

# THE STATISTICAL PROPERTIES OF RANDOM BITSTREAMS AND THE SAMPLING DISTRIBUTION OF COSINE SIMILARITY

GRAHAM L. GILLER

ABSTRACT. We summarize the statistical properties of statistics computed from independent random bitstreams including the commonly discussed *support* and *cosine similarity*. We derive the moments of the asymptotically normal approximation to the sampling distribution of the cosine similarity of independent random bitstreams and compare those computed moments to those measured by Monte-Carlo simulation. We find agreement for bitstreams of *internet scale* in length (i.e. of order 10 000 bits), and also much smaller (10 bits), and demonstrate that the expected value of the cosine similarity of *independent* bitstreams may be very significantly distant from zero. To compensate for this bias we propose a new statistic, the *Support Adjusted Cosine Similarity* or SACS.

## 1. DEFINITION OF RANDOM BITSTREAMS

A random bitstream  $(N, p)$  is a vector  $\mathbf{B}$  in  $\mathbb{B}^N$  where each component of the vector is independently drawn<sup>1</sup> from Bernoulli( $p$ ). i.e.

$$(1) \quad b_i = \mathbf{B} \cdot \mathbf{e}_i \sim \text{Bernoulli}(p) \quad \forall i \in [1, N]$$

We will represent this as the notation  $\mathbf{B} \sim \text{Bitstream}(N, p)$ . Two bitstreams of equal length,  $\mathbf{A} \sim \text{Bitstream}(N, p)$  and  $\mathbf{B} \sim \text{Bitstream}(N, q)$ , are *independent* if  $E(a_i b_j) = E(a_i)E(b_j) \quad \forall i, j \in [1, N]$ . From this it follows that

$$(2) \quad \mathbf{A} \wedge \mathbf{B} \sim \text{Bitstream}(N, pq)$$

where  $\mathbf{A} \wedge \mathbf{B}$  represents the Boolean *and* operation performed pairwise on bitstream elements and is identically equal to the vector of the pairwise products of the bitstream elements. i.e.

$$(3) \quad (\mathbf{A} \wedge \mathbf{B}) \cdot \mathbf{e}_i = a_i b_i$$

Trivially,  $\mathbf{A} \wedge \mathbf{A} = \mathbf{A}$  for any bitstream,  $\mathbf{A}$ .

## 2. THE SUPPORT AND INNER PRODUCT OF BITSTREAMS

2.1. **Definition.** In market basket and collaborative filtering analysis it is common to introduce a term called the “support”<sup>2</sup>[1] of the bitstream which represents the

---

*Date:* February 28, 2013. *Giller Investments Research Note:* 20121024/1.

<sup>1</sup>In the following we use the notation  $x \sim D$  to mean “ $x$  is distributed as  $D$ .”

<sup>2</sup>The usage has some commonalities with, but is not identical to, the usage of the term in measure theory.

count of the bits that are non-zero divided by the length of the bitstream,  $N$ . We write

$$(4) \quad \text{supp } \mathbf{B} = \frac{1}{N} \sum_{i=1}^N b_i = \bar{b}.$$

**2.2. Distribution.** From Equation 1 we see that

$$(5) \quad \mathbf{B} \sim \text{Bitstream}(N, p) \Rightarrow N \text{ supp } \mathbf{B} \sim \text{Binomial}(N, p).$$

As a scaled Binomial distribution, we immediately have

$$(6) \quad \begin{aligned} E(\text{supp } \mathbf{B}) &= p \\ \text{and } \text{Var}(\text{supp } \mathbf{B}) &= \frac{p(1-p)}{N} \end{aligned}$$

and it immediately follows that the distribution of  $\text{supp } \mathbf{B}$  has a well understood limiting Normal distribution defined solely by these moments.

**2.3. Relationship to the Inner Product.** From Equation 4, we see that the inner product of two bitstreams is proportional to the support of their logical product i.e.

$$(7) \quad \mathbf{A} \cdot \mathbf{B} = N \text{ supp } \mathbf{A} \wedge \mathbf{B} \sim \text{Binomial}(N, pq)$$

where  $\mathbf{A} \sim \text{Bitstream}(N, p)$  and  $\mathbf{B} \sim \text{Bitstream}(N, q)$ .

**2.4. Covariance of the Inner Product with its Arguments.** Let  $\mathbf{A}$  and  $\mathbf{B}$  be independent bitstreams of equal length as defined above. Then the inner products  $\mathbf{A} \cdot \mathbf{B}$  and  $\mathbf{A} \cdot \mathbf{A}$  are both Binomial random variables but they are not independent. Suppose for some  $\mathbf{A}$ ,  $N \text{ supp } \mathbf{A} = n \sim \text{Binomial}(N, p)$  i.e. there are  $n$  non-zero bits where  $n$  is a Binomial variate. The elements of  $\mathbf{B}$  at the unsupported bits are irrelevant as  $x \wedge 0 = 0 \forall x \in \mathbb{B}$  and so  $\mathbf{A} \cdot \mathbf{B} \leq n$ .

We may reorder the *arbitrary* indices of *both* vectors so that  $\mathbf{A} = (\mathbf{1}_n, \mathbf{0}_{N-n})$  and  $\mathbf{B} = (\mathbf{B}_n, \mathbf{B}_{N-n})$  where the subscripted vectors represent the partition of a vector in the first  $n$  bits,  $\mathbf{B}_n$ , and the remaining  $N - n$  bits,  $\mathbf{B}_{N-n}$  and  $\mathbf{1}$  and  $\mathbf{0}$  are the unit and zero vectors respectively. Clearly,  $\mathbf{A} \cdot \mathbf{B} = \mathbf{B}_n \cdot \mathbf{1}_n = n \text{ supp } \mathbf{B}_n$ . Since the full length vectors are independent then any similar partition of them must also be independent and so  $n \text{ supp } \mathbf{B}_n \sim \text{Binomial}(n, q)$  and *independent* of  $\mathbf{A}$  except for the dependence on  $n$ . Therefore

$$(8) \quad \begin{aligned} \mathbf{A} \cdot \mathbf{B} \mid \mathbf{A} \cdot \mathbf{A} &\sim \text{Binomial}(\mathbf{A} \cdot \mathbf{A}, q) \\ \text{where } \mathbf{A} \cdot \mathbf{A} &\sim \text{Binomial}(N, p). \end{aligned}$$

Thus  $\mathbf{A} \cdot \mathbf{B}$  is seen to be drawn from a *nested* Binomial distribution with parameters  $(N, p, q)$ .

The covariance of  $\mathbf{A} \cdot \mathbf{B}$  and  $\mathbf{A} \cdot \mathbf{A}$  is given by

$$(9) \quad \text{Cov}(\mathbf{A} \cdot \mathbf{B}, \mathbf{A} \cdot \mathbf{A}) = E(\mathbf{A} \cdot \mathbf{B} \times \mathbf{A} \cdot \mathbf{A}) - E(\mathbf{A} \cdot \mathbf{B})E(\mathbf{A} \cdot \mathbf{A}).$$

The latter two terms are known to be  $Npq$  and  $Np$  respectively. The cross term may be evaluated as

(10)

$$E(\mathbf{A} \cdot \mathbf{B} \times \mathbf{A} \cdot \mathbf{A}) = \sum_{\mathbf{A} \cdot \mathbf{A}=0}^N \mathbf{A} \cdot \mathbf{A} \Pr(\mathbf{A} \cdot \mathbf{A}|N, p) \sum_{\mathbf{A} \cdot \mathbf{B}=0}^{\mathbf{A} \cdot \mathbf{A}} \mathbf{A} \cdot \mathbf{B} \Pr(\mathbf{A} \cdot \mathbf{B}|\mathbf{A} \cdot \mathbf{A}, q)$$

and (remarkably) the entire expression reduces<sup>3</sup> to

$$(11) \quad \text{Cov}(\mathbf{A} \cdot \mathbf{B}, \mathbf{A} \cdot \mathbf{A}) = Np(1-p)q.$$

Similarly

$$(12) \quad \text{Cov}(\mathbf{A} \cdot \mathbf{B}, \mathbf{B} \cdot \mathbf{B}) = Npq(1-q).$$

### 3. THE COSINE SIMILARITY OF TWO BITSTREAMS

**3.1. Definition.** We define the *cosine similarity* of two bitstreams of equal length as

$$(13) \quad \cos(\mathbf{A}, \mathbf{B}) = \frac{\mathbf{A} \cdot \mathbf{B}}{\sqrt{\mathbf{A} \cdot \mathbf{A} \mathbf{B} \cdot \mathbf{B}}},$$

following the usual definition for an inner product space. Some properties of this metric are discussed in [3]. Unlike the cosine in Euclidean inner product space, this statistic is bounded below at 0.

**3.2. Relationship to the Support.** The properties of the support then allow us to write

$$(14) \quad \cos(\mathbf{A}, \mathbf{B}) = \frac{\text{supp } \mathbf{A} \wedge \mathbf{B}}{\sqrt{\text{supp } \mathbf{A} \text{ supp } \mathbf{B}}}.$$

**3.3. A Tighter Upper Bound.** The authors of [3] point out that it also follows from the definition of support that

$$(15) \quad \cos(\mathbf{A}, \mathbf{B}) \leq \sqrt{\frac{\text{supp } \mathbf{A}}{\text{supp } \mathbf{B}}},$$

when  $\text{supp } \mathbf{A} \leq \text{supp } \mathbf{B}$ .

**3.4. A Normal Approximation to the Sampling Distribution of the Cosine Similarity.** We will approximate the sampling distribution of the cosine similarity of equal length independent bitstreams by using the well-known ‘‘Delta Method’’[2] of taking the expectations of the first order Taylor expansion of a function of random variables about the means of the underlying variables. i.e.

$$(16) \quad Ef(\mathbf{r}) \simeq f(E\mathbf{r}) + (\nabla f) \cdot E\delta\mathbf{r}$$

$$(17) \quad \text{and } \text{Var } f \simeq (\nabla f) \cdot \text{Var } \mathbf{r} \cdot \nabla f$$

<sup>3</sup>This was simplified using *Mathematica* v6. as the expression  
`FullSimplify[Sum[AA Pr[AA, N, p] Sum[AB Pr[AB, AA, q], {AB, 0, AA}], {AA, 0, N}]]`  
 where `Pr[n., N., p.] := PDF[BinomialDistribution[N, p], n].`

Consider

$$(18) \quad f(\mathbf{r}) = \frac{x}{\sqrt{yz}},$$

where  $\mathbf{r} = (x, y, z)'$  and  $x, y$ , and  $z$  are correlated asymptotically normal random variables with a known mean vector and covariance matrix. In this case we have

$$(19) \quad \frac{\partial f}{\partial x} = \frac{f}{x}, \quad \frac{\partial f}{\partial y} = -\frac{f}{2y} \quad \text{and} \quad \frac{\partial f}{\partial z} = -\frac{f}{2z}$$

and so

$$(20) \quad E \cos(\mathbf{A}, \mathbf{B}) \simeq \sqrt{pq}$$

$$(21) \quad \text{Var} \cos(\mathbf{A}, \mathbf{B}) \simeq \frac{4 - 3p - 3q + 2pq}{4N},$$

where we have substituted in the previously developed values of the inner product expectations, variances and covariances using  $x = \mathbf{A} \cdot \mathbf{B}$ ,  $y = \mathbf{A} \cdot \mathbf{A}$  and  $z = \mathbf{B} \cdot \mathbf{B}$ . Therefore, we model  $\cos(\mathbf{A}, \mathbf{B})$  as asymptotically distributed as  $\text{Normal}(\mu, \sigma^2)$  where  $\mu$  is given by Equation 20 and  $\sigma^2$  is given by Equation 21.

The expected value of the cosine similarity for fully independent bitstreams may be quite large, for example if  $p = q = 0.5$  then it is 0.5 but for internet scale bitstreams of large lengths, e.g.  $N = 10\,000$ , the sampling error about this number is quite small in comparison (0.006 in this case). Note that as  $p, q \rightarrow 1$  then  $\sigma^2 \rightarrow 0$  whereas  $p, q \rightarrow 0$  then  $\sigma^2 \rightarrow 1/N$ .

**3.5. Numerical Simulations.** To verify the formulae of Section 3.4 we conducted a small number of Monte-Carlo experiments in *Mathematica*. We generate a random bitstream with the function

$$(22) \quad \text{RandomBitstream}[N_-, p_-] := \text{Table}[\text{If}[\text{Random}[] \leq p, 1, 0], \{N\}]$$

and generate the of cosine similarity of two independent bitstreams with

$$(23) \quad \begin{aligned} \text{Cosine}[N_-, p_-, q_-] := & \text{Module}[\{\mathbf{A}, \mathbf{B}\}, \mathbf{A} = \text{RandomBitstream}[N, p]; \\ & \mathbf{B} = \text{RandomBitstream}[N, q]; \\ & (\mathbf{A} \cdot \mathbf{B}) / \text{Sqrt}[(\mathbf{A} \cdot \mathbf{A})(\mathbf{B} \cdot \mathbf{B})]. \end{aligned}$$

One thousand trials of the cosine similarity were then generated for bitstreams of length 10 000, 100 and 10 bits and various combinations of  $p$  and  $q$ . The sample mean and variance for these trials were then compared with the predicted values. This results are show in Table 1 on page 5. Without performing any precise statistical testing, it seems fair to conclude that the formulae established are reasonable estimates of the first two moments of the sampling distribution of the cosine similarity of independent random bitmaps.

**3.6. Support Adjusted Cosine Similarity.** To compensate for the bias in the cosine similarity as discussed in Section 3.4, we propose adjusting the formula of Equation 13 by subtracting the asymptotic expected value. i.e.

$$(24) \quad \text{sacos}(\mathbf{A}, \mathbf{B}) = \cos(\mathbf{A}, \mathbf{B}) - \sqrt{\text{supp } \mathbf{A} \text{ supp } \mathbf{B}}$$

$$(25) \quad = \frac{\text{supp } \mathbf{A} \wedge \mathbf{B} - \text{supp } \mathbf{A} \text{ supp } \mathbf{B}}{\sqrt{\text{supp } \mathbf{A} \text{ supp } \mathbf{B}}}$$

where Equation 14 has been substituted into Equation 24 to give Equation 25. We note there is a parallel between the structure of this equation and that of the Pearson correlation coefficient.

## REFERENCES

- [1] S. Brin, R. Motwani, and C. Silverstein. Beyond market baskets: Generalizing association rules to correlations. *Proceedings of SIGMOD*, pages 265–276, 1997.
- [2] G.W. Oehlert. A note on the delta method. *The American Statistician*, 46:27–29, 1992.
- [3] Shiwei Zhu, Junjie Wu, Hui Xiong, and Guoping Xia. Scaling up top- $k$  cosine similarity search. *Data & Knowledge Engineering*, 70:60–83, 2011.

## TABLES

$N$	$p$	$q$	Exp. Mean	Obs. Mean	Exp. Var.	Obs. Var.
10 000	0.5	0.5	0.5	0.500091	0.000038	0.000039
10 000	0.5	0.25	0.353553	0.353659	0.000050	0.000049
10 000	0.5	0.125	0.25	0.249839	0.000056	0.000053
10 000	0.5	0.0625	0.176777	0.176880	0.000059	0.000058
100	0.5	0.5	0.5	0.499110	0.00375	0.003740
100	0.5	0.25	0.353553	0.347558	0.005	0.005328
100	0.5	0.125	0.25	0.248546	0.005625	0.005876
100	0.5	0.0625	0.176777	0.175621	0.005938	0.006224
10	0.5	0.5	0.5	0.485566	0.0375	0.041569
10	0.5	0.25	0.353553	0.330916	0.05	0.049284
10	0.5	0.125	0.25	0.287991	0.05625	0.055474
10	0.5	0.0625	0.176777	0.254209	0.059375	0.056185

TABLE 1. Results of numerical simulations of the cosine similarity between independent random bitmaps of length 10 000, 100 and 10 bits. For each parameter set, 1 000 trials were generated.