# Chapter 83
# Mining Statistically Significant Substrings based on the Chi–Square Measure

**Sourav Dutta**
*IBM Research Lab, India*

**Arnab Bhattacharya**
*Indian Institute of Technology Kanpur, India*

## ABSTRACT

*With the tremendous expansion of reservoirs of sequence data stored worldwide, efficient mining of large string databases in various domains including intrusion detection systems, player statistics, texts, and proteins, has emerged as a practical challenge. Searching for an unusual pattern within long strings of data is one of the foremost requirements for many diverse applications. Given a string, the problem is to identify the substrings that differ the most from the expected or normal behavior, i.e., the substrings that are statistically significant (or, in other words, less likely to occur due to chance alone). We first survey and analyze the different statistical measures available to meet this end. Next, we argue that the most appropriate metric is the chi-square measure. Finally, we discuss different approaches and algorithms proposed for retrieving the top-k substrings with the largest chi-square measure.*

## INTRODUCTION

Detection or identification of statistically significant sequences or mining interesting patterns from a given string has lately emerged as an important area of study (Denise et al., 2001; Ye & Chen, 2001). In such applications, we are given an input string composed of symbols from an alphabet set with a probability distribution defining the chance of occurrence of each symbol, and the aim is to find those portions of the string that deviate most from their expected nature, and are thus potent sources of hidden pattern and information. Such solutions come handy in automated monitoring systems, such as in a cluster of sensors sensing the ambient temperature for possible fire alert, or

a network server sniffing the network for intrusion detection. Also, text analysis of blogs, stock market trend deciphering, detection of protein mutation and the identification of good and bad career patches of a sports icon can be few of the target applications. It is such diverse utility that makes the study and development of this field challenging and necessary.

## STATISTICAL MODELS AND TOOLS

Establishing a relationship of the empirical or observed results of an experiment to factors affecting the system or to pure chance calls for various statistical models and measures. In such scenarios, an observation is deemed statistically significant if its presence cannot be attributed to randomness alone. The literature hosts a number of statistical models to capture the uniqueness of such observations such as *p-value* and *z-score*. In the next few sections, we discuss different important statistical tools that are used for this purpose.

Before venturing forward, we provide a formal definition of the problem.

**Problem 1:** Given a string $S$ of length $l$ comprising symbols from the alphabet set $\Sigma$ of cardinality $m$, and with a given probability distribution $P$ modeling the chance of occurrence of each symbol in $\Sigma$, the problem is to efficiently identify and extract the top-$k$ substrings that exhibit the largest deviation from the expected nature, i.e., the substrings that are most statistically significant.

It is this measure of deviation of a sequence that we will capture by using various statistical models. In the remainder of the chapter, we interchangeably use the term string with sequence and substring with subsequence.
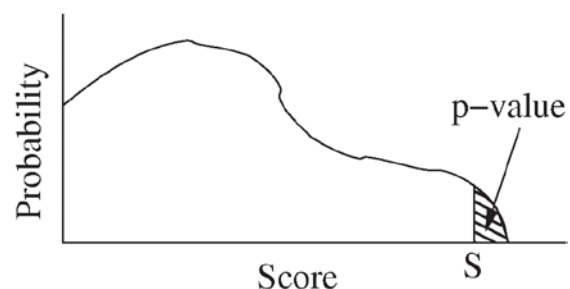
## Hypothesis Testing and P-Value

Given an observation sample $X$ (in this case a substring), with an associated score of $S(X)$, the *p-value* of $X$ is defined as the probability of obtaining a random sample with score $S(X)$ or greater under the same probability model (Bejerano et al., 2004; Regnier & Vandenbogaert, 2006). For each such observation, we test the null hypothesis $H_0$ that the substring is drawn from the given probability model $P$ against the alternate hypothesis $H_1$ that the subsequence is not drawn from the same probability distribution. The p-value measures the chance of rejecting the null hypothesis; in other words, the less the p-value, the less likely it is that the null hypothesis is true.

Figure 1 shows an example. For a particular score $S$, the shaded area represents the chance of having a sample with a score greater than the one under consideration. In other words, the p-value is the value of the cumulative density function (cdf) measured at $S$ subtracted from the total probability, i.e.,

$$pvalue(S) = 1 - cdf(S).$$

If the probability density function (pdf) of the scores is known, it is relatively simpler to compute the p-value of a particular score using the above formula. However, in most real situations, the pdf is hard to estimate or can be non-parametric. The accurate computation of the p-value then needs all the possible outcomes to be listed, their scores

*Figure 1. Computing the p-value of X with score S*

## Related Content

About Shan's Bioinformatics in Research of Biomimicry of Robot-Engineering Systems
Dina Kharicheva (2022). *International Journal of Applied Research in Bioinformatics (pp. 1-12).*
www.irma-international.org/article/shan-bioinformatics-research-biomimicry-robot/290344

Implementation of n-gram Methodology for Rotten Tomatoes Review Dataset Sentiment Analysis
Prayag Tiwari, Brojo Kishore Mishra, Sachin Kumarand Vivek Kumar (2017). *International Journal of Knowledge Discovery in Bioinformatics (pp. 30-41).*
www.irma-international.org/article/implementation-gram-methodology-rotten-tomatoes/178605

A Multi Agent Pharmacoinformatics Reference Model for the Improvement of Hospital Management
Tagelsir Mohamed Gasmelseid (2012). *Pharmacoinformatics and Drug Discovery Technologies: Theories and Applications  (pp. 187-201).*
www.irma-international.org/chapter/multi-agent-pharmacoinformatics-reference-model/64072

Towards Optimal Microarray Universal Reference Sample Designs: An In-Silico Optimization Approach
George Potamias, Sofia Kaforouand Dimitris Kafetzopoulos (2013). *Bioinformatics: Concepts, Methodologies, Tools, and Applications  (pp. 1676-1687).*
www.irma-international.org/chapter/towards-optimal-microarray-universal-reference/76141

Graphical Analysis and Visualization Tools for Protein Interaction Networks
Sirisha Gollapudi, Alex Marshall, Daniel Zadikand Charlie Hodgman (2009). *Biological Data Mining in Protein Interaction Networks (pp. 286-311).*
www.irma-international.org/chapter/graphical-analysis-visualization-tools-protein/5570