

□ 应用统计学丛书

Extremes in Random Fields

A Theory and Its Applications

随机域中的极值统计学

理论及应用 (英文版)

Benjamin Yakir 著



高等教育出版社
HIGHER EDUCATION PRESS

□ 应用统计学丛书

Extremes in Random Fields

A Theory and Its Applications

随机域中的极值统计学

理论及应用 (英文版)

Benjamin Yakir 著

SUIJIYU ZHONG DE JIZHI TONGJIXUE
LILUN JIYINGYONG

Author

Benjamin Yakir

Department of Statistics, Mount Scopus

The Hebrew University of Jerusalem

Jerusalem 91905, Israel

© 2013 Higher Education Press, 4 Dewai Dajie, 100120, Beijing, P. R. China

图书在版编目 (C I P) 数据

随机域中的极值统计学: 理论及应用 = Extremes in random fields: a theory and its applications: 英文 / (以) 亚基尔 (Yakir, B.) 著. -- 北京: 高等教育出版社, 2013. 9

ISBN 978-7-04-037817-7

I. ①随… II. ①亚… III. ①统计学 - 研究 - 英文
IV. ①C8

中国版本图书馆 CIP 数据核字 (2013) 第 187655 号

策划编辑 王丽萍

责任编辑 李华英

封面设计 姜 磊

责任印制 田 甜

出版发行 高等教育出版社
社 址 北京市西城区德外大街4号
邮政编码 100120
印 刷 北京铭成印刷有限公司
开 本 787 mm × 1092 mm 1/16
印 张 15.25
字 数 320 千字
购书热线 010-58581118

咨询电话 400-810-0598
网 址 <http://www.hep.edu.cn>
<http://www.hep.com.cn>
网上订购 <http://www.landaco.com>
<http://www.landaco.com.cn>
版 次 2013 年 9 月第 1 版
印 次 2013 年 9 月第 1 次印刷
定 价 59.00 元

本书如有缺页、倒页、脱页等质量问题, 请到所购图书销售部门联系调换

版权所有 侵权必究

物 料 号 37817-00

Extremes in Random Fields

应用统计学丛书

书号	书名	著译者
9787040378177	随机域中的极值统计学：理论及应用（英文版）	Benjamin Yakir 著
9787040322927	金融工程中的蒙特卡罗方法	Paul Glasserman 著 范韶华、孙武军 译
9787040348309	大维统计分析	白志东、郑术蓉、姜丹丹
9787040348286	结构方程模型：Mplus与应用（英文版）	王济川、王小倩 著
9787040348262	生存分析：模型与应用（英文版）	刘宪
9787040345407	MINITAB软件入门：最易学实用的统计分析教程	吴令云 等 编著
9787040321883	结构方程模型：方法与应用	王济川、王小倩、姜宝法 著
9787040319682	结构方程模型：贝叶斯方法	李锡钦 著 蔡敬衡、潘俊豪、周影辉 译
9787040315370	随机环境中的马尔可夫过程	胡迪鹤 著
9787040256390	统计诊断	韦博成、林金官、解锋昌 编著
9787040250626	R语言与统计分析	汤银才 主编
9787040247510	属性数据分析引论（第二版）	Alan Agresti 著 张淑梅、王睿、曾莉 译
9787040182934	金融市场中的统计模型和方法	黎子良、邢海鹏 著 姚沛沛 译

网上购书：academic.hep.com.cn,

www.china-pub.com, 卓越, 当当

其他订购办法:

各使用单位可向高等教育出版社读者服务部

汇款订购。书款通过邮局汇款或银行转账

均可。购书免邮费, 发票随后寄出。

单位地址: 北京西城区德外大街4号

电 话: 010-58581118/7/6/5/4

传 真: 010-58581113

通过邮局汇款:

地 址: 北京西城区德外大街4号

户 名: 高等教育出版社销售部综合业务部

通过银行转账:

户 名: 高等教育出版社有限公司

开 户 行: 交通银行北京马甸支行

银行账号: 110060437018010037603

To David

Preface

This text started as class notes for a course that I gave in the Mathematical Sciences Center (MSC) in Tsinghua University, Beijing, that got overblown and became a book. I was enjoying a sabbatical leave in the Department of Statistics and Applied Probability (DSAP) of the National University of Singapore when I was given an offer to teach a summer course in China. Of course I accepted. How could I resist the opportunity to fulfil a childhood dream of visiting China?

After accepting the proposal I had to decide what to teach. I decided to fulfil yet another dream, the dream of summarizing and unifying a subject I was writing about all my career, even before I knew what the subject was. The subject is the distribution of extremes in random fields and the analysis of statistical problems that can be formulated in relation to such extremes. Immediately after obtaining my PhD, and as a continuation of my PhD thesis, I was interested in the investigation of the average run length of the Shiryaev–Roberts change-point detection rule. Therefore, I found it natural to try to address a challenge that was presented to me by David Siegmund during a barbecue meal that he prepared for me in his yard. The challenge was to develop a simpler method for analyzing this average run length. In an attempt to attack this problem I began experimenting with the likelihood ratio identity, one of David’s favorite techniques, and followed the road that eventually led me to writing this book.

The original problem was the investigation of average run length in a sequential change-point detection problem.¹ However, the basic technique that was developed turned out to be useful for the investigation of a relatively wide array of different statistical problems that involve the distribution of maxima.² Among other things, David and I used the method in order to investigate the significance level of sequence alignment, for the computation of the false detection rate in

¹ Yakir B., Pollak M. A new representation for the renewal-theoretic constant appearing in asymptotic approximations of large deviations. *Ann. Appl. Probab.* **8**, 749–774 (1998).

² Grossman S., Yakir B. Large deviations for global maxima of independent superadditive processes with negative drift and an application to optimal sequence alignment. *Bernoulli* **5**, 829–845 (2004).

Siegmund D.O., Yakir B. Approximate p-values for local sequence alignments. *Ann. Statist.* **28**, 657–680 (2000).

Siegmund D.O., Yakir B. Statistical analysis of direct identity by descent mapping. *Ann. Hum. Genet.* **67**, 464–470 (2003).

Siegmund D.O., Yakir B. Correction note: Approximate p-values for local sequence alignments. *Ann. Statist.* **31**, 1027–1031 (2003).

scanning statistic, for producing more efficient ways of simulation, etc. Each application required this modification or that trick in order to apply the basic principle. However, after 20 years of repeating the same argument even I was able to identify the pattern. The thrust of this book is a description of the pattern and the demonstration of its usefulness in the analysis of nontrivial statistical problems.

The basic argument relies on a likelihood ratio identity that uses a sum of likelihood ratios. This identity translates the original problem that involves the approximation of a vanishingly small probability to a problem that calls for the summation of approximations of expectations. The expectations are with respect to alternative distributions in which the event in question is much more likely to occur. Moreover, by carefully selecting the alternative distributions one may separate the leading term in the probability from the expectations that form the sum, enabling the investigation to concentrate on finer effects.

The method is useful since it does not rely on the ordering of the parameter set and it does not require the normal distribution. In many applications, some of them are presented in the book, a natural formulation of the model calls for the use of collections of random variables that are parameterized not by subsets of the real line. Frequently, the normal assumption may fit the limit in a central limit formulation but may not fit as a description of the extreme tail. In all such cases an alternative to the methods that are usually advocated in the literature are required. The method we present is such an alternative which we felt others may benefit from by knowing about.

This is why we wrote the book. But who is the target audience? This is a tough call. Even if I may state otherwise, the book requires a relatively advanced knowledge in probability as background, perhaps at the level of Durrett's book.³ Prior knowledge in statistics is an advantage. Indeed, there is an appendix that lists theorems and results and can be used as reference for the statements that are made in the book. Still, I guess that this book is not an easy read even for experts, and much less so for students.

With this warning in mind, I hope that the effort that is required in reading the book will be rewarding. Definitely, for an expert who wants to add yet another method to his toolbox but also for a student who wants to become an expert.

Seigmund D.O., Yakir B. Significance level in interval mapping. In *Development of Modern Statistics and Related Topics*, Series in Biostatistics, Volume 1. World Scientific Publishing, River Edge, NJ, 10–19 (2003).

Shi J., Siegmund D.O., Yakir B. Importance sampling for estimating p-values in linkage analysis. *JASA* **102**, 929–937 (2007).

Yakir B. On the average run length to false alarm in surveillance problems which possess an invariance structure. *Ann. Statist.* **26**, 1198–1214 (1998).

Yakir B. Approximation of the p-value in a multipoint linkage analysis using grandparent grandchild pairs and partially informative markers. *Nonlinear Anal.* **47**, 1973–1984 (2001).

Yakir B. Discussion on “Is average run length to false alarm always an informative criterion?” by Yajun Mei. *Sequential Analysis* **27**, 406–410 (2008).

³ Durrett R. *Probability: Theory and Examples* (2nd Edition). Duxbury Press, Belmont, CA (1995).

For such students, the book can be used as a basis for an advanced seminar. Reading chapters of the book can be used as a primer for a student who is then required to analyze a new problem that was not digested for him/her in the book. This is how I intend to use this book with my students.

The teacher can start such a course by discussing Chapters 1–4 that give the basic background and demonstrate the technique. Chapter 5 is more technical and can be skipped, unless the main interest is in the mathematical details. From the second part of the book it is probably recommended to go over Chapter 6, which is of an intermediate level of difficulty, and then read some or all of Chapters 7–10 depending on the interests of the teacher and the students and on the time constraints.

Acknowledgments

My first acknowledgments are to environments, and especially the people who enabled these environments. The first half of the book was written mainly in DSAP. I know of no better place to do this type of scientific work. I will always be grateful. The second place is MSC. Without them I do not know when, if at all, this book would have been written. Next, I would like to recognize the financial support that I got from the Israel Science Foundation (Grant No. 325/09) and from the US–Israel Binational Science Foundation (Grant No. 2006101). This support was instrumental for the development of the original work that led to the applications that are presented in the second part of the book.

Some of the people that gave me a helping hand I would like to mention by name. Unfortunately, I cannot give the names of the anonymous reviewers who made very useful suggestions on the first draft of the book and helped me improve it. But I can give the name of the editor from Higher Education Press, Liping Wang. Thanks to her this is a book and not just class notes. Also I would like to thank Yuval Nardi, Moshe Pollak, Ton Dieker and Nancy Zhang who coauthored with David Siegmund and myself some of the works that are related directly to the content of the book.

And finally there is David Siegmund. The work presented in this book is basically our joint work. The only reason that we do not share authorship is the fact that I wanted to dedicate this book to him as my modest contribution to the celebration of his career and his accomplishments and as an appreciation for what he gave me. It is not appropriate for a book to be dedicated to one of its authors. So here it is: this is for you, David.

February 2013

Benjamin Yakir, Jerusalem, Israel

Contents

Preface	xi
Acknowledgments	xv
Part I THEORY	1
1 Introduction	3
1.1 Distribution of extremes in random fields	3
1.2 Outline of the method	7
1.3 Gaussian and asymptotically Gaussian random fields	9
1.4 Applications	11
2 Basic examples	15
2.1 Introduction	15
2.2 A power-one sequential test	15
2.3 A kernel-based scanning statistic	24
2.4 Other methods	38
3 Approximation of the local rate	41
3.1 Introduction	41
3.2 Preliminary localization and approximation	43
3.2.1 Localization	43
3.2.2 A discrete approximation	46
3.3 Measure transformation	51
3.4 Application of the localization theorem	55
3.4.1 Checking Condition I*	57
3.4.2 Checking Condition V*	57
3.4.3 Checking Condition IV*	58
3.4.4 Checking Condition II*	59
3.4.5 Checking Condition III*	63
3.5 Integration	67

4	From the local to the global	71
4.1	Introduction	71
4.2	Poisson approximation of probabilities	72
4.3	Average run length to false alarm	78
5	The localization theorem	87
5.1	Introduction	87
5.2	A simplified version of the localization theorem	88
5.3	The localization theorem	90
5.4	A local limit theorem	95
5.5	Edge effects and higher order approximations	100
Part II	APPLICATIONS	103
6	Nonparametric tests: Kolmogorov–Smirnov and Peacock	105
6.1	Introduction	105
6.1.1	Classical analysis of the Kolmogorov–Smirnov test	106
6.1.2	Peacock’s test	108
6.2	Analysis of the one-dimensional case	109
6.2.1	Preliminary localization	110
6.2.2	An approximation by a discrete grid	111
6.2.3	Measure transformation	114
6.2.4	The asymptotic distribution of the local field and the global term	115
6.2.5	Application of the localization theorem and integration	117
6.2.6	Checking the conditions of the localization theorem	119
6.3	Peacock’s test	120
6.4	Relations to scanning statistics	123
7	Copy number variations	125
7.1	Introduction	125
7.2	The statistical model	127
7.3	Analysis of statistical properties	131
7.3.1	The alternative distribution	131
7.3.2	Preliminary localization and approximation	132
7.3.3	Measure transformation	132
7.3.4	The localization theorem and the local limit theorem	133
7.3.5	Checking Condition V*	137
7.3.6	Checking Condition II*	137
7.4	The false discovery rate	140
8	Sequential monitoring of an image	143
8.1	Introduction	143
8.2	The statistical model	146

8.3	Analysis of statistical properties	148
8.3.1	Preliminary localization	149
8.3.2	Measure transformation, the localization theorem, and integration	155
8.3.3	Checking the conditions of the localization theorem	157
8.3.4	Checking Condition V*	157
8.3.5	Checking Condition IV*	158
8.3.6	Checking Condition II*	159
8.4	Optimal change-point detection	161
9	Buffer overflow	165
9.1	Introduction	165
9.2	The statistical model	169
9.2.1	The process of demand from a single source	169
9.2.2	The integrated process of demand	171
9.3	Analysis of statistical properties	172
9.3.1	The large deviation factor	172
9.3.2	Preliminary localization	174
9.3.3	Approximation by a cruder grid	175
9.3.4	Measure transformation	179
9.3.5	The localization theorem	180
9.3.6	Integration	183
9.3.7	Checking the conditions of the localization theorem	184
9.3.8	Checking Condition IV*	184
9.3.9	Checking Condition V*	185
9.3.10	Checking Condition II*	185
9.4	Heavy tail distribution, long-range dependence, and self-similarity	186
10	Computing Pickands' constants	191
10.1	Introduction	191
10.1.1	The double-sum method	192
10.1.2	The method based on the likelihood ratio identity	193
10.1.3	Pickands' constants	195
10.2	Representations of constants	196
10.3	Analysis of statistical error	199
10.4	Enumerating the effect of local fluctuations	204
Appendix:	Mathematical background	209
A.1	Transforms	209
A.2	Approximations of sum of independent random elements	211
A.3	Concentration inequalities	214
A.4	Random walks	215
A.5	Renewal theory	215
A.6	The Gaussian distribution	216

x	CONTENTS	
	A.7 Large sample inference	217
	A.8 Integration	218
	A.9 Poisson approximation	219
	A.10 Convexity	220
	References	221
	Index	223

Part I

THEORY

1

Introduction

1.1 Distribution of extremes in random fields

The aim of this book is to present a method for analyzing the tail distribution of extreme values in random fields. A random field can be considered as a collection of random variables $\{X_t : t \in T\}$, indexed by a set of parameters T . The index set T may be quite complex. However, in the applications that we will analyze in this book it will typically turn out that T is a ‘nice’ subset of \mathbb{R}^d , the d -dimensional space of real numbers.

In some statistical applications one is interested in probabilities such as:

$$P\left(\sup_{t \in T} X_t \geq x\right),$$

the probability that the maximum of the random field exceeds a threshold x , for large values of x . There are only a few special cases in which the problem of computing such probabilities has an exact solution. In all other cases one is forced to use numerical methods, such as simulations, or to apply asymptotic approximations in order to evaluate the probability. This book concentrates on the application of the proposed method for producing asymptotic analytical expansions of the probability. Nonetheless, some elements in the method may, and have been, applied in order to simulate numerical evaluations more efficiently. An application that illustrates the usefulness of the method in the context of simulations is presented in the second part of the book.

As a motivating example consider scanning statistics. Scanning statistics are used in order to detect rare signals in an environment contaminated by random noise. For example, let us assume measurements that are taken in a one-dimensional environment. Each measurement is associated with a point in the environment and the points are equally spaced. For the most part, the expected values of the observations are fixed at some baseline level throughout