

REU in Random Walk

Xiaoyu Liu

November 21, 2022

Contents

1 Adapted K-S Test for discrete sampling	1
---	----------

Note 6: Adapted K-S Test for discrete sampling

Sun 20 Nov 2022 18:32

1 Adapted K-S Test for discrete sampling

Motivation

Simulation of continuous random processes are discrete in nature. Two factors primarily contribute to the discretization:

1. Numerical precision limit.
2. The approximation of continuous process via discrete process.

Discretization gives rise to the **oversampling effect**: When we generate a sample, the sample we get is from a discretized version of the underlying distribution. Hence, if we increasing the size of the sample, the empirical cdf cannot converge uniformly to the true underlying distribution. This brings huge problem to goodness-of-fit tests, since their consistency relies on Glivenko-Cantelli Theorem, the fact that empirical cdf must uniformly converge to true distribution.

Therefore, we need to develop techniques to perform goodness-of-fit tests on discrete simulations and avoid the oversampling effect.

Problem Background

Partition. We define a partition P of \mathbb{R} as a countable set of distinct real numbers. We say P is finer or equal to P' if $P' \subset P$, and write $P > P'$.

Discretization. Let F be a continuous real-valued distribution. Let P be a partition of \mathbb{R} . We say $F^{(P)}$ is a discretization of F with respect to partition P if $\tilde{F}(x) = F(y)$, where $y \in P$ minimizes $|x - y|$. We may verify that $F^{(P)}$ is indeed discrete, since it takes values in $F(P)$ which has at most countably many values.

Contribution

We have developed and implemented techniques to perform k-s test to samples from a discretized distribution. The techniques are applicable to the two following cases:

1. (One sample test) The underlying distribution F_0 is continuous real-valued. We have a sample $\mathbf{X} = X_1, \dots, X_n$ drawn independently from some distribution $F_X^{(P)}$ that is discretization of continuous F_X with respect to P . We want to test if $F_0 \equiv F_X$.
2. (Two sample test) We have samples \mathbf{X} and \mathbf{Y} from discretizations $F_X^{(P)}$ and $F_Y^{(P')}$ with respect to partitions P and P' , where $P \leq P'$. We want to test if $F_X \equiv F_Y$.

Development

Our null hypothesis in the one-sample case is:

$$H_0 : F_0 \equiv F_X.$$

Under the null hypothesis, we must have

$$H'_0 : F_0^{(P)} \equiv F_X^{(P)}.$$

That is, the discretizations of F_0 and F_X are equal.

We see that $H_0 \Rightarrow H'_0$, so whenever we reject H'_0 we must also reject H_0 . Since our sample is from $F_X^{(P)}$, we have just enough information to test H'_0 . Hence, the decision criterion is set to reject H_0 if we reject H'_0 .

On one hand, we can control type I error if we assume we can control type I error in testing H'_0 , since

$$\begin{aligned} & \Pr[H_0 | \text{reject } H_0] \\ &= \Pr[H_0 | \text{reject } H'_0] \\ &= \Pr[H_0 | H'_0 \wedge \text{reject } H'_0] \Pr[H'_0 | \text{reject } H'_0] \\ &\quad + \Pr[H_0 | \neg H'_0 \wedge \text{reject } H'_0] \Pr[\neg H'_0 | \text{reject } H'_0] \\ &= \Pr[H_0 | H'_0 \wedge \text{reject } H'_0] \Pr[H'_0 | \text{reject } H'_0] \\ &\leq \Pr[H'_0 | \text{reject } H'_0]. \end{aligned}$$

On the other hand, we cannot control type II error, and the method above yields upper bound 1. Intuitively, even if H'_0 is true, we can say nothing about probability that H_0 is true, unless we have a notion of measure within the space of all continuous distributions. But we are treating F_0 as an *arbitrary* distribution, so we cannot obtain meaningful probability bound.

What remains is to find a way to test H'_0 . K-S test is not directly applicable to discretizations, since it only works for continuous distributions. Nevertheless, we can use **linear interpolation** to transform $F_0^{(P)}$ and $F_X^{(P)}$ to continuous piecewise-linear distributions. Then the K-S test would be applicable. It is easy to see that

Claim. There is a 1-1 correspondence between discrete distributions wrt P and piecewise linear functions wrt P . Moreover, $F_0^{(P)} \equiv F_X^{(P)}$ iff $\tilde{F}_0^{(P)} \equiv \tilde{F}_X^{(P)}$.

For two-sample case, the null hypothesis is:

$$H_0 : F_X \equiv F_Y.$$

We weaken the null hypothesis to the partition P (the coarser partition).

$$H'_0 : F_X^{(P)} = F_Y^{(P)}.$$

Then the same reasoning holds. We also justify that the decision power of the sample depends only on the coarsest of P and P' . (TODO)

Linear Interpolation

Calculation of discretization of a continuous distribution is straightforward from the definition. What we need to develop here is a method to perturb the sample in the way that performs linear interpolation on its underlying distribution.

Let X_1, X_2, \dots, X_m be samples from $F_X^{(P)}$. We define the noisy sample \tilde{X}_i as follows: if $X_i = x_i \in P$, then let \tilde{X}_i be a random sample from the distribution with pdf

$$\tilde{f}_i(x) = \begin{cases} 0 & x < \frac{x_{i-1} + x_i}{2} \\ \frac{1}{x_i - x_{i-1}} & x < x_i \\ \frac{1}{x_{i+1} - x_i} & x < \frac{x_i + x_{i+1}}{2} \\ 0 & x \geq \frac{x_i + x_{i+1}}{2} \end{cases}.$$

Then $\tilde{X}_1, \dots, \tilde{X}_m$ is a sample from $\tilde{F}_X^{(P)}$, the linear interpolation of $F_X^{(P)}$.

In the special case when $P = \{q\Delta x + r : q \in \mathbb{Z}\}$ for some fixed $\Delta x > 0$ and $0 \leq r < \Delta x$, we have

$$\tilde{X}_i \equiv_d X_i + \text{Unif}(-\Delta x/2, \Delta x/2).$$

Application to SSRW simulation

For SSRW simulation, we take F_0 to be the cdf of standard normal distribution. Because the discreteness is due to independent steps of same length, the partition P falls into the special case. It remains to find Δx and r .

Suppose we simulate n steps of a SSRW. Then the range of scaled distribution is

$$\frac{n[-1, 1]}{\sqrt{n}} = [-\sqrt{n}, \sqrt{n}].$$

When n is even, due to parity of steps, there are $n+1$ possible final positions after n steps. So the range is evenly divided into n intervals. Hence the resolution, Δx , is

$$\Delta x = \frac{2\sqrt{n}}{n} = \frac{2}{\sqrt{n}}.$$

Since 0 gets mapped to 0 during the rescaling, $r = 0$.

Using these knowledge of P , we can calculate the discretized-interpolated $\tilde{F}_0^{(P)}$ and noised samples.

Result

See notebook:

Further steps

- When studying edge reinforced random walks, where discrete simulation is applied extensively, the adapted k-s test can give reliable goodness-of-fit test.