4 | 8 | 256 | 4 |
| Block size | 6 | 6 | 6 | 6 | 6 |
| $p$ | 7 | 6 | 7 | 401 | n/a |
| $q$ | 4 | 4 | 4 | 256 | n/a |

Table 2.3: Top parameters and special metrics

Finally, `Trials` represents the total number of steps needed in the binary search, combined over the 500 test blocks ($500 = \texttt{n\_nodes}$). For each test block, the goal is to find the closest neighbor block in the digit sequence (640,000 digits). The `Trials` value depends on the base, the block size, the number of digits in the sequence, and the number of test blocks. A very small value means that most of the test blocks are already present in the digit sequence: for these test blocks, the binary search is not needed. In the table, a value of 2197 means that on average, each test block requires $2197/500 = 4.39$ trials before finding the closest neighbor

block. The lower `Trials`, the faster the computations. Because I use an optimized binary search, it is already 3–4 times faster than the standard method. The generation of random digits is also just as fast as Numpy.

Regarding the project steps, I now provide answers to selected questions. Let $N$ be the number of digits in base $b$ in the sequence, and $\pi$ be the probability for a run to be of length $L > 0$, for any digit $d < b$. We are dealing with a binomial distribution of parameter $(N, \pi)$. Thus,

$$\rho(L, d) = N\pi, \quad \sigma^2(L, d) = N\pi(1 - \pi), \quad \text{with } \pi = \left(\frac{b-1}{b}\right)^2 \cdot \left(\frac{1}{b}\right)^L,$$

Note that are $\rho(L, d)$ and $\sigma^2(L, d)$ are respectively the expectation and variance of the binomial distribution. This answers the first question in Step 2. For the second question about the binary search, see lines $184 - 197$ in the code. If you set A=1 and B=1 respectively in lines $187$ and $188$, it becomes a standard binary search, with computational complexity $O(\log_2 n)$ in all cases. With my choice of A and B, it is similar to interpolation search [Wiki], which is $O(\log_2 \log_2 n)$ for the average case when the underlying distribution is uniform. See also the recent article on interpolated binary search [23].

### 2.3.3 Python implementation

The code is also on GitHub, here.

```python
1  import numpy as np
2  from collections import OrderedDict
3
4  #--- [1] Functions to generate random digits
5
6  def get_next_digit(x, p, q, tau, step, base, option = "Home-Made"):
7
8      if option == "Numpy":
9          digit = np.random.randint(0, base)
10     elif option == "Home-Made":
11         x = ((p * x) // q) % tau # integer division for big int
12         tau += step
13         digit = x % base
14     return(digit, x, tau)
15
16
17 def update_runs(digit, old_digit, run, max_run, hash_runs):
18
19     if digit == old_digit:
20         run += 1
21     else:
22         if (old_digit, run) in hash_runs:
23             hash_runs[(old_digit, run)] += 1
24         else:
25             hash_runs[(old_digit, run)] = 1
26         if run > max_run:
27             max_run = run
28         run = 1
29     return(run, max_run, hash_runs)
30
31
32 def update_blocks(digit, m, block, base, block_size, hash_blocks):
33
34     if m < block_size:
35         block = base * block + digit
36         m += 1
37     else:
38         if block in hash_blocks:
39             hash_blocks[block] += 1
40         else:
41             hash_blocks[block] = 1
42         block = 0
43         m = 0
44     return(m, block, hash_blocks)
45
```

```python
46
47  def update_cp(digit, k, buffer, max_lag, N, cp_data):
48
49      # processing k-th digit starting at k=2
50      # buffer stores the last max_lag digits
51      # cp stands for the cross-product part (in autocorrelation)
52
53      mu = cp_data[0]
54      cnt = cp_data[1]
55      cp_vals = cp_data[2]
56      cp_cnt = cp_data[3]
57
58      buffer[(k-2) % max_lag] = digit
59      mu += digit
60      cnt += 1
61
62      for lag in range(max_lag):
63          if k-2 >= lag:
64              cp_vals[lag] += digit * buffer[(k-2-lag) % max_lag]
65              cp_cnt[lag] += 1
66
67      cp_data = (mu, cnt, cp_vals, cp_cnt)
68      return(cp_data, buffer)
69
70
71  def generate_digits(N, x0, p, q, tau, step, base, block_size, max_lag, option =
        "Home-Made"):
72
73      # Main function. Also produces output to test randomnes:
74      #    - hash_runs and max_run: input for the run_test function
75      #    - hash_blocks: input for the block_test function
76      #    - cp_data input for the correl_test function
77
78      # for run and block test
79      hash_runs = {}
80      hash_blocks = {}
81      block = 0
82      m = 0
83      run = 0
84      max_run = 0
85      digit = -1
86
87      # for correl_test
88      mu = 0
89      cnt = 0
90      buffer = np.zeros(max_lag)
91      cp_vals = np.zeros(max_lag) # cross-products for autocorrel
92      cp_cnt = np.zeros(max_lag)
93      cp_data = (mu, cnt, cp_vals, cp_cnt)
94
95      x = x0
96
97      for k in range(2, N):
98
99          old_digit = digit
100         (digit, x, tau) = get_next_digit(x, p, q, tau, step, base, option)
101         (run, max_run, hash_runs) = update_runs(digit, old_digit, run, max_run, hash_runs)
102         (m, block, hash_blocks) = update_blocks(digit, m, block, base, block_size,
               hash_blocks)
103         (cp_data, buffer) = update_cp(digit, k, buffer, max_lag, N, cp_data)
104
105     print("----------------")
106     print("PRNG = ", option)
107     print("block_size (digits per block), digit base: %d, %d", block_size, base)
108     if option == "Home-Made":
109         print("p, q: %d, %d" %(p, q))
```

```python
110        print(len(str(x)), "decimal digits in last x")
111    return(hash_runs, hash_blocks, max_run, cp_data)
112
113
114 #--- [2] Functions to perform tests of randomness
115
116 def run_test(base, max_run, hash_runs):
117
118    # For each run, chi2 has approx. chi2 distrib. with base degrees of freedom
119    # This is true assuming the digits are random
120
121    print()
122    print("Digit ", end = " ")
123    if base <= 8:
124        for digit in range(base):
125            print("%8d" %(digit), end =" ")
126    print("  Exp", end = " ")
127    print("  Avg", end = " ") # count average over all digits
128    print(" Chi2", end = " ") # degrees of freedom = base
129    print(" norm", end = " ")
130    print("\n")
131
132    for run in range(1, max_run+1):
133
134        print("Run %3d" % (run), end = " ")
135        prob = ((base-1)/base)**2 * (1/base)**run
136        exp = N*prob
137        var = N*prob*(1-prob)
138        avg = 0
139        chi2 = 0
140
141        for digit in range(0, base):
142            key = (digit, run)
143            count = 0
144            if key in hash_runs:
145                count = hash_runs[key]
146                avg += count
147                chi2 += (count - exp)**2 / var
148            if base <= 8:
149                print("%8d" %(count), end =" ")
150
151        avg /= base
152        norm = (chi2 - base) / (2*base)
153        print("%8d" %(int(0.5 + exp)), end =" ")
154        print("%8d" %(int(0.5 + avg)), end =" ")
155        print("%8.2f" %chi2, end =" ")
156        print("%8.4f" %norm, end =" ")
157        print()
158
159    return()
160
161
162 def get_closest_blocks(block, list, trials):
163
164    # Used in block_test
165    # Return (left_block, right_block) with left_block <= block <= right_block,
166    #  - If block is in list, left_block = block = right_block
167    #  - Otherwise, left_block = list(left), right_block = list(right);
168    #            they are the closest neighbors to block, found in list
169    #  - left_block, right_block are found in list with weighted binary search
170    #  - list must be ordered
171    # trials: to compare spped of binary search with weighted binary search
172
173    found = False
174    left = 0
175    right = len(list) - 1
```

```python
176     delta = 1
177     old_delta = 0
178
179     if block in list:
180         left_block = block
181         right_block = block
182
183     else:
184         while delta != old_delta:
185             trials += 1
186             old_delta = delta
187             A = max(list[right] - block, 0) # in standard binary search: A = 1
188             B = max(block - list[left], 0) # in standard binary search: B = 1
189             middle = (A*left + B*right) // (A + B)
190             if list[middle] > block:
191                 right = middle
192             elif list[middle] < block:
193                 left = middle
194             delta = right - left
195
196         left_block = list[middle]
197         right_block = list[min(middle+1, len(list)-1)]
198
199     return(left_block, right_block, trials)
200
201
202 def true_cdf(block, max_block):
203     # Used in block_test
204     # blocks uniformly distributed on {0, 1, ..., max_block}
205     return((block + 1)/(max_block + 1))
206
207
208 def block_test(hash_blocks, n_nodes, base, block_size):
209
210     # Approximated KS (Kolmogorov-Smirnov distance) between true random and PRNG
211     # Computed only for blocks with block_size digits (each digit in base system)
212     # More nodes means better approximation to KS
213
214     hash_cdf = {}
215     hash_blocks = OrderedDict(sorted(hash_blocks.items()))
216     n_blocks = sum(hash_blocks.values())
217     count = 0
218     trials = 0 # total number of iterations in binary search
219
220
221     for block in hash_blocks:
222         hash_cdf[block] = count + hash_blocks[block]/n_blocks
223         count = hash_cdf[block]
224
225     cdf_list = list(hash_cdf.keys())
226     max_block = base**block_size - 1
227     KS = 0
228     trials = 0 # total number of iterations in binary search
229     arr_cdf = [] # theoretical cdf, values
230     arr_ecdf = [] # empirical cdf (PRMG), values
231     arr_arg = [] # arguments (block number) associated to cdf or ecdf
232
233     for k in range(0, n_nodes):
234
235         block = int(0.5 + k * max_block / n_nodes)
236         (left_block, right_block, trials) = get_closest_blocks(block, cdf_list, trials)
237         cdf_val = true_cdf(block, max_block)
238         ecdf_lval = hash_cdf[left_block] # empirical cdf
239         ecdf_rval = hash_cdf[right_block] # empirical cdf
240         ecdf_val = (ecdf_lval + ecdf_rval) / 2 # empirical cdf
241         arr_cdf.append(cdf_val)
```

```python
242         arr_ecdf.append(ecdf_val)
243         arr_arg.append(block)
244         dist = abs(cdf_val - ecdf_val)
245         if dist > KS:
246             KS = dist
247
248     return(KS, arr_cdf, arr_ecdf, arr_arg, trials)
249
250
251 def autocorrel_test(cp_data, max_lag, base):
252
253     mu = cp_data[0]
254     cnt = cp_data[1]
255     cp_vals = cp_data[2]
256     cp_cnt = cp_data[3]
257
258     mu /= cnt
259     t_mu = (base-1) / 2
260     var = cp_vals[0]/cp_cnt[0] - mu*mu
261     t_var = (base*base -1) / 12
262     print()
263     print("Digit mean: %6.2f (expected: %6.2f)" % (mu, t_mu))
264     print("Digit var : %6.2f (expected: %6.2f)" % (var, t_var))
265     print()
266     print("Digit autocorrelations: ")
267     for k in range(max_lag):
268         autocorrel = (cp_vals[k]/cp_cnt[k] - mu*mu) / var
269         print("Lag %4d: %7.4f" %(k, autocorrel))
270
271     return()
272
273
274 #--- [3] Main section
275
276 # I tested (p, q) in {(3, 2), (7, 4), (13, 8), (401, 256)}
277
278 N = 64000        # number of digits to generate
279 p = 7            # p/q must > 1, preferably >= 1.5
280 q = 4            # I tried q = 2, 4, 8, 16 and so on only
281 base = q         # digit base, base <= q (base = 2 and base = q work)
282 x0 = 50001       # seed to start the bigint sequence in PRNG
283 tau = 41197      # co-seed of home-made PRNG
284 step = 37643     # co-seed of home-made PRNG
285 digit = -1       # fictitious digit before creating real ones
286 block_size = 6   # digits per block for the block test; must be integer > 0
287 n_nodes = 500    # number of nodes for the block test
288 max_lag = 3      # for autocorrel test
289 seed = 104       # needed only with option = "Numpy"
290 np.random.seed(seed)
291
292 #- [3.1] Home-made PRNG with parameters that work
293
294 p, q, base, block_size = 7, 4, 4, 6
295
296 # generate random digits, home-made PRNG
297 (hash_runs, hash_blocks, max_run, cp_data) = generate_digits(N, x0, p, q, tau, step,
298                                                 base, block_size, max_lag,
299                                                 option="Home-Made")
300 # run_test
301 run_test(base,max_run,hash_runs)
302
303 # block test
304 (KS, arr_cdf, arr_ecdf1, arr_arg1, trials) = block_test(hash_blocks, n_nodes,
305                                             base ,block_size)
306 # autocorrel_test
307 autocorrel_test(cp_data, max_lag, base)
```

33

```
308
309  print()
310  print("Trials = ", trials)
311  print("KS = %8.5f\n\n" %(KS))
312
313
314  #- [3.2] Home-made, with parameters that don't work
315
316  p, q, base, block_size = 6, 4, 4, 6
317
318  # generate random digits, home-made PRNG
319  (hash_runs, hash_blocks, max_run, cp_data) = generate_digits(N, x0, p, q, tau, step,
320                                                 base, block_size, max_lag,
321                                                 option="Home-Made")
322  # run_test
323  run_test(base,max_run,hash_runs)
324
325  # block test
326  (KS, arr_cdf, arr_ecdf2, arr_arg2, trials) = block_test(hash_blocks, n_nodes,
327                                          base ,block_size)
328  # autocorrel_test
329  autocorrel_test(cp_data, max_lag, base)
330
331  print()
332  print("Trials = ", trials)
333  print("KS = %8.5f\n\n" %(KS))
334
335  #- [3.3] Home-made, another example of failure
336
337  p, q, base, block_size = 7, 4, 8, 6
338
339  # generate random digits, home-made PRNG
340  (hash_runs, hash_blocks, max_run, cp_data) = generate_digits(N, x0, p, q, tau, step,
341                                                 base, block_size, max_lag,
342                                                 option="Home-Made")
343  # run_test
344  run_test(base,max_run,hash_runs)
345
346  # block test
347  (KS, arr_cdf, arr_ecdf3, arr_arg3, trials) = block_test(hash_blocks, n_nodes,
348                                          base ,block_size)
349  # autocorrel_test
350  autocorrel_test(cp_data, max_lag, base)
351
352  print()
353  print("Trials = ", trials)
354  print("KS = %8.5f\n\n" %(KS))
355
356  #- [3.4] Home-made, using a large base
357
358  p, q, base, block_size = 401, 256, 256, 6
359
360  # generate random digits, home-made PRNG
361  (hash_runs, hash_blocks, max_run, cp_data) = generate_digits(N, x0, p, q, tau, step,
362                                                 base, block_size, max_lag,
363                                                 option="Home-Made")
364  # run_test
365  run_test(base,max_run,hash_runs)
366
367  # block test
368  (KS, arr_cdf, arr_ecdf4, arr_arg4, trials) = block_test(hash_blocks, n_nodes,
369                                          base ,block_size)
370  # autocorrel_test
371  autocorrel_test(cp_data, max_lag, base)
372
373  print()
```

34

```python
374    print("Trials = ", trials)
375    print("KS = %8.5f\n\n" % (KS))
376
377
378    #- [3.5] Numpy PRNG (Mersenne twister)
379
380    # here p, q are irrelevant
381    base, block_size = 4, 6
382
383    # generate random digits, home-made PRNG
384    (hash_runs, hash_blocks, max_run, cp_data) = generate_digits(N, x0, p, q, tau, step,
385                                                      base, block_size, max_lag,
386                                                      option="Numpy")
387    # run_test
388    run_test(base,max_run,hash_runs)
389
390    # block test
391    (KS, arr_cdf, arr_ecdf5, arr_arg5, trials) = block_test(hash_blocks, n_nodes,
392                                                 base ,block_size)
393    # autocorrel_test
394    autocorrel_test(cp_data, max_lag, base)
395
396    print()
397    print("Trials = ", trials)
398    print("KS = %8.5f\n\n" % (KS))
399
400
401    #--- [4] Scatterplot cdf (true random) versus ecdf (based on the two PRNGs)
402
403    import matplotlib.pyplot as plt
404    import matplotlib as mpl
405
406    mpl.rcParams['axes.linewidth'] = 0.5
407    plt.rcParams['xtick.labelsize'] = 7
408    plt.rcParams['ytick.labelsize'] = 7
409
410    arr_cdf = np.array(arr_cdf)
411    delta_ecdf1 = (np.array(arr_ecdf1) - arr_cdf)
412    delta_ecdf2 = (np.array(arr_ecdf2) - arr_cdf)
413    delta_ecdf3 = (np.array(arr_ecdf3) - arr_cdf)
414    delta_ecdf4 = (np.array(arr_ecdf4) - arr_cdf)
415    delta_ecdf5 = (np.array(arr_ecdf5) - arr_cdf)
416
417    # print()
418    # print("blocks (arguments) used to compute ecdf1:\n")
419    # print(arr_arg1)
420
421    plt.plot(delta_ecdf1, linewidth = 0.4, color = 'red', alpha = 1)
422    plt.plot(delta_ecdf2, linewidth = 0.3, color = 'blue', alpha = 0.5)
423    plt.plot(delta_ecdf3, linewidth = 0.3, color = 'darkorange', alpha = 1)
424    plt.plot(delta_ecdf4, linewidth = 0.4, color = 'purple', alpha = 1)
425    plt.plot(delta_ecdf5, linewidth = 0.4, color = 'darkgreen', alpha = 1)
426    plt.axhline(y = 0.0, linewidth = 0.5, color = 'black', linestyle = 'dashed')
427    plt.legend(['HM1: good', 'HM2: bad', 'HM3: bad', 'HM4: good', 'Numpy'],
428               loc='lower right', prop={'size': 7}, )
429    plt.show()
```

# Bibliography

[1] Adel Alamadhi, Michel Planat, and Patrick Solé. Chebyshev's bias and generalized Riemann hypothesis. *Preprint*, pages 1–9, 2011. arXiv:1112.2398 [Link]. 65

[2] K. Binswanger and P. Embrechts. Longest runs in coin tossing. *Insurance: Mathematics and Economics*, 15:139–149, 1994. [Link]. 58

[3] Ramiro Camino, Christian Hammerschmidt, and Radu State. Generating multi-categorical samples with generative adversarial networks. *Preprint*, pages 1–7, 2018. arXiv:1807.01202 [Link]. 94

[4] Fida Dankar et al. A multi-dimensional evaluation of synthetic data generators. *IEEE Access*, pages 11147–11158, 2022. [Link]. 93

[5] Antónia Földes. The limit distribution of the length of the longest head-run. *Periodica Mathematica Hungarica*, 10:301–310, 1979. [Link]. 59

[6] Louis Gordon, Mark F. Schilling, and Michael S. Waterman. An extreme value theory for long head runs. *Probability Theory and Related Fields*, 72:279–287, 1986. [Link]. 58

[7] Vincent Granville. *Statistics: New Foundations, Toolbox, and Machine Learning Recipes*. Data Science Central, 2019. 24

[8] Vincent Granville. *Synthetic Data and Generative AI*. MLTechniques.com, 2022. [Link]. 100

[9] Vincent Granville. Feature clustering: A simple solution to many machine learning problems. *Preprint*, pages 1–6, 2023. MLTechniques.com [Link]. 82

[10] Vincent Granville. Generative AI: Synthetic data vendor comparison and benchmarking best practices. *Preprint*, pages 1–13, 2023. MLTechniques.com [Link]. 75

[11] Vincent Granville. Generative AI technology break-through: Spectacular performance of new synthesizer. *Preprint*, pages 1–16, 2023. MLTechniques.com [Link]. 12, 15, 121

[12] Vincent Granville. *Gentle Introduction To Chaotic Dynamical Systems*. MLTechniques.com, 2023. [Link]. 24, 60, 64

[13] Vincent Granville. How to fix a failing generative adversarial network. *Preprint*, pages 1–10, 2023. ML-Techniques.com [Link]. 14

[14] Vincent Granville. Massively speed-up your learning algorithm, with stochastic thinning. *Preprint*, pages 1–13, 2023. MLTechniques.com [Link]. 13, 82

[15] Vincent Granville. Smart grid search for faster hyperparameter tuning. *Preprint*, pages 1–8, 2023. ML-Techniques.com [Link]. 13, 80, 82

[16] Vincent Granville. *Synthetic Data and Generative AI*. Elsevier, 2024. [Link]. 38, 43, 44, 45, 47, 48, 54, 55, 57, 63, 74, 75, 79, 80, 82, 83, 85, 93, 94

[17] Elisabeth Griesbauer. *Vine Copula Based Synthetic Data Generation for Classification*. 2022. Master Thesis, Technical University of Munich [Link]. 82

[18] Emil Grosswald. Oscillation theorems of arithmetical functions. *Transactions of the American Mathematical Society*, 126:1–28, 1967. [Link]. 65

[19] Adam J. Harper. Moments of random multiplicative functions, II: High moments. *Algebra and Number Theory*, 13(10):2277–2321, 2019. [Link]. 65

[20] Adam J. Harper. Moments of random multiplicative functions, I: Low moments, better than squareroot cancellation, and critical multiplicative chaos. *Forum of Mathematics, Pi*, 8:1–95, 2020. [Link]. 65

[21] Adam J. Harper. Almost sure large fluctuations of random multiplicative functions. *Preprint*, pages 1–38, 2021. arXiv [Link]. 65

[22] Zsolt Karacsony and Jozsefne Libor. Longest runs in coin tossing. teaching recursive formulae, asymptotic theorems and computer simulations. *Teaching Mathematics and Computer Science*, 9:261–274, 2011. [Link]. 59

[23] Adnan Saher Mohammed, Şahin Emrah Amrahov, and Fatih V. Çelebi. Interpolated binary search: An efficient hybrid search algorithm on ordered datasets. *Engineering Science and Technology*, 24:1072–1079, 2021. [Link]. 29

[24] Tamas Mori. The a.s. limit distribution of the longest head run. *Canadian Journal of Mathematics*, 45:1245–1262, 1993. [Link]. 59

[25] Michel Planat and Patrick Solé. Efficient prime counting and the Chebyshev primes. *Preprint*, pages 1–15, 2011. arXiv:1109.6489 [Link]. 65

[26] M.S. Schmookler and K.J. Nowka. Bounds on runs of zeros and ones for algebraic functions. *Proceedings 15th IEEE Symposium on Computer Arithmetic*, pages 7–12, 2001. ARITH-15 [Link]. 58

[27] Mark Shilling. The longest run of heads. *The College Mathematics Journal*, 21:196–207, 2018. [Link]. 58

[28] Chang Su, Linglin Wei, and Xianzhong Xie. Churn prediction in telecommunications industry based on conditional Wasserstein GAN. *IEEE International Conference on High Performance Computing, Data, and Analytics*, pages 186–191, 2022. IEEE HiPC 2022 [Link]. 93

[29] Terence Tao. Biases between consecutive primes. *Tao's blog*, 2016. [Link]. 65

[30] Ruonan Yu, Songhua Liu, and Xinchao Wang. Dataset distillation: A comprehensive review. *Preprint*, pages 1–23, 2022. Submitted to IEEE PAMI [Link]. 80

# Index