# A Note on Chebyshev Inequality: To Explain or to Predict

By Amaresh Das

*Southern University at New Orleans, United States*

*Abstract-* The question is: What proportion of the total probability of a random varriable *X* lies within a certain interval of the mea $\mu$? What is the probability of being hit by a meteor greater in size than five times the standard deviation above the mean? Because it can be applied to completely arbitrary distributions(unknown except for mean and variables), the inequality generally gives a poor bound compared to what might be deduced if more aspects are known about the distribution involved.

*Keywords:* eucladian norm, monotonic function, jensen inequality.

*GJSFR-F Classification:* MSC 2010: 11Y16

ANOTEONCHEBYSHEVINEQUALITYTOEXPLAINORTOPREDICT

*Strictly as per the compliance and regulations of:*

# A Note on Chebyshev Inequality: To Explain or to Predict

Amaresh Das

*Abstract-* The question is: What proportion of the total probability of a random variable $X$ lies within a certain interval of the mea $\mu$? What is the probability of being hit by a meteor greater in size than five times the standard deviation above the mean? Because it can be applied to completely arbitrary distributions(unknown except for mean and variables), the inequality generally gives a poor bound compared to what might be deduced if more aspects are known about the distribution involved.

*Keywords: eucladian norm, monotonic function, jensen inequality.*

## I. Introduction

This note introduces a useful inequality called the Chebyshev (often called Tchebysheff) inequality but its explanation requires a good familiarity in calculus. One of the most difficult and important tasks before the statistician is to discover the probability distribution that may be involved in any problem. If one is unsure what the underlying distribution is, it is comforting to know that there are more universal inequalities that may give us some useful information; At the very least, the Chebyshev inequality allows one to bound how far away from the mean the random variable could be.

## II. Exploring the Inequality

*Theorem 1*

Let $\mu(X)$ be a nonnegative function of the random variable $X$. If $E[\mu(\mathrm{x})]$ exists, then, for every positive constant $C$

$$\Pr[\mu(X) \geq C] \leq \frac{[\mu(\mathrm{x})]}{\mathrm{C}} \tag{1}$$

The proof[1] is given when the random variable[2] $X$ of the continuous type, but, the proof can be adapted to the discrete case if we replace integrals by sums. Let $A = [\mathrm{x}; \mu(\mathrm{x}) \geq \mathrm{C}]$ and let $f(\mathrm{x})$ denote the pdf. of $X$ Then

---

*Author: Southern University at New Orleans and University of New Orleans. e-mail: adas2@cox.net*

[1] See any standard text, for example, Mood and Grsybill [5]

[2] Chebyshev's inequality is usually stated for random variables, but can be generalized to a statement about measure spaces.

Fix t and let $A_t$ be defined as $A_t = [x \in X \mid f(x) \geq t]$ and $\mathrm{I_A}$ be the indicator function of the set $\mathrm{A}_t$. Then it is not difficult to see for any $t$, $0 < g(t)\, \mathrm{I_A} \leq g(f(x))\, \mathrm{I_A}$ since $g$ is not decreasing on the range of $t$ and, therefore,

$g(t)\,\mu(A_t) = \int_A g \circ f\, d\mu \leq \int_x g\, \text{Of}\, d\mu$. The desired monotonicity follows from dividing the above inequality by $g(t)$

$$E \mu (\text{x})] = \int_{\infty}^{\infty} \mu(x) \, f(x) \, dx = f(x) \, dx =$$
$$\int_A \mu(x) \, f(x) \, dx + \int_{A^*} \mu(x) f(x0 \, dx \tag{2}$$

Since each of the integrals in the extreme right-hand member of the preceding equation is nonnegative, the left-hand side member is greater than or equal to either of them. In particular,

$$E\left[\mu(X)\right] \geq \int_A \mu(\text{x}) f(\text{x}) \, dx \tag{3}$$

However, if $\text{x} \in A$ then $\mu(x) \geq C$; accordingly, the right-hand member of the preceding inequality is not increased if we replace $\mu(x)$ by C. Thus

$$E\left[\mu(x)\right] \geq C \int_A f(\text{x}) \, dx, \tag{4}$$

Since

$$\int_A f(\text{x}) \, dx = \Pr(X \in A) = Pr[\mu(X) \geq C \tag{5}$$

It follows that

$$E[\mu(X)] \geq C \Pr\left[\mu(X) \geq C\right] \tag{6}$$

which is the desired result.

The theorem is a generalization of an inequality which is often called Chebyshev's inequality.[3] One can establish the inequality this way[4]:

*Theorem 2*

Let the random variable $X$ have a distribution of probability about which we assume only that there is a finite variance $\sigma^2$. This of course, implies that there is a mean $\mu$. Then for every $K > 0$

$$\Pr\left(\left|X - \mu\right| \geq k\sigma\right) \leq \frac{1}{k^2} \tag{7}$$

In theorem 1 take $\mu(X) = (X - \mu)^2$ and $C = k^2 \sigma^2$ Then we have

$$\Pr\left[(X - \mu)^2 \geq k^2\sigma^2\right] \geq k^2 \sigma^2 \leq \frac{E[(X - \mu)^2]}{k^2 \sigma^2} \tag{8}$$

---

[3] Ferentinos[1] has shown that for a vector $X = (x_1, x_2 \ldots) \mu = (\mu_1, \mu_2 \ldots)$ variance $\sigma^2 = (\sigma_1^2, \sigma_2^2 \ldots)$ and

The Euclidian norm $\|.\|$ that $\Pr(\|x - \mu\| \geq k \|\sigma\| \leq \frac{1}{k^2}$

[4] Symmetry of the distribution decreases the inequality's bounds by a factor of 2 while unimodality sharpens the bounds by a factor of 4/9.Because the mean and the mode in a unimodal distribution differ by at most $\sqrt{3}$ standard deviations at most 5% of a symmetrical unimodal distribution lies outside $(2\sqrt{10} + 3\sqrt{3})/3$ standard deviations of the mean (approximately 3.840 standard deviations). This is sharper than the bounds provided by the Chebyshev inequality (approximately 4.472 standard deviations).These bounds on the mean are less sharp than those that can be derived from symmetry of the distribution alone which shows that at most 5% of the distribution lies outside approximately 3.162 standard deviations of the mean. The known Vysochanskiï–Petunin inequality further sharpens this bound by showing that for such a distribution that at most 5% of the distribution lies outside $4\sqrt{5}/3$ (approximately 2.981) standard deviations of the mean. See Kotz et al [ 3]

Since the numerator of the right-hand side member of the preceding inequality is $\sigma^2$, the inequality may be written as

$$\Pr\left( \ |X - \mu| \ \geq k\sigma \ \right) \ \leq \frac{1}{k^2} \tag{9}$$

which is the desired result. Obviously we should take the positive number k to be greater than one to have the inequality of consequence. Because it can be applied to completely arbitrary distributions (unknown except for mean and variance), the inequality generally gives a poor bound compared to what might be deduced if more aspects are known about the distribution

*Table 1:* Chebyshev Inequality

| k | Min. % within k stan deviations of mean | Max. % beyond k stand deviations from mean |
|---|---|---|
| 1 | 0% | 100% |
| $\sqrt{2}$ | 50% | 50% |
| 1.5 | 55.56% | 44.44% |
| 2 | 75% | 25% |
| 3 | 88.8889% | 11.1111% |
| 4 | 93.75% | 6.25% |
| 5 | 96% | 4% |
| 6 | 97.2222% | 2.7778% |
| 7 | 97.9592% | 2.0408% |
| 8 | 98.4375% | 1.5625% |
| 9 | 98.7654% | 1.2346% |
| 10 | 99% | 1% |

Although Chebyshev inequality enables you to find an answer to the questions we raised at the very outset, it comes to the rescue in offering at least an appropriate answer.

*Table 2:* Chebyshev Inequality Bounds and Actual Bounds

| K value | Chebyshev | Gaussian | Chi-square | t |
|---|---|---|---|---|
| 2 | .45 | .05412 | .05241 | .05312 |
| 3 | 10 | .00231 | .00511 | .03121 |
| 3 | .07 | .00001 | .00313 | .00613. |

How tight is the broadly applicable inequality? We can calculate Chebyshev inequality and contrast that value with the exact calculation obtained from knowing

---

[5] There are many extensions to Chebyshev inequality, for example, Chebyshev inequality of exponential version. Inequalities for bounded variables, inequalities in the multivariate case, or its use in infinite dimensional case; see [Stellato, *et al*. [7], Lal [4]] There may be integral inequality, too. An extension to higher moments is also possible.

If $f.g\ ;[a,b]\rightarrow \mathbf{R}$ are less monotonic functions of the same monotonicity, then

$$\frac{1}{b-a}\int_a^b dx \geq \left[\frac{1}{b-a}\ \int_a^b f(\,x\,)\,dx\right]\left[\frac{1}{b-a}\right]\ \int_a^b g(x)\,dx$$

the probability function as indicated as follows: The resilts are tabulated in the above Table 2.

$$\Pr\{\,|x|\ge k\,\sigma \;=\; \int_{-\alpha}^{-k\sigma} df\,(x) + \int_{k\sigma}^{\infty} df\,(x) \tag{11}$$

The Chebyshev theorem typically provides rather loose bounds[6;]. However, these bounds cannot in general (remaining true for arbitrary distributions) be improved upon. The bounds are sharp for the following example: for any $k \ge 1$,

$$x = \begin{bmatrix} 1 \text{ with probability } \dfrac{1}{2k^2} \\[2mm] 0 \text{ with probability } 1 - \dfrac{1}{2k^2} \\[2mm] 1 \text{ with probability } \dfrac{1}{2k^2} \end{bmatrix} \tag{12}$$

*Exercise 1*

Is it possible to find an upper bound for this integral?

$$\int_{0}^{A} (\,A - x\,)\; P(\; x0\; dx \tag{13}$$

*Hint;* Lower is easy to find by using Markow's Inequality but how to find the Upper bound?

## III.   Concluding Remark

This note discusses the Chebyshev inequality as a very app approximate but universally applicable upper bound on probability. The Chebyshev inequality allows us to bound how far away from the mean the random variable could be. It is rather remarkable that one can find inequalities on probability that will hold for any distribution.

### References  Références Referencias

1. Dasgupta, A ( 2000)  Best Constants in Chebyshev Inequalities withVarious Applications, *Biometrika* 5 (1),pp 186 -200.
2. Ferentinos K  (1982), 'On Tchebysheff  Type Inequalities'.  *Trabajos Estadist Investigacion Oper,* 33, pp 125 -132.
3. Kotz, Samuel, Balakrishnan N, and John Norman Continuous *Multivariate Distribution, Models and Application,* 2012 Houghton- Miffler.
4. Lal, D N (1955) 'A Note on a Form of Tchebusheff's Inequality Two or More Variables', *Sankhya,* 15 (3) Indian Statistical Institute, Calcutta, pp  300 -320.
5. Mood M and Grayhill, *An  Introduction to the  Theory of  Statistics*, McGraw Hill, 1963.

---

If  *f* and *g* are of opposite monotonicity, then the above) inequality works in the opposite way. This inequality is related to Jensen's inequality or Kantorovich's inequality.

[6] Although Chebyshev's inequality may not be necessarily true for finite samples. Samuelson's inequality states that all values of a sample will lie within $\sqrt{(N-1)}$ standard deviations of the mean. Chebyshev's bound improves as the sample size increases. When $N = 10$, Samuelson's inequality states that all members of the sample lie within 3 standard deviations of the mean: in contrast Chebyshev's states that 99.5% of the sample lies within 13.5789 standard deviations of the mean. When $N = 100$, Samuelson's inequality states that all members of the sample lie within approximately 9.9499 standard deviations of the mean: Chebyshev's states that 99% of the sample lies within 10 standard deviations of the mean. See Dasgupta[1]. DasGupta has determined a set of best possible bounds for a normal distribution for this inequality. Steliga and Szynal [6] have extended these bounds to the Pareto distribution.

6. SteglaK, and Szynal (2012) "On Markov Type Inequalities' *International Journal of Pure and Applied Mathematics,* 58 (2), pp 137- 152.

7. StellatoB, P Bart, Goulart P (2016) 'Multivariate Chebysheff Inequality with Estimable Mean and Variance', *The American Statistician*, ISSN 0003 - 1306.

Notes

Notes

30

This page is intentionally left blank