



New Approaches to Uniformity Testing

Winston Li

Rutgers University

Mentor: Periklis Papakonstantinou

July 17, 2025



Problem

Given a string of bits X , can we determine if X is **random**?



Problem

Given a string of bits X , can we determine if X is **random**?

Hypothesis Testing

Write $X := X_1X_2 \dots X_n$. If X is random, then for all $i \in [n]$, $X_i \sim \text{Bern}(\frac{1}{2})$



Problem

Given a string of bits X , can we determine if X is **random**?

Hypothesis Testing

Write $X := X_1X_2 \dots X_n$. If X is random, then for all $i \in [n]$, $X_i \sim \text{Bern}(\frac{1}{2})$

Frequency Test (NIST)

Define $S_n = X_1 + X_2 + \dots + X_n$. By Central Limit Theorem, for large enough n (too big), $S_n \rightarrow \mathcal{N}(\frac{n}{2}, \frac{n}{4})$. Performing a z-test, we can compute the probability for the value of S_n occurring.



Problem

Given a string of bits X , can we determine if X is **random**?

Hypothesis Testing

Write $X := X_1X_2 \dots X_n$. If X is random, then for all $i \in [n]$, $X_i \sim \text{Bern}(\frac{1}{2})$

Frequency Test (NIST)

Define $S_n = X_1 + X_2 + \dots + X_n$. By Central Limit Theorem, for large enough n (too big), $S_n \rightarrow \mathcal{N}(\frac{n}{2}, \frac{n}{4})$. Performing a z-test, we can compute the probability for the value of S_n occurring.

Question

Can we do better?



Chernoff-Hoeffding Bounds

Let $X = \sum_{i=1}^n X_i$, where X_i are independent random variables distributed in $[0, 1]$. For all $t > 0$,

$$\Pr[X > \mathbb{E}[X] + t] \leq e^{-2t^2/n}$$

$$\Pr[X < \mathbb{E}[X] - t] \leq e^{-2t^2/n}$$



Chernoff-Hoeffding Bounds

Let $X = \sum_{i=1}^n X_i$, where X_i are independent random variables distributed in $[0, 1]$. For all $t > 0$,

$$\Pr[X > \mathbb{E}[X] + t] \leq e^{-2t^2/n}$$

$$\Pr[X < \mathbb{E}[X] - t] \leq e^{-2t^2/n}$$

Frequency Test

Recall $S_n = \sum_{i \in [n]} X_i$ where $X_i \sim \text{Bern}(\frac{1}{2})$. Then $\mathbb{E}[X] = \frac{n}{2}$, so $\Pr[|S_n - \frac{n}{2}| > t] \leq 2e^{-2t^2/n}$. For $\alpha \in (0, 1)$, take $t = \sqrt{\frac{n}{2}(\log 2 - \log \alpha)}$. The probability S_n is farther than t away from $\frac{n}{2}$ is less than α .



Block Frequency Test

Let X be a bit string of length $n \times m$. Write $X = X_1 \dots X_n$, where X_i is a substring of length m . Let $\pi_i = \frac{1}{m} \sum_{j=1}^m X_{ij}$ be the proportion of bits in X_i that are 1. Then $S = 4m \sum_{i=1}^n (\pi_i - \frac{1}{2})^2$ is a χ^2 distribution.



Block Frequency Test

Let X be a bit string of length $n \times m$. Write $X = X_1 \dots X_n$, where X_i is a substring of length m . Let $\pi_i = \frac{1}{m} \sum_{j=1}^m X_{ij}$ be the proportion of bits in X_i that are 1. Then $S = 4m \sum_{i=1}^n (\pi_i - \frac{1}{2})^2$ is a χ^2 distribution.

Chernoff Bounds (χ^2)

Let $Y_i = \frac{2\pi_i - m}{\sqrt{m}}$, so $S = \sum_{i=1}^n Y_i^2$. Applying Chernoff bounds,

$$M_{Y_i}(\lambda) \leq (1 - 2\lambda)^{-1/2} \quad M_S(\lambda) \leq (1 - 2\lambda)^{-n/2}$$
$$\Pr[S \geq t] \leq \min_{\lambda} e^{-\lambda t} (1 - 2\lambda)^{-n/2} = \left(\frac{t}{n}\right)^{n/2} e^{(n-t)/2}$$



Block Frequency Test

Let X be a bit string of length $n \times m$. Write $X = X_1 \dots X_n$, where X_i is a substring of length m . Let $\pi_i = \frac{1}{m} \sum_{j=1}^m X_{ij}$ be the proportion of bits in X_i that are 1. Then $S = 4m \sum_{i=1}^n (\pi_i - \frac{1}{2})^2$ is a χ^2 distribution.

Chernoff Bounds (χ^2)

Let $Y_i = \frac{2\pi_i - 1}{\sqrt{m}}$, so $S = \sum_{i=1}^n Y_i^2$. Applying Chernoff bounds,

$$M_{Y_i}(\lambda) \leq (1 - 2\lambda)^{-1/2} \quad M_S(\lambda) \leq (1 - 2\lambda)^{-n/2}$$
$$\Pr[S \geq t] \leq \min_{\lambda} e^{-\lambda t} (1 - 2\lambda)^{-n/2} = \left(\frac{t}{n}\right)^{n/2} e^{(n-t)/2}$$

Problem

Notice m is not included in the bound. For χ^2 tests, this scales with class count, but not string length.



Definition

The **statistical difference** between two random variables X, Y is $\Delta(X, Y) := \frac{1}{2}|X - Y|_1$. X and Y are ε -close if $\Delta(X, Y) \leq \varepsilon$.



Definition

The **statistical difference** between two random variables X, Y is $\Delta(X, Y) := \frac{1}{2} \|X - Y\|_1$. X and Y are ε -close if $\Delta(X, Y) \leq \varepsilon$.

Definition

Let X be a random variable.

The **Shannon Entropy** is $H_{Sh}(X) = \mathbb{E}_{x \sim X}[-\log \Pr[X = x]]$.

The **min-entropy** is $H_\infty(X) = \min_x (-\log \Pr[X = x])$.



Definition

The **statistical difference** between two random variables X, Y is $\Delta(X, Y) := \frac{1}{2} \|X - Y\|_1$. X and Y are ε -close if $\Delta(X, Y) \leq \varepsilon$.

Definition

Let X be a random variable.

The **Shannon Entropy** is $H_{Sh}(X) = \mathbb{E}_{x \sim X}[-\log \Pr[X = x]]$.

The **min-entropy** is $H_\infty(X) = \min_x (-\log \Pr[X = x])$.

Definition

A random variable X is a k -source if $H_\infty(X) \geq k$, or $\Pr[X = x] \leq 2^{-k}$.



Definition

The **statistical difference** between two random variables X, Y is $\Delta(X, Y) := \frac{1}{2} \|X - Y\|_1$. X and Y are ε -close if $\Delta(X, Y) \leq \varepsilon$.

Definition

Let X be a random variable.

The **Shannon Entropy** is $H_{Sh}(X) = \mathbb{E}_{x \sim X}[-\log \Pr[X = x]]$.

The **min-entropy** is $H_\infty(X) = \min_x (-\log \Pr[X = x])$.

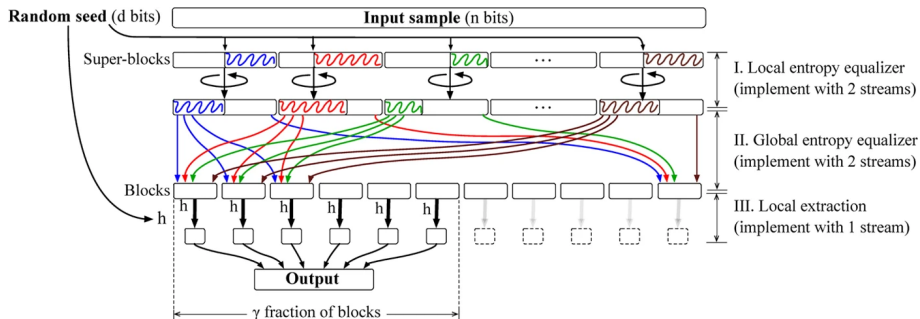
Definition

A random variable X is a k -source if $H_\infty(X) \geq k$, or $\Pr[X = x] \leq 2^{-k}$.

Definition

$\text{Ext} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ is a (k, ε) -**extractor** if for every k -source X , $\text{Ext}(X, U_d)$ is ε close to U_m .

Randomized Re-Bucketing Extractor





Definition

A **Toeplitz Matrix** is a matrix with constant diagonals.

$$A = \begin{bmatrix} a_0 & a_1 & a_2 & a_3 \\ a_{-1} & a_0 & a_1 & a_2 \\ a_{-2} & a_{-1} & a_0 & a_1 \end{bmatrix} \quad A_{i,j} = A_{i+1,j+1}$$



Definition

A **Toeplitz Matrix** is a matrix with constant diagonals.

$$A = \begin{bmatrix} a_0 & a_1 & a_2 & a_3 \\ a_{-1} & a_0 & a_1 & a_2 \\ a_{-2} & a_{-1} & a_0 & a_1 \end{bmatrix} \quad A_{i,j} = A_{i+1,j+1}$$

Property 1

Toeplitz matrices over \mathbb{F}_2 form a family of **universal hash functions**.



Definition

A **Toeplitz Matrix** is a matrix with constant diagonals.

$$A = \begin{bmatrix} a_0 & a_1 & a_2 & a_3 \\ a_{-1} & a_0 & a_1 & a_2 \\ a_{-2} & a_{-1} & a_0 & a_1 \end{bmatrix} \quad A_{i,j} = A_{i+1,j+1}$$

Property 1

Toeplitz matrices over \mathbb{F}_2 form a family of **universal hash functions**.

Property 2

An $m \times n$ Toeplitz matrix is generated by a $m + n - 1$ dimension vector.

$$A \in \mathbb{F}_2^{m \times n} \quad w = (a_{1-m}, \dots, a_{m+n-1})$$
$$Av = \bigotimes_{i=1}^n v_i w_{i+m}$$

Source Coding Theorem (Shannon)

For a constant stream of bits, the maximum compression ratio asymptotically approaches the Shannon entropy of the source.

Source Coding Theorem (Shannon)

For a constant stream of bits, the maximum compression ratio asymptotically approaches the Shannon entropy of the source.

Experiment

Fix some computable sequence $f : \mathbb{N} \rightarrow \{0, 1\}$ and proportion $p \in [0, 1]$. Use f to poison a bit string $X := X_1 \dots X_n$, generating $Y := Y_1 \dots Y_n$,

$$Y_i = \begin{cases} X_i & \text{with probability } p \\ f(i) & \text{with probability } 1 - p \end{cases}$$

We expect $H_{Sh}(X) = n$ and $H_{Sh}(Y) \approx nH_b(\frac{1-p}{2})$. Expect $|\text{compress}(Y)|$ to be similar to $H_{Sh}(Y)$.



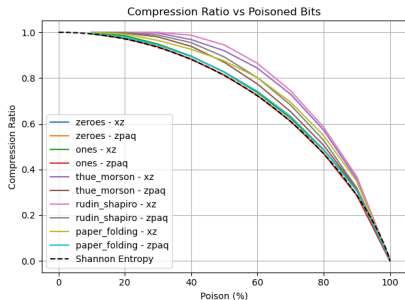
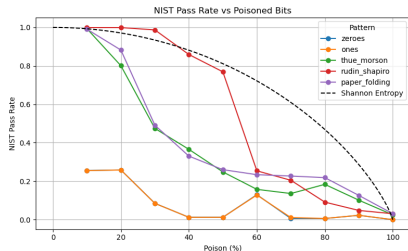
Setup

- Extract seed from 2.6GB of audio files using BIWZ extractor.
- Extract data from 64GB movie file with RRB extractor and seed.
- Poison extracted data using a fixed sequence.
- Compress poisoned file with xz and zpaq at max setting.
- Run NIST tests on poisoned file.



Setup

- Extract seed from 2.6GB of audio files using BIWZ extractor.
- Extract data from 64GB movie file with RRB extractor and seed.
- Poison extracted data using a fixed sequence.
- Compress poisoned file with xz and zpaq at max setting.
- Run NIST tests on poisoned file.





RRB Parameterization

Given a bit string of length 67,108,864 (8 MiB), split it into 50 substrings. Run the NIST test suite on each substring. The program takes an unusually long time computing a DFT (Discrete Fourier Transform).

Truncating the string to $50 \times 1,000,000$ bits, the tests run instantly.



RRB Parameterization

Given a bit string of length 67,108,864 (8 MiB), split it into 50 substrings. Run the NIST test suite on each substring. The program takes an unusually long time computing a DFT (Discrete Fourier Transform).

Truncating the string to $50 \times 1,000,000$ bits, the tests run instantly.

Lightning Strikes!

It turns out $\lfloor \frac{67108864}{50} \rfloor = 1342177$ is a prime number. For a bit string of length $n = p_1 p_2 \dots p_k$, NIST's DFT runs in $O(n \sum_k p_i)$. This is around $O(n \log n)$ if n has many factors, but $O(n^2)$ if n is prime.



RRB Parameterization

Given a bit string of length 67,108,864 (8 MiB), split it into 50 substrings. Run the NIST test suite on each substring. The program takes an unusually long time computing a DFT (Discrete Fourier Transform).

Truncating the string to $50 \times 1,000,000$ bits, the tests run instantly.

Lightning Strikes!

It turns out $\lfloor \frac{67108864}{50} \rfloor = 1342177$ is a prime number. For a bit string of length $n = p_1 p_2 \dots p_k$, NIST's DFT runs in $O(n \sum_k p_i)$. This is around $O(n \log n)$ if n has many factors, but $O(n^2)$ if n is prime.

Solution

Splitting the bit string into 64 substrings, or dropping a bit from each substring resolves this problem.



- I would like to thank my mentor Professor Periklis Papakonstantinou for guiding me through this project.
- I would like to thank Larry for being Larry.
- This work was carried out as part of the 2025 DIMACS REU program at Rutgers University, supported by NSF grant CCF-2247342

References

- [1] Boaz Barak, Russell Impagliazzo, and Avi Wigderson. “Extracting Randomness Using Few Independent Sources”. In: *SIAM Journal on Computing* 36.4 (2006), pp. 1095–1118. DOI: [10.1137/S0097539705447141](https://doi.org/10.1137/S0097539705447141).
- [2] Lawrence Bassham et al. *A Statistical Test Suite for Random and Pseudorandom Number Generators for Cryptographic Applications*. en. 2010.
- [3] Devdatt Dubhashi and Alessandro Panconesi. *Concentration of Measure for the Analysis of Randomized Algorithms*. 1st. USA: Cambridge University Press, 2009. ISBN: 0521884276.
- [4] Periklis A. Papakonstantinou, David P. Woodruff, and Guang Yang. “True Randomness from Big Data”. en. In: *Scientific Reports* 6.1 (2016). DOI: [10.1038/srep33740](https://doi.org/10.1038/srep33740).
- [5] C. E. Shannon. “A mathematical theory of communication”. In: *The Bell System Technical Journal* 27.3 (1948), pp. 379–423. DOI: [10.1002/j.1538-7305.1948.tb01338.x](https://doi.org/10.1002/j.1538-7305.1948.tb01338.x).

