

# Off the Beaten Path Tutorial: Stochastic Processes and Simulations – Volume 1

Vincent Granville, Ph.D.  
[Data Shaping Solutions, LLC](#)

Anacortes, WA, February 2022

**Note:** External links (in blue) and internal references (in red) are clickable throughout this document. Key-words highlighted in orange are indexed; those in red are both indexed and in the glossary section.

## Contents

<b>About this Textbook</b>	<b>2</b>
Target Audience . . . . .	3
About the Author . . . . .	4
<b>1 Poisson-binomial or Perturbed Lattice Process</b>	<b>4</b>
1.1 Definitions . . . . .	5
1.2 Point Count and Interarrival Times . . . . .	6
1.3 Limiting Distributions, Speed of Convergence . . . . .	7
1.4 Properties of Stochastic Point Processes . . . . .	8
1.4.1 Stationarity . . . . .	8
1.4.2 Ergodicity . . . . .	8
1.4.3 Independent Increments . . . . .	8
1.4.4 Homogeneity . . . . .	9
1.5 Transforming and Combining Multiple Point Processes . . . . .	9
1.5.1 Marked Point Process . . . . .	9
1.5.2 Rotation, Stretching, Translation and Standardization . . . . .	9
1.5.3 Superimposition and Mixing . . . . .	10
1.5.4 Hexagonal Lattice, Nearest Neighbors . . . . .	11
<b>2 Applications</b>	<b>11</b>
2.1 Modeling Cluster Systems in Two Dimensions . . . . .	12
2.1.1 Generalized Logistic Distribution . . . . .	13
2.1.2 Illustrations . . . . .	14
2.2 Infinite Random Permutations with Local Perturbations . . . . .	16
2.3 Probabilistic Number Theory and Experimental Maths . . . . .	16
2.3.1 Poisson Limit of the Logistic-binomial Distribution, with Applications . . . . .	17
2.3.2 Perturbed Version of the Riemann Hypothesis . . . . .	19
2.4 Videos: Fractal Supervised Classification and Riemann Hypothesis . . . . .	21
2.4.1 Dirichlet Eta Function . . . . .	21
2.4.2 Fractal Supervised Classification . . . . .	23
<b>3 Statistical Inference, Machine Learning, and Simulations</b>	<b>23</b>
3.1 Model-free Tests and Confidence Regions . . . . .	24
3.1.1 Methodology and Example . . . . .	24
3.1.2 Periodicity and Amplitude of Point Counts . . . . .	28
3.1.3 A New Test of Independence . . . . .	29
3.2 Estimation of Core Parameters . . . . .	31
3.2.1 Intensity and Scaling Factor . . . . .	31
3.2.2 Model Selection to Identify $F$ . . . . .	32
3.2.3 Theoretical Values Obtained by Simulations . . . . .	33
3.3 Hard-to-Detect Patterns and Model Identifiability . . . . .	34
3.4 Spatial Statistics, Nearest Neighbors, Clustering . . . . .	35
3.4.1 Stochastic Residues . . . . .	35
3.4.2 Inference for Two-dimensional Processes . . . . .	35
3.4.3 Clustering Using GPU-based Image Filtering . . . . .	38

3.4.4	Black-box Elbow Rule to Detect Outliers and Number of Clusters . . . . .	40
3.5	Boundary Effect . . . . .	43
3.5.1	Quantifying some Biases . . . . .	44
3.5.2	Extreme Values . . . . .	45
3.6	Poor Random Numbers and Other Glitches . . . . .	47
3.6.1	A New Type of Pseudo-random Number Generator . . . . .	47
<b>4</b>	<b>Theorems</b>	<b>48</b>
4.1	Notations . . . . .	48
4.2	Link between Interarrival Times and Point Count . . . . .	48
4.3	Point Count Arithmetic . . . . .	49
4.4	Link between Intensity and Scaling Factor . . . . .	49
4.5	Expectation and Limit Distribution of Interarrival Times . . . . .	50
4.6	Convergence to the Poisson Process . . . . .	51
4.7	The Inverse or Hidden Model . . . . .	52
4.8	Special Cases with Exact Formula . . . . .	53
4.9	Fundamental Theorem of Statistics . . . . .	53
<b>5</b>	<b>Exercises, with Solutions</b>	<b>55</b>
5.1	Full List . . . . .	55
5.2	Probability Distributions, Limits and Convergence . . . . .	55
5.3	Features of Poisson-binomial Processes . . . . .	59
5.4	Lattice Networks, Covering Problems, and Nearest Neighbors . . . . .	61
5.5	Miscellaneous . . . . .	65
<b>6</b>	<b>Source Code, Data, Videos, and Excel Spreadsheets</b>	<b>69</b>
6.1	Interactive Spreadsheets and Videos . . . . .	70
6.2	Source Code: Point Count, Interarrival Times . . . . .	71
6.2.1	Compute $E[N(B)]$ , $\text{Var}[N(B)]$ and $P[N(B) = 0]$ . . . . .	72
6.2.2	Compute $E[T]$ , $\text{Var}[T]$ and $E[T^r]$ . . . . .	73
6.2.3	Produce random deviates for various $F$ 's . . . . .	74
6.2.4	Compute $F(x)$ for Various $F$ 's . . . . .	74
6.3	Source Code: Radial Cluster Simulation . . . . .	75
6.4	Source Code: Nearest Neighbor Distances . . . . .	75
6.5	Source Code: Detection of Connected Components . . . . .	79
6.6	Source Code: Visualizations, Density Maps . . . . .	81
6.6.1	Visualizing the Nearest Neighbor Graph . . . . .	81
6.6.2	Clustering and Density Estimation via Image Filtering . . . . .	82
6.7	Source Code: Production of the Videos . . . . .	85
6.7.1	Dirichlet Eta Function . . . . .	85
6.7.2	Fractal Supervised Clustering . . . . .	86
	<b>Glossary</b>	<b>89</b>
	<b>List of Figures</b>	<b>90</b>
	<b>List of Tables</b>	<b>90</b>
	<b>References</b>	<b>90</b>
	<b>Index</b>	<b>94</b>

## About this Textbook

This scratch course on stochastic processes covers significantly more material than usually found in traditional books or classes. The approach is original: I introduce a new yet intuitive type of random structure called perturbed lattice or Poisson-binomial process, as the gateway to all the stochastic processes. Such models have started to gain considerable momentum recently, especially in sensor data, cellular networks, chemistry, physics and engineering applications. I present state-of-the-art material in simple words, in a compact style, including new research developments and open problems. I focus on the methodology and principles, providing the reader with solid foundations and numerous resources: theory, applications, illustrations, statistical inference, refer-

ences, glossary, educational spreadsheet, source code, stochastic simulations, original exercises, videos and more.

Below is a short selection highlighting some of the topics featured in the textbook. Some are research results published here for the first time.

GPU clustering	Fractal supervised clustering in GPU (graphics processing unit) using image filtering techniques akin to neural networks, automated black-box detection of the number of clusters, unsupervised clustering in GPU using density (gray levels) equalizer
Inference	New test of independence, spatial processes, model fitting, dual confidence regions, minimum contrast estimation, oscillating estimators, mixture and surperimposed models, radial cluster processes, exponential-binomial distribution with infinitely many parameters, generalized logistic distribution
Nearest neighbors	Statistical distribution of distances and Rayleigh test, Weibull distribution, properties of nearest neighbor graphs, size distribution of connected components, geometric features, hexagonal lattices, coverage problems, simulations, model-free inference
Cool stuff	Random functions, random graphs, random permutations, chaotic convergence, perturbed Riemann Hypothesis (experimental number theory), attractor distributions in extreme value theory, central limit theorem for stochastic processes, numerical stability, optimum color palettes, cluster processes on the sphere
Resources	27 Exercises with solution expanding the theory and methods presented in the textbook, well documented source code and formulas to generate various deviates and simulations, simple recipes (with source code) to design your own data animations as MP4 videos – see ours on YouTube

This first volume deals with point processes in one and two dimensions, including spatial processes and clustering. The next volume in this series will cover other types of stochastic processes, such as Brownian-related and random, chaotic dynamical systems. The point process which is at the core of this textbook is called the Poisson-binomial process (not to be confused with a binomial nor a Poisson process) for reasons that will soon become apparent to the reader. Two extreme cases are the standard Poisson process, and fixed (non-random) points on a lattice. Everything in between is the most exciting part.

## Target Audience

College-educated professionals with an analytical background (physics, economics, finance, machine learning, statistics, computer science, quant, mathematics, operations research, engineering, business intelligence), students enrolled in a quantitative curriculum, decision makers or managers working with data scientists, graduate students, researchers and college professors, will benefit the most from this textbook. The textbook is also intended to professionals interested in automated machine learning and artificial intelligence.

It includes many original exercises requiring out-of-the-box thinking, and offered with solution. Both students and college professors will find them very valuable. Most of these exercises are an extension of the core material. Also, a large number of internal and external references are immediately accessible with one click, throughout the textbook: they are highlighted respectively in red and blue in the text. The material is organized to facilitate the reading in random order as much as possible and to make navigation easy. It is written for busy readers.

The textbook includes full source code, in particular for simulations, image processing, and video generation. You don't need to be a programmer to understand the code. It is well documented and easy to read, even for people with little or no programming experience. Emphasis is on good coding practices. The goal is to help you quickly develop and implement your own machine learning applications from scratch, or use the ones offered in the textbook. The material also features professional-looking spreadsheets allowing you to perform interactive statistical tests and simulations in Excel alone, without statistical tables or any coding. The code, data sets, videos and spreadsheets are available on my GitHub repository.

## About the Author

Vincent Granville, PhD is a pioneering data scientist and machine learning expert, co-founder of Data Science Central (acquired by a publicly traded company in 2020), former VC-funded executive, author and patent owner. Vincent's past corporate experience includes Visa, Wells Fargo, eBay, NBC, Microsoft, CNET, InfoSpace and other Internet startup companies (one acquired by Google). Vincent is also a former post-doct from Cambridge University, and the National Institute of Statistical Sciences (NISS). He is currently publisher at [DataShaping.com](https://DataShaping.com). He makes a living as an independent researcher working on stochastic processes, dynamical systems, experimental math and probabilistic number theory.

Vincent published in Journal of Number Theory, Journal of the Royal Statistical Society (Series B), and IEEE Transactions on Pattern Analysis and Machine Intelligence, among others. He is also the author of multiple books, including "Statistics: New Foundations, Toolbox, and Machine Learning Recipes", "Applied Stochastic Processes, Chaos Modeling, and Probabilistic Properties of Numeration Systems" with a combined reach of over 250,000, as well as "Becoming a Data Scientist" published by Wiley. For details, see my Google Scholar profile, [here](#).

## 1 Poisson-binomial or Perturbed Lattice Process

I introduce here one of the simplest point process models. The purpose is to illustrate, in simple English, the theory of point processes using one of the most elementary and intuitive examples, keeping applications in mind. Many other point processes will be covered in the next sections, both in one and two dimensions. Key concepts, soon to be defined, include:

Category	Description	Book sections
Top parameters	Intensity $\lambda$ – granularity of the process	<a href="#">4.4</a> , <a href="#">3.2.1</a>
	Scaling factor $s$ – quantifies point repulsion or mixing	<a href="#">3.1.1</a> , <a href="#">3.2.1</a>
	Distribution $F$ – location-scale family, with $F_s(x) = F(x/s)$	<a href="#">1.1</a> , <a href="#">3.2.2</a>
Properties	Stationarity and ergodicity	<a href="#">1.4</a> , <a href="#">5.3</a>
	Homogeneity and anisotropy	<a href="#">1.4.4</a>
	Independent increments	<a href="#">1.4.3</a> , <a href="#">3.1.3</a>
Core distributions	Interarrival times $T$	<a href="#">1.2</a> , <a href="#">4.2</a>
	Nearest neighbor distances	<a href="#">3.4</a> , <a href="#">5.4</a>
	Point count $N(B)$ in a set $B$	<a href="#">4.3</a> , <a href="#">5.3</a>
	Point distribution (scattering, on a set $B$ )	<a href="#">1.2</a>
Type of process	Marked point process	<a href="#">1.5.1</a>
	Cluster point process	<a href="#">2.1</a> , <a href="#">2.1.2</a>
	Mixtures and interlacings (superimposed processes)	<a href="#">1.5.3</a> , <a href="#">3.4.3</a>
Topology	Lattice space (index space divided by $\lambda$ )	<a href="#">2.1</a> , <a href="#">4.7</a>
	State space (where the points are located)	<a href="#">2.1</a>
	Index space (hidden space of point indices: $\mathbb{Z}$ or $\mathbb{Z}^2$ )	<a href="#">4.7</a> , <a href="#">2.2</a>
Other concepts	Convergence to stationary Poisson point process	<a href="#">1.3</a> , <a href="#">4.6</a>
	Boundary effects	<a href="#">3.5</a>
	Dimension (of the state space)	<a href="#">1.2</a>
	Model identifiability	<a href="#">3.3</a>

I also present several probability distributions that are easy to sample from, including logistic, uniform, Laplace and Cauchy. I use them in the simulations. I also introduce new ones such as the [exponential-binomial distribution](#) (the distribution of interarrival times), and a new type of [generalized logistic distribution](#). One of the core distributions is the [Poisson-binomial](#) with an infinite number of parameters. The Poisson-binomial process is named after that distribution, attached to the [point count](#) (a random variable) counting the number