

Limit Theorems as Blessing of Dimensionality: Neural-Oriented Overview

Vladik Kreinovich  and Olga Kosheleva 

University of Texas at El Paso, El Paso TX 79968, USA; vladik@utep.edu, olgak@utep.edu

* Correspondence: vladik@utep.edu (V.K.)

1 **Abstract:** As a system becomes more complex, at first, its description and analysis becomes more
 2 complicated. However, a further increase in the system’s complexity often makes this analysis
 3 simpler. A classical example is Central Limit Theorem: when we have a few independent sources
 4 of uncertainty, the resulting uncertainty is very difficult to describe, but as the number of such
 5 sources increases, the resulting distribution gets close to an easy-to-analyze normal one – and
 6 indeed, normal distributions are ubiquitous. We show that such limit theorems often make
 7 analysis of complex systems easier – i.e., lead to blessing of dimensionality phenomenon – for all
 8 the aspects of these systems: the corresponding transformation, the system’s uncertainty, and the
 9 desired result of the system’s analysis.

10 **Keywords:** limit theorems; curse and blessing of dimensionality; neural networks

11 1. Introduction: From Curse of Dimensionality to Blessing of Dimensionality

12 **First, a curse.** Often, the more we analyze a system, the more accurately we want to
 13 predict its behavior – the more factors we need to take into account, the more complex
 14 the system’s behavior.

15 In some cases, real-life data is intrinsically low-dimensional: most of the factors
 16 can be reduced to a few of them. However, in many other real-life situations, all these
 17 factors are important. As a result, as a system’s description becomes more complex,
 18 analyzing this system becomes more complicated. This phenomenon is known as *curse*
 19 *of dimensionality*.

20 **Then, a blessing.** Interestingly, often, a further increase in the system’s complexity often
 21 makes this analysis simpler. Following [1], we will call this phenomenon *blessing of*
 22 *dimensionality*.

23 **Example.** A classical example of this first-curse-then-blessing phenomenon is the joint
 24 effect of many random phenomena. When we know the probability distribution of each
 25 phenomenon, in principle, we can compute their joint effect – but, as the number of
 26 these phenomena becomes larger and larger, the corresponding computations become
 27 more and more complicated. At first glance, this is a classical example of the curse of
 28 dimensionality.

29 However, as the number of these phenomena increases further, we start seeing the
 30 effect of the Central Limit Theorem (see, e.g., [2]), according to which, under reasonable
 31 conditions, the joint effect of many small independent random phenomena is close to
 32 Gaussian. The resulting distribution becomes very close to the easy-to-analyze Gaussian
 33 distribution – and this is one of the main reasons why normal (= Gaussian) distributions
 34 are ubiquitous.

35 **Other examples.** In the last decade, many other examples of blessing-of-dimensionality
 36 appeared, both in the general analysis of complex systems (see, e.g., [1,3–7]) and, specifi-
 37 cally, in the analysis of neural networks; see, e.g. [8–11].

Citation: Kreinovich, V.; Kosheleva, O. Limit Theorems as Blessing of Dimensionality: Neural-Oriented Overview. *Entropy* **2021**, *1*, 0. <https://doi.org/>

Received:
 Accepted:
 Published:

Publisher’s Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Copyright: © 2021 by the authors. Submitted to *Entropy* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

38 **Are these lucky examples or a general trend?** At first glance, it may appear that all these
39 examples are lucky breaks in the dark world of curse-of-dimensionality phenomena. So,
40 a natural question is: is this pessimistic viewpoint correct – or blessing-of-dimensionality
41 results are ubiquitous?

42 **It is a general trend.** In this paper, we show that the above pessimistic viewpoint is –
43 well – unnecessarily pessimistic. Actually, as we will show, similar limit theorems are
44 ubiquitous – and their use can (and do) help in data processing in general – and, in
45 particular, when using neural networks to process data.

46 While most above-cited blessing of dimensionality results are related to a statistical
47 description of some phenomenon, we show that there are other limit theorems that are
48 related to non-random phenomena.

49 We also show that limit theorems help explain the surprising empirical success of
50 many techniques, from traditional neural networks to convex techniques and clustering.

51

52 **Caution: blessing-of-dimensionality is not a panacea.**

- 53 • The fact that limit theorems can explain some empirical successes does not mean, of
54 course, that these blessing-of-dimensionality results are the only reason for these
55 empirical successes: sometimes, as we have mentioned, the multi-dimensional data
56 is actually intrinsically low-dimensional.
- 57 • The fact that limit theorems *often* make data processing easier does not mean that as
58 the data complexity increases, the analysis *always* becomes simpler: many problems
59 remain complex. At present, there is no clear general understanding of when the
60 blessing of dimensionality occurs and where it does not occur. It would be nice to
61 find such an understanding.

62 **What we do in this paper.** In this paper, we review, in an expository mathematics format,
63 several published results (some of them our own) showing that limit theorems can
64 simplify the analysis of complex systems in general and neural networks in particular.

65 Our main interest is in applications to neural networks, so when a theorem has
66 such applications, we explicitly mention them – but we mention other applications as
67 well. The number of neural applications of limit theorems is, at present, not large, but
68 we hope that papers like this one – that explain how such theorems are successfully used
69 in other applications – will encourage interested readers to develop new applications of
70 these blessing-of-dimensionality results to neural networks.

71 The intended audience of this paper are readers with a conceptual understanding
72 of the mathematics involved, not necessarily with a specialist knowledge. Readers
73 interested in more detailed discussions and/or exact formulations and proofs are wel-
74 come to look at the corresponding papers listed in the bibliography. In these papers,
75 the corresponding discussions, formulations, and proofs are presented in all necessary
76 detail.

77 The general study of blessing-of-dimensionality phenomena has started only a few
78 decades ago, there are still more open problems than results – and available results are
79 mostly breakthroughs in different directions, not yet forming a very coherent picture.
80 Good news is that there are already many such results, and their applications already
81 over many areas. We hope that by listing these results and some of their applications,
82 we will encourage interested readers to get involved in the related research – and that,
83 together, we will make this phenomenon even more ubiquitous.

84 **How this paper is structured.** We start, in Section 2, with classifying sources of di-
85 mensionality into spatial and temporal. Such a distinction is well known in neural
86 network applications; in this section, we extend it to the general case of complex systems.
87 Section 3 deals with spatial dimensionality, of which the dimensionality correspond-
88 ing to the Central Limit Theorem is one of the examples. We start, in Subsection 3.1,
89 with a new application of the Central Limit Theorem. In Subsection 3.2, we consider

90 generalizations of Central Limit Theorem to other types of probability distributions. In
 91 Subsection 3.3, we consider limit theorems corresponding to the case when we do not
 92 know the corresponding probabilities, when we only know the set of possible values
 93 of the corresponding quantity or quantities. Subsection 3.4 lists related open questions.
 94 Finally, Section 4 deals with limit theorems related to temporal dimensionality.

95 2. Two Main Sources of Dimensionality: Spatial and Temporal

96 To provide an adequate analysis of the situation, let us first observe that in general,
 97 there are two main sources of dimensionality:

- 98 • First, at each moment of time, there is usually a large number of phenomena –
 99 located, in general, at different points in space – that need to be taken into account.
 100 Even if we use a few parameters to describe each of these phenomena, overall, we
 101 will need a very large number of parameters to describe all these phenomena – and
 102 thus, the dimensionality of the problems grows. We will call this dimensionality
 103 of *spatial origin*, or simply *spatial dimensionality*, for short. The above-mentioned
 104 Central Limit Theorem is a good example of spatial dimensionality.
- 105 • Also, there may be parameters describing the history of the analyzed phenomenon –
 106 which also affect its current state. What naturally comes to mind is that the values
 107 of physical quantities change with time. In some cases, we observe these changes
 108 and we can analyze the corresponding time series. In other cases, we only observe
 109 the final results of these changes: e.g., inside a sensor, the original value may be
 110 transformed many times, and what we get as a resulting signal is the result of all
 111 these past transformations. In yet other cases, what changes are the simulated values
 112 – e.g., when we apply iterative algorithms. We will call the resulting dimensionality
 113 of *temporal origin*, or simply *temporal dimensionality*.

114 And, of course, in many real-life phenomena, we have both spatial and temporal sources
 115 of dimensionality which are difficult to separate. A neural-related example of such
 116 phenomena is traveling waves; see, e.g., [12,13].

117 In this paper, we will mention the limit theorems related to both spatial and tempo-
 118 ral sources of dimensionality – and we hope that these results can be extended to the
 119 phenomena where both sources are intertwined.

120 *Comment.* Limit theorems are often somewhat complicated to understand and prove. In
 121 our experience, a better understanding of a complex multi-dimensional phenomenon
 122 is usually achieved if we consider easier-to-analyze few-dimensional particular cases
 123 or analogues. For limit theorems, a natural few-dimensional analogues are iterative
 124 methods in numerical mathematics, such as:

- Newton's iterative method

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

125 for finding the solution to the equation $f(x) = 0$ or

- the gradient descent method

$$x_i^{(k+1)} = x_i^{(k)} - \alpha \cdot \frac{\partial f}{\partial x_i} \Big|_{x=x^{(k)}}$$

126 for finding the minimum of a function $f(x)$; we mention this method, since back-
 127 propagation, the main way neural networks learn, is, from the mathematical view-
 128 point, exactly gradient descent – with additional computational simplifications; see,
 129 e.g., [14,15].

130 In both examples, convergence is not guaranteed, and the results explaining when there
 131 is convergence are often difficult to prove. However, what is much easier to prove is that

132 if there is a convergence, then the limit satisfies the desired property – e.g., for Newton’s
133 method the limit value x satisfies the property $f(x) = 0$. In some cases, the limit value
134 only satisfies part of this desired property: for example, for the gradient descent method,
135 the limit is always a stationary point but not necessarily the desired global minimum of
136 the objective function $f(x)$. Indeed:

- 137 • For Newton’s method, if $x^{(k)} \rightarrow x$, then, in the limit, we get $x = x - \frac{f(x)}{f'(x)}$, which
138 implies that $f(x) = 0$.
- 139 • For the gradient descent, if $x^{(k)} \rightarrow x$, then, in the limit, we get $x_i = x_i - \alpha \cdot \frac{\partial f}{\partial x_i}$,
140 which implies that $\frac{\partial f}{\partial x_i} = 0$. Thus, the limit point is always a stationary point,
141 which is a necessary (but, as is well known, not sufficient) condition for it being the
142 location of the minimum.

143 Similarly to these cases, in this paper, we will concentrate not so much on the conditions
144 under which the processes converge, but rather on the description of the limit cases *when*
145 there is convergence.

146 3. Dimensionality of Spatial Origin

147 As we have mentioned, the standard Central Limit Theorem is an example of what
148 we called dimensionality of spatial origin. While many consequences of this theorem are
149 well known, as we will show, there are many aspects of this theorem which still need
150 exploring. So, the first thing we will consider – in the first subsection of this section – is
151 what are the less known consequences of the Central Limit Theorem.

152 Of course, the limit distribution does not have to be normal: as we have mentioned,
153 the convergence to the normal distribution happens only under certain conditions. For
154 situations when these conditions are not satisfied, there are more general limit theorems.
155 Applications of these more general theorems – mostly to uncertainty quantification – is
156 what we will overview in the second subsection of this section.

157 All this assumes that we know the probability distributions that we are trying
158 to combine. But what if we do not know the probabilities, what if we only know the
159 corresponding range of possible values – and we do not know the probabilities of
160 different points from this range? This situation is discussed in the third subsection of
161 this section.

162 This section ends with related open questions.

163 3.1. Not-Well-Known Consequences of the Central Limit Theorem

164 **Why are many things in the world discrete?** Outside quantum physics, most physical
165 processes are continuous, most probability distributions are continuous – so what we
166 should observe should be continuous as well. However, in reality, many things in the
167 real world are discrete. We do not have weather continuously changing from sunny to
168 rain: most of the time, we either have a sunny day or a rainy day. Yes, it is possible to
169 have hybrid animals like mules, but most of the time, animals we see fall into one of the
170 precise categories.

171 In many specific examples, there is a specific explanation for this discreteness – e.g.,
172 Darwin’s Theory of Evolution explains that only mutations which are beneficial to the
173 individual survive, and all intermediate stages between two beneficial states become
174 extinct fast. However, the very fact that the same discreteness phenomenon appears in
175 many different application areas seems to be an indication that discreteness is a general
176 phenomenon that must have a general explanation.

177 Discreteness is observed in machine learning as well: when we use a neural network
178 (or any similar tool) for classification, what this network actually produces are continuous
179 numbers – that can be converted, e.g., to degrees to which the object belongs to different
180 categories. However, usually, we do not return these degrees to the user. What we

181 usually do at the end is select one of these categories (e.g., the most probable one) – and
 182 in most cases, this is exactly the desired classification, cat or dog, car or not-a-car, disease
 183 or healthy, and this is usually exactly what the users want.

184 This discreteness definitely helps when making decisions – instead of a continuum
 185 of possible values, we need to deal with only a few discrete ones. So, this discreteness
 186 can be viewed as an example of a blessing of dimensionality.

187 But why are we mentioning this discreteness? At first glance, it may seem to be
 188 unrelated to the Central Limit Theorem – which is all about the normal distribution,
 189 which is, of course, absolutely continuous. Interestingly, there is a relation. Let us
 190 describe it.

191 **This puzzling discreteness has been observed before.** Of course, we are not the first
 192 ones who noticed that, in spite of the the fact that many processes are continuous, what
 193 we observe is often discrete. For example, B. S. Tsirelson noticed in [16] that in many
 194 cases, when we reconstruct a signal from noisy data, and we assume that the resulting
 195 signal belongs to a certain class, the reconstructed signal is often an *extreme* point from
 196 this class – i.e., is one of the discrete extreme points. In other words, the result is as
 197 discrete as our assumptions allow. For example:

- 198 • when we assume that the reconstructed signal is monotonic, the reconstructed
 199 function is often (piece-wise) constant;
- 200 • if we additionally assume that the signal is one time differentiable, the result is
 201 usually one time differentiable but rarely twice differentiable, etc.

202 **Tsirelson’s explanation.** Out of many papers that mention the puzzling discreteness, we
 203 cited [16] – because this paper not only *mentions* the fact of discreteness, it also provides
 204 an *explanation* for this discreteness, and this explanation is closely related to the Central
 205 Limit Theorem (see also [17]).

206 Indeed, when we extract a signal from a mixture with Gaussian noise, then the
 207 *maximum likelihood* estimation (a traditional statistical technique; see, e.g., [2]) means that
 208 out of all possible signals from the given class of signals, we look for the signal which
 209 is the closest (in the least squares — i.e., in effect, Euclidean – metric) to the observed
 210 “signal + noise” combination.

In particular, if the signal is determined by finitely many (say, d) parameters, we
 must look for a signal $\vec{s} = (s_1, \dots, s_d)$ from the a priori set $A \subseteq \mathbb{R}^d$ that is the closest (in
 the usual Euclidean sense) to the observed values

$$\vec{o} = (o_1, \dots, o_d) = (s_1 + n_1, \dots, s_d + n_d),$$

211 where n_i denotes the (unknown) values of the noise.

Since the noise is Gaussian, we can conclude that the average value of $(n_i)^2$ is close
 to σ^2 , where σ is the standard deviation of the noise. In other words, we can conclude
 that

$$(n_1)^2 + \dots + (n_d)^2 \approx d \cdot \sigma^2.$$

In geometric terms, this means that the distance

$$\sqrt{\sum_{i=1}^d (o_i - s_i)^2} = \sqrt{\sum_{i=1}^d n_i^2}$$

212 between \vec{s} and \vec{o} is $\approx \sigma \cdot \sqrt{d}$. Let us denote this distance $\sigma \cdot \sqrt{d}$ by ε .

213 For simplicity of explanation, let us consider the case when $d = 2$, and when A is a
 214 convex polygon. When the point \vec{o} corresponding to observations is itself inside the set
 215 A , then this point is its own closest point in the set A . Let us consider the case when the
 216 point \vec{o} is outside the set A . We can divide all points \vec{o} which are outside the set A and

217 which are ε -close to A into several zones depending on what part of A is the closest to \vec{o} :
 218 one of the *sides* (1-D faces), or one of the *vertices*.

219 Geometrically, the set of all points \vec{o} for which the closest point $a \in A$ belongs to the
 220 *side* e is bounded by the straight line segments orthogonal (perpendicular) to e . The total
 221 length of this set is therefore equal to the length of this particular side; hence, the total
 222 length of the set of all the points that are the closest to the sides is equal to the *perimeter*
 223 of the polygon. This total length thus does not depend on ε at all.

224 On the other hand, the overall length of the set of all the points \vec{o} at the distance ε
 225 from A grows with the increase in ε ; this length grows approximately as the circumfer-
 226 ence of a circle, i.e., as $\text{const} \cdot \varepsilon$.

227 When ε increases, the (constant) perimeter of the polygon A is a vanishing part of
 228 the overall length. Hence, for large ε :

- 229 • the fraction of the points that are the closest to one of the sides tends to 0, while
- 230 • the fraction of the points \vec{o} for which the *closest* point from the set A is one of A 's
 231 *vertices* tends to 1.

232 Thus, with high probability, the reconstructed signal corresponds to one of the vertices
 233 (extreme points) of the set A .

234 Similar arguments can be repeated for any dimension d . For the same noise level σ ,
 235 when d increases, the distance $\varepsilon = \sigma \cdot \sqrt{d}$ also increases, and therefore, for large d , for
 236 "almost all" observed points \vec{o} , the reconstructed signal is one of the extreme points of
 237 the *a priori* set A .

238 Much less probable is that the reconstructed signal \vec{s} belongs to the 1-dimensional
 239 face of the set A , even less probable that \vec{s} belongs to a 2-D face, etc.

240 **Methodological consequence.** So, when the dimension increases, we have a clear
 241 example of blessing of dimensionality: instead of having to consider a continuum of
 242 possible states, we only have to deal with a much smaller discrete set of extreme points –
 243 vertices of the corresponding polyhedron.

244 So, all observed phenomena falls into a few clusters – exactly as we observe in many
 245 cases.

246 *Comment.* This idea helps even in the quantum case. Namely, in quantum physics,
 247 there is a known paradox formulated by Schroedinger himself (the author of the main
 248 equation of quantum physics): while in quantum physics, we can have a superposition
 249 of any two states, how come we never see a superposition of two macro-states, e.g., of
 250 the state in which a cat is alive and the state in which the same cat is dead? This is indeed
 251 a serious problem, it was one of the reasons why Einstein did not believe that quantum
 252 physics is an adequate description of reality; see, e.g., [18–20].

253 Strictly speaking, this is not a paradox in the purely logical sense of the word – it is
 254 just a contradiction between our intuition and the predictions of quantum theory. Many
 255 features of quantum physics are counter-intuitive, but usually, such counter-intuitive
 256 features are about the micro-world of elementary particles, not about the usual macro-
 257 size objects. The above idea makes this contradiction less troubling, because it implies
 258 that with very high probability, we will observe one of the two original states and not
 259 their convex combination (i.e., in this case, not their superposition).

260 **Resulting discreteness is only approximate.** Of course, as with every probabilistic
 261 phenomenon, the above conclusion about discreteness is only approximate: we do not
 262 necessarily get one of the vertices, we get a point which is *close* to one of the vertices.
 263 This is why we did not write that all observed phenomena *coincide* with one of the few
 264 cases – we wrote that all observed phenomena fall into a few *clusters*. Within each cluster,
 265 we still have continuous changes – e.g., we can have cats of different length, different
 266 weight, etc.

267 3.2. *Uncertainty Quantification and Probabilistic Limit Theorems – Including Theorems Beyond*
 268 *Normal Distributions*

269 **Need for data processing.** What are the main objectives of science and engineering? We
 270 want to *understand* the world – i.e., to learn the values of the quantities that characterized
 271 the current state of the world. We want to *predict* the future state of the world – i.e.,
 272 we want to predict the future values of the corresponding quantities. And finally, we
 273 want to *change* the world – we want to find the design parameters that satisfy given
 274 specifications, we want to find the control values that will lead a system to the desired
 275 state, etc.

276 Some quantities that describe the world we can directly measure: e.g., the distance
 277 between two houses on the same street. For many other quantities, we cannot measure
 278 them directly: e.g., the distance to a nearby star. And we clearly cannot directly measure
 279 the future values of the quantities or the adequate value of control parameters. All these
 280 quantities have to be estimated based on the known information about the world – i.e.,
 281 based on the results of measuring some measurable quantities.

282 To estimate a desired quantity y , we need to know the relation $y = f(x_1, \dots, x_n)$
 283 between this quantity and measurable quantities x_1, \dots, x_n . Sometimes, we know an
 284 explicit analytical expression for this relation. In many other cases, we just know an
 285 algorithm that computes y from the values x_i . This algorithm can include a numerical
 286 solution of a complex system of non-linear differential equations – as when we predict
 287 tomorrow’s weather. The algorithm can also be a neural network trained to estimate the
 288 desired value y based on the known values x_1, \dots, x_n .

289 **Need for uncertainty quantification.** Whether we use neural networks or other algo-
 290 rithms for data processing, the inputs to all these algorithms are real numbers. These
 291 real numbers usually come from measurements, and measurements are never absolutely
 292 accurate; see, e.g., [21]. There is always noise. As a result, the measurement results \tilde{x}_i are,
 293 in general, somewhat different from the actual (unknown) values x_i of the corresponding
 294 quantities, and the difference $\Delta x_i \stackrel{\text{def}}{=} \tilde{x}_i - x_i$ – known as *measurement error* – is, in general,
 295 different from 0. So, when we apply the data processing algorithm f to the measurement
 296 results, the algorithm’s output $\tilde{y} = f(\tilde{x}_1, \dots, \tilde{x}_n)$ is, in general, different from the value
 297 $y = f(x_1, \dots, x_n)$ that we would have obtained if we knew the actual values x_i .

298 In practice, it is important to know how close is our estimate \tilde{y} to the desired value
 299 y , i.e., in other words, how big can the difference $\Delta y \stackrel{\text{def}}{=} \tilde{y} - y$ be. For example, suppose
 300 that we are prospecting for oil, and our estimate \tilde{y} for the amount of oil y in the given
 301 region is 150 million ton. Then, if the accuracy is ± 10 million tons, this estimate is good
 302 news, and we can start exploiting this region. On the other hand, if it is 150 ± 200 , then
 303 maybe there is no oil at all, so before we invest a lot of money into digging deep wells,
 304 we better perform more measurements to make sure that this money will not be wasted.

305 Estimating Δy is one the most important aspects of *uncertainty quantification*.

Possibility of linearization. We are interested in estimating the quantity

$$\Delta y = f(\tilde{x}_1, \dots, \tilde{x}_n) - f(x_1, \dots, x_n) = f(\tilde{x}_1, \dots, \tilde{x}_n) - f(\tilde{x}_1 - \Delta x_1, \dots, \tilde{x}_n - \Delta x_n).$$

Measurements are usually reasonable accurate, so the measurement errors Δx_i are relatively small. For small values Δx_i , their squares $(\Delta x_i)^2$ are much smaller than the values themselves – and can therefore be usually safely ignored. For example, if $\Delta x_i \approx 10\%$, then $(\Delta x_i)^2 \approx 1\% \ll \Delta x_i$. Thus, a reasonable idea is to expand the above expression for Δy in Taylor series and ignore terms which are quadratic (or of higher order) in terms of the measurement errors Δx_i . As a result, we get a linear dependence:

$$\Delta y \approx \sum_{i=1}^n c_i \cdot \Delta x_i, \text{ where } c_i \stackrel{\text{def}}{=} \frac{\partial f}{\partial x_i}.$$

307 *Comment.* This linearization – replacing the generic dependence with a linear one – is a
 308 usual idea in applications. Actually, it one of the main ideas in many applications; see,
 309 e.g., [22].

310 **Here, the Central Limit Theorem can help.** Let us first consider an important case –
 311 typically described in textbooks – when we know the probability distribution of each
 312 measurement error Δx_i . Usually, each measuring instrument is *calibrated* – if it has a *bias*,
 313 i.e., if the mean value $E[\Delta x_i]$ of the measurement error is not 0, we simply subtract this
 314 mean value from all the measurement results and thus, reduce it to 0.

315 In many practical applications, the number n of inputs is large, and the role of
 316 each of these inputs is relatively small. For example, one of the important data when
 317 prospecting for oil is seismograms – several-times-a-second recordings of the seismic
 318 signal. There are thousands of the corresponding values, and the effect of each indi-
 319 vidual value of the result of data processing is indeed small. The measurement errors
 320 corresponding to different measurements are usually reasonably independent. Thus,
 321 we are under the condition of the Central Limit Theorem – so we can conclude that the
 322 desired estimation error Δy is normally distributed.

A normal distribution is uniquely determined by its mean μ and its standard deviation σ . When each measurement error Δx_i has mean value 0, the mean value of their linear combination Δy is also 0, and the variance σ of this linear combination can be determined from the known fact that the variance of the sum of independent random variables is equal to the sum of variances:

$$\sigma^2 = \sum_{i=1}^n c_i^2 \cdot \sigma_i^2.$$

323

How can we actually estimate σ ? In principle, we can directly use the above formula to estimate the standard deviation σ of the approximation error Δy . The main computational difficulty is that the data processing algorithm f is usually very complicated (especially in case of neural networks), so it is not possible to compute the partial derivatives analytically. We can, however, use the fact that a partial derivative is defined as the limit of the ratios

$$\frac{\partial f}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(\tilde{x}_1, \dots, \tilde{x}_{i-1}, \tilde{x}_i + h, \tilde{x}_{i+1}, \dots, \tilde{x}_n) - \tilde{y}}{h},$$

and thus, for a sufficiently small h , the value of the ratio is very close to the desired partial derivative. Thus, we can estimate c_i as

$$c_i \approx \frac{f(\tilde{x}_1, \dots, \tilde{x}_{i-1}, \tilde{x}_i + h, \tilde{x}_{i+1}, \dots, \tilde{x}_n) - \tilde{y}}{h}.$$

324 The problem with this idea is that it takes too long. Indeed, if we have several
 325 thousand inputs, then, to compute all the corresponding values c_i , we need to call the
 326 data processing algorithm f (which often takes hours to compute) $n + 1$ times: one time
 327 to compute \tilde{y} and n time to compute the corresponding n ratios c_i . For several thousand
 328 inputs, this is not realistic.

Good news is that we can instead use Monte-Carlo techniques: instead of computing n partial derivatives, we can simply emulate, certain number of times K , measurement errors $\delta x_i^{(k)}$ which are normally distributed with standard deviation σ_i , and compute the differences

$$\delta y^{(k)} = \tilde{y} - f(\tilde{x}_1 - \delta x_1^{(k)}, \dots, \tilde{x}_n - \delta x_n^{(k)}).$$

329 By the same logic as before, the differences $\delta y^{(k)}$ are normally distributed with the
 330 desired standard deviation σ . Thus, from a sample of K values, we can estimate σ with
 331 accuracy $\approx 1/\sqrt{K}$ [2]. So, if we want to estimate σ with relative accuracy $1/\sqrt{K} \approx 20\%$, it

332 is sufficient to call the algorithm f $K = 25$ times – which is much smaller than thousands
333 needed for exact estimation.

334 **So what?** Why are we spending so much time on the ideas that are well known to many
335 readers? Because this will prepare readers to something that – unfortunately – not too
336 many readers know: that we can use limit theorems beyond normal distributions to
337 cover other realistic cases of uncertainty quantification.

338 **Need for interval uncertainty.** In the previous text, we assumed that for each measure-
339 ment, we know the probability distribution of the corresponding measurement error.
340 The usual way to find this distribution is to *calibrate* the given measuring instrument
341 (MI), i.e., to compare its results with the results of a “standard” (= much more accurate)
342 measuring instrument. Since the standard measuring instrument (SMI) is much more
343 accurate than the one we are calibrating, we can safely ignore SMI’s measurement errors
344 (in comparison with MI’s measurement errors), and take the results measured by SMI as
345 true values.

346 However, there are two important cases when calibration is not done. The first is
347 the case of state-of-the-art measurements, when the MI that we have is the best there is.
348 It would be great if near the Hubble telescope, there would fly a 5 times more accurate
349 instrument for measuring the stars’ locations, but this telescope is the best we have.
350 Similarly, in geophysics, oil prospecting companies use the best measuring instruments
351 they can find – these instruments are expensive, but digging a well in the location where
352 there is no oil would be much more expensive. In this case, there is no SMI to compare,
353 so we cannot calibrate our MI.

354 Another case is manufacturing and other practical applications. In this case, in
355 principle, we can calibrate every single measuring instrument and determine its prob-
356 ability distribution. However, nowadays, many sensors are cheap – e.g., kids playing
357 with robots buy distance sensors for a few bucks. However, calibrating a sensor means
358 utilizing a standard measuring instrument, which is usually much more expensive to
359 use. The companies usually cannot afford to calibrate all their sensors. Instead, we
360 have to rely on the information provided by the manufacturers of the corresponding
361 measuring instruments.

362 The manufacturer of the MI also has the option to calibrate it – but since this
363 calibration costs a lot, the calibrated sensors, with certified probability distributions of
364 measurement errors, cost much more. It is much cheaper to buy a sensor for which
365 only the minimum of necessary information is provided. In practice, this means that the
366 only information that we have about the measurement error Δx is an upper bound Δ
367 on its absolute value: $|\Delta x| \leq \Delta$. (At least such an upper bound needs to be provided –
368 otherwise, it is not a measuring instrument, it is a wild guess.)

369 Once we know the upper bound Δ_i on the absolute value $|\Delta x_i| = |\tilde{x}_i - x_i|$ of each
370 measurement error, then, based on the measurement result \tilde{x}_i , the only information we
371 gain about the actual (unknown) value x_i of the corresponding quantity is that this value
372 belongs to the interval $[x_i, \bar{x}_i] \stackrel{\text{def}}{=} [\tilde{x}_i - \Delta_i, \tilde{x}_i + \Delta_i]$. Because of this fact, such a situation
373 is known as *interval uncertainty*.

374 **Is the corresponding distribution Gaussian?** If we carefully eliminated all major
375 sources of measurement error, then only small factors remain that affect the measure-
376 ment error. Thus, due to the Central Limit Theorem, we can safely conclude that the
377 distribution of the measurement error is close to Gaussian. Will that help? Not really:
378 since we did not do the calibration, we do not know what is the bias. In principle, the
379 bias can take any value from $-\Delta_i$ and Δ_i , so the fact that we have a normal distribution
380 will not decrease the interval of uncertainty.

Uncertainty quantification: case of interval uncertainty. Under interval uncertainty,
the only thing we can conclude about the value $y = f(x_1, \dots, x_n)$ that we would have
obtained if we used the actual (unknown) values of the quantities x_i is that it belongs

to the *range* $[\underline{y}, \bar{y}]$ of possible values of the function f when x_i are in the corresponding intervals:

$$[\underline{y}, \bar{y}] = \{f(x_1, \dots, x_n) : x_i \in [\underline{x}_i, \bar{x}_i] \text{ for all } i\}.$$

381 The problem of computing this interval is known as the problem of *interval computation*;
382 see, e.g., [23,24].

383 In general, this problem is NP-hard [25] – which means that, unless $P = NP$ (which
384 most computer scientists do not believe to be possible), no feasible algorithm is possible
385 for solving all particular cases of this problem. However, in the linearized case, a feasible
386 algorithm *is* possible. Indeed, since the expression $\sum_i c_i \cdot \Delta x_i$ is linear (thus monotonic) in
387 the variables Δx_i , its largest value is attained:

- 388 • for $c_i > 0$, when the value Δx_i is the largest, i.e., when $\Delta x_i = \Delta_i$, and
- 389 • for $c_i < 0$, when the value Δx_i is the smallest, i.e., when $\Delta x_i = -\Delta_i$.

Thus, the largest possible value Δ of Δy is equal to

$$\Delta = \sum_{i=1}^n |c_i| \cdot \Delta_i.$$

390 Similarly, one can easily show that the smallest possible value of Δy is equal to $-\Delta$.

391 **How to estimate uncertainty in the interval case.** How can we compute this sum Δ ?
392 We can directly use this formula – i.e., use numerical differentiation to compute all the
393 partial derivatives c_i and then compute the sum. However, as we have mentioned earlier,
394 in many practical situations, this approach is not realistic. What can we do?

Another limit distribution comes to the rescue. As we have mentioned, the convergence to a normal distribution only happens under certain conditions. In other cases, we may have convergence to other so-called *infinitely divisible* distributions [2]. One of such distributions is the *Cauchy distribution*, in which the probability density $\rho(x)$ has the following form:

$$\rho(x) = \text{const} \cdot \frac{1}{1 + \left(\frac{x}{\Delta}\right)^2},$$

395 for some parameter Δ .

An important feature of the Cauchy distribution is that if we have several independent Cauchy distributed random variables r_i with parameters Δ_i , then their linear combination $\sum_i c_i \cdot r_i$ is also Cauchy distributed, with parameter $\Delta = \sum_i |c_i| \cdot \Delta_i$ – which is exactly the value that we want to compute. This feature leads to the following Monte-Carlo method for computing Δ : we emulate, certain number of times K , measurement errors $\delta x_i^{(k)}$ which are Cauchy distributed with parameters Δ_i , and compute the differences

$$\delta y^{(k)} = \tilde{y} - f\left(\tilde{x}_1 - \delta x_1^{(k)}, \dots, \tilde{x}_n - \delta x_n^{(k)}\right).$$

396 Then, due to the above feature, the differences $\delta y^{(k)}$ are Cauchy distributed with the
397 desired parameter Δ . Thus, to a sample of K values, we can apply, e.g., the maximum
398 likelihood method [2], and thus estimate Δ with accuracy $\approx 1/\sqrt{K}$. Similarly to the case
399 of normal distributions, this drastically speeds up computations: if we want to estimate
400 Δ with relative accuracy 20%, it is sufficient to call the algorithm f 25 times – which is
401 much smaller than thousands of times needed for exact estimation.

402 This method has been successfully used in many applications; see, e.g., [26].

403 *Comment.* Note that, in contrast to many simulation techniques, the use of Cauchy
404 distribution in interval-related uncertainty quantification is *not* a realistic simulation:

- 405 • the actual measurement error is always located *inside* the interval $[-\Delta, \Delta]$, while

- 406 • the Cauchy-distributed random variable has a non-zero probability to be anywhere,
407 in particular, *outside* the interval.

408 3.3. What If We Have No Information about Probabilities

409 **Formulation of the problem.** What if we know that the disturbance $x = (x_1, \dots, x_n)$
410 is a joint effect of several independent small ones: $x = x^{(1)} + \dots + x^{(N)}$, where about
411 each component $x^{(i)}$, we only know the set $X^{(i)}$ of its possible values – and we do not
412 have any information about probabilities of different points within each set. The only
413 constraint is that all the points from each set $X^{(i)}$ are small, i.e., that for some small
414 values $\varepsilon > 0$, the length $\|x^{(i)}\|$ of each vector $x^{(i)} \in X^{(i)}$ does not exceed ε . We will call
415 such sets ε -small.

416 In this case, the set X of all possible values of the sum x is the set of all possible
417 sums $x^{(1)} + \dots + x^{(N)}$, where $x^{(i)} \in X^{(i)}$ for all i . In mathematics, the set of all such
418 sums is known as the *Minkowski sum* of the sets $X^{(i)}$. The Minkowski sum is usually
419 denoted by $X^{(1)} + \dots + X^{(N)}$.

420 What can we say about such set X ?

421 **1-D case.** The 1-D case $n = 1$ was studied in [27]. This paper showed that if a set X is the
422 Minkowski sum of several ε -small closed sets, then it is ε -close to some interval $I = [a, b]$,
423 i.e.:

- 424 • every point from the set X is ε -close to some point from the interval I , and
425 • every point from the interval I is ε -close to some point from the set X .

426 In the limit $\varepsilon \rightarrow 0$, we conclude that the Minkowski sum tends to the interval.

427 To be more precise, the following results were proven:

428 **Theorem 1.** *If a set $S \subseteq \mathbb{R}$ is a Minkowski sum of δ -small closed sets, then S is δ -close to an*
429 *interval.*

430 **Theorem 2.** *If a set $S \subseteq \mathbb{R}$ can be, for every $\delta > 0$, represented as a Minkowski sum of finitely*
431 *many δ -small closed sets, then S is an interval.*

432 *Comment.* This limit theorem is similar, in formulation, to the Central Limit Theorem
433 and its generalizations: it shows that if a quantity can be represented as the sum of
434 many small components, then the set of all possible values of this quantity is close to an
435 interval – and the smaller the components, the closer is the resulting set to an interval.

436 Similarly to the fact that the original Central Limit Theorem explains the real-life
437 ubiquity of normal distributions, this limit theorem explains the ubiquity of interval
438 uncertainty; see, e.g., [21,23,24].

439 **General case.** It is well known that every convex set X containing 0 can be represented,
440 for every $\varepsilon > 0$, as a Minkowski sum of ε -small sets: indeed, it is sufficient to take
441 $X^{(i)} = N^{-1} \cdot X$ for a sufficiently large N , then:

- 442 • the inclusion $X \subseteq X^{(1)} + \dots + X^{(N)}$ follows from the fact that each element x can
443 be represented as the sum $x = N^{-1} \cdot x + \dots + N^{-1} \cdot x$; and
444 • the opposite inclusion $X^{(1)} + \dots + X^{(N)} \subseteq X$ follows from the fact that the set X is
445 convex and thus, once the elements $x^{(1)}, \dots, x^{(N)}$ belong to this set, their convex
446 combination $N^{-1} \cdot x^{(1)} + \dots + N^{-1} \cdot x^{(N)}$ also belongs to X .

447 Whether the opposite is true – i.e., whether only convex sets can be represented as sums
448 of small sets – remained an open problem. This problem – first formulated in [27] – was
449 resolved in [28], where the following result was proven:

450 **Theorem 3.** *If a set $X \subseteq \mathbb{R}^n$ can be represented, for each $\varepsilon > 0$, as a Minkowski sum of ε -small*
451 *closed sets, then this set X is convex.*

452 To be more precise, this paper proved the following result:

453 **Theorem 4.** *For every $\gamma > 0$, if a set $X \subset \mathbb{R}^n$ of diameter < 1 is δ -close to a Minkowski sum of*
454 *sets of diameter $\leq \varepsilon$, then X is γ -close to a convex set, for $\delta = \gamma/3$ and $\varepsilon = \gamma^2/(20n)$.*

455 *Comment.* This limit theorem explains the ubiquity of convex set in real-life problems.
 456 This is very good news, since it is known that convexity makes many computational
 457 problems easier to solve; see, e.g., [29].

458 3.4. Important Open Questions

459 **What if we only have partial information about probabilities?** In the above, we first
 460 considered cases where we know the probability distributions of the aggregated factors
 461 before moving to those in which when we only know the ranges, and we have no
 462 information about the probability of different values from these ranges. These are two
 463 extreme situations – either we know everything about the probabilities, or we have no
 464 information about these probabilities at all. In practice, we often have intermediate
 465 situations, when we have *partial* information about the probabilities. It is therefore
 466 desirable to extend the limit results from both extreme cases to the such intermediate
 467 situations as well.

468 **Possible approach and natural generalizations of the Central Limit Theorem.** When
 469 we know all the probabilities, then for uncertainty quantification, we can use Monte-
 470 Carlo approach with normal distributions. When we only know the upper bounds, we
 471 can use Cauchy distributions. What if for some components, we know the probabilities,
 472 and for others, we only know bounds? The resulting random variable is the sum of two
 473 partial sums, for which the first partial sum can be handled by the normal distribution,
 474 while the second partial sum can be handled by the Cauchy distribution. In this case,
 475 it seems reasonable to use the distributions corresponding to the sum of normally and
 476 Cauchy distributed random variables.

477 The family of such distributions is also a natural limit – the limit of sums in which
 478 the first partial sum tends to normal distribution and the second partial sum tends to
 479 the Cauchy one. Such mixed distributions are not covered by the usual limit theorems,
 480 which only consider 2-parametric limit families of probability distributions: e.g., a
 481 normal distribution is determined by two parameters – the mean and standard deviation
 482 of the normal distribution. Sums would require more parameters: we need mean and
 483 standard deviation of the normal part and the parameter Δ of the Cauchy part.

484 Possible generalizations of the traditional limit theorems to such multi-parametric
 485 families have been analyzed in [30]. It turns out that, in general, in this case, the resulting
 486 distribution is equivalent to the distribution of the sum of several different infinitely
 487 divisible distributions: e.g., to the sum of normally and Cauchy distributed variables. So
 488 maybe other distributions of this type can be used for uncertainty quantification in other
 489 cases when we only have partial information about probabilities?

What if we are interested in the extreme case? Very often, we are interested in the
 extreme case: e.g., when we design a bridge, we want it to withstand the strongest
 possible winds that can happen in this area. In such situations, we are interested not
 in the summary effect of several random variables, but rather in the largest value
 $x = \max(x_1, \dots, x_n)$ of several random variables x_i – e.g., variables describing the wind
 on different days. When all these variables are identically distributed, then, similarly
 to the Central Limit Theorem, we have a finite-parametric family of distributions that
 represents the distribution of such extreme events; see, e.g., [31–38]. Such results are
 known as *Extreme Value Theory*. The most widely used result is that if the random vari-
 ables x_i are independent and identically distributed, then, under reasonable conditions,
 as n increases, the cumulative distribution function of the maximum x of these variables
 tends to one of the three distribution functions: Gumbel law

$$F(x) = \exp\left(-\exp\left(-\frac{x-b}{a}\right)\right),$$

Fréchet law

$$F(x) = \exp\left(-\left(\frac{x-b}{a}\right)^{-\alpha}\right) \text{ for } x > b,$$

and Weibull law

$$F(x) = \exp\left(-\left|\frac{x-b}{a}\right|^\alpha\right) \text{ for } x < b.$$

490 This result is actively used in practice, e.g., in reliability engineering, to estimate the
491 probability of an extreme event.

492 The above result holds when all the variables x_i are identically distributed. In reality,
493 the distributions of the corresponding values x_i are, in general, somewhat different. So,
494 a natural question is: can we extend the Extreme Value Theory to such more general
495 case? A similar generalization is possible for the Central Limit Theorem: it holds for the
496 sum $x = x_1 + \dots + x_n$ even when the distributions of different variables x_i are different.
497 However, no such extension is known for the Extreme Value Theory. The absence of such
498 general extension is not caused by our inability to prove the corresponding result: it can
499 be shown that, if we simply remove the restriction that all variables x_i are identically
500 distributed, then the set of all limit distributions is no longer finite-dimensional; see [39].

501 Due to the practical importance of the Extreme Value Theory, an important question
502 emerges: since in a *general* case, we have an infinite-dimensional family of limit distribu-
503 tions, can we find *specific* cases when distributions are different but a finite-dimensional
504 family of limit distributions is still possible?

505 4. Dimensionality of Temporal Origin

Case study. Let us consider the case of a simple hardware sensor, in which the input
 x – e.g., intensity of light – generates a signal that goes through multiple layers until it
produces the final electric signal. When passing through these layers, the signal under-
goes a sequence of transformations. These transformations are, in general, nonlinear. In
mathematical terms, this means that the resulting transformation $f(x)$ of the original
real value x to the 1-D sensor output $f(x)$ is a composition of several different nonlinear
functions

$$f(x) = f_n(f_{n-1}(\dots f_2(f_1(x)) \dots)).$$

506 We can consider the sensor as a whole, with the transformation function $f(x)$. We
507 can divide it into several layers and consider the overall value-to-signal transformation
508 $f(x)$ as a composition of transformations corresponding to different layers. Each of these
509 layers can be viewed as several sub-layers, so the corresponding value n can be very
510 large – and transformations $f_i(x)$ corresponding to all these very thin sub-layers are
511 close to identity $f_i(x) \approx x$.

512 In the Central Limit Theorem, we took into account that the random variable x is
513 equal to the sum $x = x_1 + \dots + x_n$ of a large number of small independent random
514 variables, and we used the fact that under reasonable conditions, in the limit when $n \rightarrow$
515 ∞ , the distribution of this sum tends to a distribution from a known finite-parametric
516 family – namely, to a normal distribution. The limit means that when n is large, the
517 distribution of the sum x is close to Gaussian.

518 In our case, we consider a composition of a large number n of functions $f_i(x)$
519 which are close to identity. It is reasonable to look for situations in which, under some
520 conditions, when n increases, such compositions would also tend to functions from some
521 finite-parametric family. How can we describe the corresponding limit functions?

522 **Let us formulate this idea in precise terms.** As we have mentioned earlier, in this paper,
523 we do not focus on *conditions* when there is a convergence, we only focus on the resulting
524 limit. In line with this approach, let us assume that we have a finite-parametric family F
525 of limit functions.

526 If we have two sequences of transformations:

- 527 • a sequence f_i whose compositions tend to some function $f \in F$ and

528 • a sequence g_i whose composition tends to some function $g \in F$,
 529 then in the case when we first apply all f_i -transformations and then all g_i -transformations,
 530 then the resulting limit function $g(f(x))$ should also belong to the family F . Thus, the
 531 desired family F of all possible limit functions should be closed under composition.

532 Most transformations in sensors are reversible. So, if we limit ourselves to such
 533 transformations, and instead of first applying f_1 , then f_2 , etc., we change the direction
 534 of signal processing and first apply f_n^{-1} , then f_{n-1}^{-1} , etc., then, in the limit, instead of the
 535 original limit function f we will get the inverse function $f^{-1}(x)$. So, the class F of all
 536 possible limit functions should contain, with each function f , its inverse function as well.
 537 So, the class F must be closed under composition and inverse. Such classes are known
 538 as *transformation groups*.

539 Also, linear transformations are ubiquitous. Thus, it make sense to consider finite-
 540 parametric groups that contain all linear transformations. What are these groups?

541 **Enter Norbert Wiener.** Interestingly, the answer to this question is related to Norbert
 542 Wiener, the father of cybernetics. As he describes in his pioneering monograph [40] on
 543 cybernetics, when he started working on engineering problems, at first, he trusted exact
 544 mathematical models much more than vague biological analogies. And then, when
 545 he came up with a draft design of a system for automatic vision, a neurophysiologist
 546 colleague Arturo Rosenblueth – who saw the corresponding picture – asked him with
 547 surprise since when Wiener has become interested in human vision: because it turned
 548 out that what Wiener came up with after many thoughts and tries was exactly the scheme
 549 implemented in human vision. This experience lead to Wiener's idea of *cybernetics*, a
 550 science studying both engineering and biological systems, in which one of the main
 551 ideas is that since we the humans are the product of billion years of improving evolution,
 552 our biology should be close to optimal – and thus simulating this biology can be very
 553 helpful in engineering.

554 In some cases, this optimality was indeed confirmed. In some other cases, Wiener
 555 became so confident in the related optimality that he made several mathematical hy-
 556 potheses based on this confidence. For example, he learned, from Dr. Rosenblueth, that
 557 when we get closer and closer to an object, there are several clearly distinct phases
 558 in our visual perception (which, by the way, again fits with the above explanation of
 559 discreteness):

- 560 • When the object is very far, all we see is a formless blurb – in other words, ob-
 561 jects obtained from one another by arbitrary smooth transformations cannot be
 562 distinguished.
- 563 • When the object gets closer, we can detect whether it is smooth or has sharp angles.
 564 We may see a circle as an ellipse, a square as a rhombus (diamond). At this stage,
 565 images obtained by a projective transformation are indistinguishable.
- 566 • When the object gets even closer, we can detect which lines are parallel but we may
 567 not yet detect the angles. For example, we are not sure whether what we see is a
 568 rectangle or a parallelogram. This stage corresponds to affine transformation.
- 569 • Then, we have a stage of similarity transformations – when we detect the shape but
 570 cannot yet detect its size.
- 571 • Finally, when the object is close enough, we can detect both its shape and its size.

572 Each stage can be thus described by an appropriate transformation group. So, Wiener
 573 conjectured that if there was a group intermediate between, e.g., all projective and all
 574 continuous transformations, our vision mechanism – the result of millions of years of
 575 improving evolution – would have used it. Thus, he formulated a hypothesis that such
 576 intermediate transformation groups are not possible [40].

577 Many mathematicians did not take this hypothesis too seriously – while they
 578 appreciated Wiener's engineering ideas, they thought that he was going too far in his
 579 analogies. But other mathematicians took it seriously – and, two decades after the
 580 first edition of Wiener's book, they came up with a formal proof that, indeed, under

581 reasonable conditions, there is only one transformation group that contains all linear (=
 582 affine) transformations and some non-linear ones: namely, the group of all projective
 583 transformations [41,42].

The general proof is very complicated – e.g., the paper [42] consists of more than 100 pages of dense mathematics. But good news is that at present, we are only interested in the transformations of 1D signals. In this case, projective transformations are nothing else but fractional-linear ones

$$f(x) = \frac{a \cdot x + b}{c \cdot x + d},$$

584 and the corresponding proof can be shortened to a few pages; see, e.g., [43,44].

585 So, we arrive at the following conclusion.

586 **So what are the limit transformations?** We have shown that limit transformations form
 587 a finite-parametric transformation group that contains all linear transformations, and
 588 that all transformations from such a group are fractional linear – with linear ones being
 589 a particular case.

590 Thus, we conclude that all limit transformations are fractional-linear.

591 **A similar conclusion can be made about all possible reasonable transformations.** In-
 592 stead of looking for *limit* transformations, we can consider a different problem: to
 593 describe a class of all transformations which are, in some sense, *reasonable*. Linear trans-
 594 formations are reasonable: shift corresponds to the changing the starting point and a
 595 multiplication by a number corresponds to changing a measuring unit. A good example
 596 of both transformations are transformation between Celsius and Fahrenheit temperature
 597 scales.

598 It is also natural to conclude that a composition of two reasonable transformations is
 599 reasonable, and that a transformation which is inverse to a reasonable transformation is
 600 also reasonable. If we want to use computers to deal with reasonable transformations, it
 601 also makes sense to require that the reasonable transformations form a finite-parametric
 602 family – since in a computer, we can only stored finitely many parameter values.

603 Thus, the class of all reasonable transformations forms a finite-parametric trans-
 604 formation group containing all linear transformations. So, we conclude that every
 605 reasonable transformation is fractional linear.

606 **What are the implications for neural networks.** Artificial neural networks – a perfect
 607 example of Wiener’s belief that emulating biological systems can be beneficial – are
 608 formed of *neurons*. In a neuron, first, we form a linear combination x of the inputs x_i ,
 609 and then we apply some non-linear transformation $y = s(x)$ to this linear combination.
 610 In neural networks, this nonlinear transformation is known as an *activation function*.

Which activation function should we use? The first nonlinear neurons use *sigmoid*
 activation function

$$s(x) = \frac{1}{1 + \exp(-x)},$$

611 because, in the first approximation, this is how signals are processed in biological
 612 neurons; see, e.g., [14]. This activation function worked very well – much better than
 613 other activation functions that have been tried. This activation function is still often used
 614 in some layers of deep neural networks [15], where they are also very successful. How
 615 can we explain this success?

A possible explanation comes from the fact that, as we have mentioned earlier, all
 inputs come with noise. The simplest case is when, for each measurement, we just have a
 constant noise $n_i = \text{const}$, when instead of the actual values x_i , the measurement results
 are shifted by this value n_i , to $x_i + n_i$. As a result, the linear combination x is also shifted
 by some constant n (which is the similar linear combination of noises n_i):

$$x \rightarrow x + n.$$

616 We do not know the exact value of this noise – if we knew, we could simply subtract
 617 it from all the measured values. It is therefore reasonable to require that the result of
 618 applying the activation functions should be insensitive to this noise as much as possible.

619 Of course, we cannot simply require that $s(x + n) = s(x)$ for all x and n – this
 620 would imply that the function $s(x)$ is a constant that does not depend on the input at
 621 all. This makes sense: for example, the formula $d = v \cdot t$ showing that the distance
 622 can be obtained by multiplying velocity and time does not change when we change
 623 the unit of time, e.g., from hours to seconds. However, this invariance does not mean
 624 that the formula remains exactly the same when we change the unit of time: to keep
 625 the formula the same, we also need to apply an appropriate transformation to velocity
 626 as well: namely, replace the values in km/h with a value in km/sec. Similarly here,
 627 a natural idea is to require that if we apply a shift $x \rightarrow x' = x + n$ to the input, the
 628 formula remains the same if we apply an appropriate transformation to y as well, i.e.,
 629 that $y' = s(x')$, where $y' = T(y)$ for some reasonable transformation T .

In other words, we conclude that for every value n , there exists some reasonable
 transformation T_n for which $s(x') = T_n(y)$. Here, $x' = x + n$, and $y = s(x)$, so $s(x + n) =$
 $T_n(s(x))$. We have already concluded that reasonable transformations are fractional
 linear, thus we have

$$s(x + n) = \frac{a(n) \cdot x + b(n)}{c(n) \cdot x + d(n)}$$

630 for some values $a(n)$ through $d(n)$. To describe all the functions $s(x)$ that have this
 631 property, we can differentiate both side of this equation by n and take $n = 0$. The
 632 resulting differential equation can then be explicitly solved; see, e.g., [43,45,46]. The
 633 generic monotonic solution to this equation indeed differs from the sigmoid activation
 634 functions only by linear transformations of x and y .

635 This explains why the sigmoid activation function indeed works well in *many*
 636 application problems.

637 *Comment.* Of course, this does not mean that this activation function works best in *all*
 638 practical applications. For example, in most layers of deep neural networks, a different
 639 activation function $s(x) = \max(0, x)$ – known as *rectified linear* activation function –
 640 works much better. Interestingly, similar invariance ideas can explain the use of the
 641 rectified linear activation function – as well many other empirically successful features
 642 of deep learning algorithms; see, e.g., [46].

643 5. Conclusions

644 In this paper, we showed that limit theorems – similar to the Central Limit Theorem
 645 from statistics – make analysis of complex systems easier – i.e., lead to the blessing-of-
 646 dimensionality phenomenon. We showed that this simplification happens for all the
 647 aspects of these systems:

- 648 • for the corresponding transformations – as shown, e.g., by the description of all
 649 possible limit and/or reasonable transformations, and by the resulting theoretical
 650 explanation of the efficiency of sigmoid activation functions;
- 651 • for the system's uncertainty – as shown, e.g., by the use of limit distributions such
 652 as normal and Cauchy to make uncertainty quantification more efficient, and by the
 653 use of limit theorems to explain the ubiquity of interval uncertainty, and
- 654 • the desired result of the system's analysis – as shown, e.g., by a limit-theorem-based
 655 explanation of why it is usually possible to meaningfully classify objects into a small
 656 finite number of classes.

657 **Author Contributions:** Both authors contributed equally to this paper. Both authors have read
 658 and agreed to the published version of the manuscript.

659 **Funding:** This work was supported in part by the National Science Foundation grants 1623190 (A
 660 Model of Change for Preparing a New Generation for Professional Practice in Computer Science),
 661 and HRD-1834620 and HRD-2034030 (CAHSI Includes). It was also supported by the program of

- 662 the development of the Scientific-Educational Mathematical Center of Volga Federal District No.
663 075-02-2020-1478.
- 664 **Acknowledgments:** The authors are greatly thankful to Alexander Gorban for his encouragement
665 and valuable discussions, and to the anonymous referees for important suggestions.
- 666 **Conflicts of Interest:** The authors declare no conflict of interest.

References

1. Kainen, P.C. Utilizing geometric anomalies of high dimension: when complexity makes computations easier, In: Warwick, K.; Kárný, M.M. (eds.), *Computer-Intensive Methods in Control and Signal Processing*, Springer: New York, 1997, pp. 283–294.
2. Sheskin, D.J. *Handbook of Parametric and Non-Parametric Statistical Procedures*, Chapman & Hall/CRC: London, UK, 2011.
3. Donoho, D.L. High-dimensional data analysis: the curses and blessings of dimensionality, *Proceedings of the American Mathematical Society Conference on Math Challenges of the 21st Century*, Los Angeles, California, August 6–12, 2020.
4. Gorban, A.N.; Tyukin, I.Y.; Romanenko, I. The blessing of dimensionality: separation theorems in the thermodynamic limit, *IFAC-PapersOnLine* **2016**, *49*(24), 64–69.
5. Gorban, A.N.; Golubkov, V.; Grechuk, B.; Mirkes, E.M.; Tyukin, I.Y. Correction of AI systems by linear discriminants: probabilistic foundations, *Information Sciences* **2018**, *466*, 303–322.
6. Gorban, A.N.; Tyukin, I.Y. Blessing of dimensionality: mathematical foundations of the statistical physics of data, *Philosophical Transactions of the Royal Society, Series A* **2018**, *376*, Article 20170237.
7. Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*, Cambridge University Press: Cambridge, UK, 2018.
8. Gorban, A.N.; Makarov, V.A.; Tyukin, I.Y. The unreasonable effectiveness of small neural ensembles in high-dimensional brain, *Physics of Life Reviews* **2019**, *29*, 55–88.
9. Kreinovich, V. The heresy of unheard-of simplicity: Comment on [8], *Physics of Life Review* **2019**, *29*, 93–95.
10. Grechuk, B.; Gorban, A.N.; Tyukin, I.Y. General stochastic separation theorems with optimal bounds, *Neural Networks* **2021**, *138*, 33–56.
11. Tyukin, I.Y.; Higham, D.J.; Gorban, A.N. On adversarial examples and stealth attacks in artificial intelligence systems, *Proceedings of the International Joint Conference on Neural Networks IJCNN'2020*, Glasgow, UK, July 19–24, 2020, pp. 1–6.
12. Alexander, D.M.; Jurica, P.; Trengove, C.; Nikolaev, A.R.; Gepshtein, S.; Zvyagintsev, M.; Mathiak, K.; Schulze-Bonhage, A.; Ruescher, J.; Ball, T.; van Leeuwen, C. Traveling waves and trial averaging: The nature of single-trial and averaged brain responses in large-scale cortical signals, *Neuroimage* **2013**, *73*, 95–112.
13. Alexander, D.M.; Trengove, C.; van Leeuwen, C. Donders is dead: cortical traveling waves and the limits of mental chronometry in cognitive neuroscience, *Cognitive Processing* **2015**, *16*, 365–375.
14. Bishop, C.M. *Pattern Recognition and Machine Learning*, Springer: New York, 2006.
15. Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*, MIT Press: Cambridge, Massachusetts, 2016.
16. Tsirel'son, B.S. A geometrical approach to maximum likelihood estimation for infinite-dimensional Gaussian location. I, *Theory of Probability and its Applications* **1982**, *27*, 411–418.
17. Nguyen, H.T.; Wu, B.; Kreinovich, V. Our reasoning is clearly fuzzy, so why is crisp logic so often adequate?, *International Journal of Intelligent Technologies and Applied Statistics (IJITAS)* **2015**, *8*(2), 133–137.
18. Einstein, A. *Collected Papers of Albert Einstein*, Princeton University Press: Princeton, New Jersey, 2009.
19. Schlipp, P.A. *Albert Einstein: Philosopher-Scientist*, MJF Books: New York, 2001.
20. Kumar, M. *Quantum: Einstein, Bohr, and the Great Debate about the Nature of Reality*, W. W. Norton & Company: New York, 2011.
21. Rabinovich, S.G. *Measurement Errors and Uncertainties: Theory and Practice*, Springer: New York, 2005.
22. Feynman, R.; Leighton, R.; Sands, M. *The Feynman Lectures on Physics*, Addison Wesley: Boston, Massachusetts, 2005.
23. Moore, R.E.; Kearfott, R.B.; Cloud, M.J. *Introduction to Interval Analysis*, SIAM: Philadelphia, 2009.
24. Mayer, G. *Interval Analysis and Automatic Result Verification*, de Gruyter: Berlin, 2017.
25. Kreinovich, V.; Lakeyev, A.; Rohn, J.; Kahl, P. *Computational Complexity and Feasibility of Data Processing and Interval Computations*, Kluwer: Dordrecht, 1998.
26. Kreinovich, V.; Ferson, S. A new Cauchy-based black-box technique for uncertainty in risk analysis, *Reliability Engineering and Systems Safety* **2004**, *85*(1–3), 267–279.
27. Kreinovich, V. Why intervals? A simple limit theorem that is similar to limit theorems from statistics", *Reliable Computing* **1995**, *1*(1), 33–40.
28. Roginskaya, M.M.; Shulman, V.S. On Minkowski sums of many small sets, *Functional Analysis and Its Applications* **2018**, *52*(3), 233–235.
29. Nocedal, G.; Wright, S.J. *Numerical Optimization*, Springer: New York, 2006.
30. Urenda, J.C.; Kosheleva, O.; Kreinovich, V. How to describe measurement errors: a natural generalization of the Central Limit Theorem beyond normal (and other infinitely divisible) distributions, In: Pavese, F.; Forbes, A.B.; Zhang, N.F.; Chunovkina, A.G. (eds.), *Advanced Mathematical and Computational Tools in Metrology and Testing XII*, World Scientific: Singapore, to appear.
31. Embrechts, P.; Klüppelberg, C.; Mikosch, T. *Modelling Extremal Events for Insurance and Finance*, Springer Verlag: Berlin, 1997.

32. Kotz, S.; Nadarajah, S. *Extreme Value Distributions: Theory and Applications*, Imperial College Press: London, UK, 2000.
33. Coles, S. *An Introduction to Statistical Modeling of Extreme Values*, Springer Verlag: London, 2001.
34. Beirlant, J.; Goegebeur, Y.; Segers, J.; Teugels, J. *Statistics of Extremes: Theory and Applications*, Wiley: New York, 2004.
35. de Haan, L.; Ferreira, A. *Extreme Value Theory: An Introduction*, Springer Verlag: New York, 2006.
36. Resnick, S.I. *Extreme Values, Regular Variation and Point Processes*, Springer Verlag: New York, 2008.
37. Novak, S.Y. *Extreme Value Methods with Applications to Finance*, Chapman & Hall/CRC Press: London, 2011.
38. Gumbel, E.J. *Statistics of Extremes*, Dover: New York, 2013.
39. Kreinovich, V.; Nguyen, H.T.; Sriboonchitta, S.; Kosheleva, O. Modeling extremal events is not easy: why the extreme value theorem cannot be as general as the central limit theorem. In: Kreinovich, V. (ed.), *Uncertainty Modeling*, Springer Verlag: Cham, Switzerland, 2017, 123–134.
40. Wiener, N. *Cybernetics, or Control and Communication in the Animal and the Machine*, 3rd edition, MIT Press: Cambridge, Massachusetts, 1962.
41. Guillemin, V.M.; Sternberg, S. An algebraic model of transitive differential geometry, *Bulletin of American Mathematical Society* **1964**, 70(1), 16–47.
42. Singer, I.M.; Sternberg, S. Infinite groups of Lie and Cartan, Part 1, *Journal d'Analyse Mathématique* **1965**, XV, 1–113.
43. Nguyen, H.T.; Kreinovich, V. *Applications of Continuous Mathematics to Computer Science*, Kluwer: Dordrecht, 1997.
44. Zapata, F.; Kosheleva, O.; Kreinovich, V. Wiener's conjecture about transformation groups helps predict which fuzzy techniques work better, *Proceedings of the 2014 Annual Conference of the North American Fuzzy Information Processing Society NAFIPS'2014*, Boston, Massachusetts, June 24–26, 2014.
45. Kreinovich, V.; Quintana, C. Neural networks: what non-linearity to choose?, *Proceedings of the 4th University of New Brunswick Artificial Intelligence Workshop*, Fredericton, N.B., Canada, 1991, 627–637.
46. Kreinovich, V.; Kosheleva, O. Optimization under uncertainty explains empirical success of deep learning heuristics", In: Pardalos, P.; Rasskazova, V.; Vrahatis, M.N. (eds.), *Black Box Optimization, Machine Learning and No-Free Lunch Theorems*, Springer: Cham, Switzerland, 2021, pp. 195–220.