

Review

High-Dimensional Brain in a High-Dimensional World: Blessing of Dimensionality

Alexander N. Gorban^{1,2,*} , Valery A. Makarov^{2,3}  and Ivan Y. Tyukin^{1,2} 

¹ Department of Mathematics, University of Leicester, Leicester LE1 7RH, UK; I.Tyukin@le.ac.uk

² Laboratory of Advanced Methods for High-Dimensional Data Analysis, Lobachevsky University, 603022 Nizhny Novgorod, Russia; vmakarov@ucm.es

³ Instituto de Matemática Interdisciplinar, Faculty of Mathematics, Universidad Complutense de Madrid, Avda Complutense s/n, 28040 Madrid, Spain

* Correspondence: a.n.gorban@le.ac.uk

Received: 9 December 2019; Accepted: 6 January 2020; Published: 9 January 2020

Abstract: High-dimensional data and high-dimensional representations of reality are inherent features of modern Artificial Intelligence systems and applications of machine learning. The well-known phenomenon of the “curse of dimensionality” states: many problems become exponentially difficult in high dimensions. Recently, the other side of the coin, the “blessing of dimensionality”, has attracted much attention. It turns out that generic high-dimensional datasets exhibit fairly simple geometric properties. Thus, there is a fundamental tradeoff between complexity and simplicity in high dimensional spaces. Here we present a brief explanatory review of recent ideas, results and hypotheses about the blessing of dimensionality and related simplifying effects relevant to machine learning and neuroscience.

Keywords: artificial intelligence; mistake correction; concentration of measure; discriminant; data mining; geometry

1. Introduction

During the last two decades, the curse of dimensionality in data analysis was complemented by the blessing of dimensionality: if a dataset is essentially high-dimensional then, surprisingly, some problems get easier and can be solved by simple and robust old methods. The curse and the blessing of dimensionality are closely related, like two sides of the same coin. The research landscape of these phenomena is gradually becoming more complex and rich. New theoretical achievements and applications provide a new context for old results. The single-cell revolution in neuroscience, phenomena of grandmother cells and sparse coding discovered in the human brain meet the new mathematical ‘blessing of dimensionality’ ideas. In this mini-review, we aim to provide a short guide to new results on the blessing of dimensionality and to highlight the path from the curse of dimensionality to the blessing of dimensionality. The selection of material and angle of view is based on our own experience. We are not trying to cover everything in the subject of review, but rather fill in the gaps in existing tutorials and surveys.

R. Bellman [1] in the preface to his book, discussed the computational difficulties of multidimensional optimization and summarized them under the heading “curse of dimensionality”. He proposed to re-examine the situation, not as a mathematician, but as a “practical man” [2], and concluded that the price of excessive dimensionality “arises from a demand for too much information”. Dynamic programming was considered a method of dimensionality reduction in the optimization of a multi-stage decision process. Bellman returned to the problem of dimensionality reduction many times in different contexts [3]. Now, dimensionality reduction is an essential element of the engineering (the “practical man”) approach to mathematical modeling [4]. Many model reduction

methods were developed and successfully implemented in applications, from various versions of principal component analysis to approximation by manifolds, graphs, and complexes [5–7], and low-rank tensor network decompositions [8,9].

Various reasons and forms of the curse of dimensionality were classified and studied, from the obvious combinatorial explosion (for example, for n binary Boolean attributes, to check all the combinations of values we have to analyze 2^n cases) to more sophisticated distance concentration: in a high-dimensional space, the distances between randomly selected points tend to concentrate near their mean value, and the neighbor-based methods of data analysis become useless in their standard forms [10,11]. Many “good” polynomial time algorithms become useless in high dimensions.

Surprisingly, however, and despite the expected challenges and difficulties, common-sense heuristics based on the simple and the most straightforward methods “can yield results which are almost surely optimal” for high-dimensional problems [12]. Following this observation, the term “blessing of dimensionality” was introduced [12,13]. It was clearly articulated as a basis of future data mining in the Donoho “Millennium manifesto” [14]. After that, the effects of the blessing of dimensionality were discovered in many applications, for example in face recognition [15], in analysis and separation of mixed data that lie on a union of multiple subspaces from their corrupted observations [16], in multidimensional cluster analysis [17], in learning large Gaussian mixtures [18], in correction of errors of multidimensional machine learning systems [19], in evaluation of statistical parameters [20], and in the development of generalized principal component analysis that provides low-rank estimates of the natural parameters by projecting the saturated model parameters [21].

Ideas of the blessing of dimensionality became popular in signal processing, for example in compressed sensing [22,23] or in recovering a vector of signals from corrupted measurements [24], and even in such specific problems as analysis and classification of EEG patterns for attention deficit hyperactivity disorder diagnosis [25].

There exist exponentially large sets of pairwise almost orthogonal vectors (‘quasiorthogonal’ bases, [26]) in Euclidean space. It was noticed in the analysis of several n -dimensional random vectors drawn from the standard Gaussian distribution with zero mean and identity covariance matrix, that all the rays from the origin to the data points have approximately equal length, are nearly orthogonal and the distances between data points are all about $\sqrt{2}$ times larger [27]. This observation holds even for exponentially large samples (of the size $\exp(an)$ for some $a > 0$, which depends on the degree of the approximate orthogonality) [28]. Projection of a finite data set on random bases can reduce dimension with preservation of the ratios of distances (the Johnson–Lindenstrauss lemma [29]).

Such an intensive flux of works ensures us that we should not fear or avoid large dimensionality. We just have to use it properly. Each application requires a specific balance between the extraction of important low-dimensional structures (‘reduction’) and the use of the remarkable properties of high-dimensional geometry that underlie statistical physics and other fundamental results [30,31].

Both the curse and the blessing of dimensionality are the consequences of the measure concentration phenomena [30–33]. These phenomena were discovered in the development of the statistical backgrounds of thermodynamics. Maxwell, Boltzmann, Gibbs, and Einstein found that for many particles the distribution functions have surprising properties. For example, the Gibbs theorem of ensemble equivalence [34] states that a physically natural microcanonical ensemble (with fixed energy) is statistically equivalent (provides the same averages of physical quantities in the thermodynamic limit) to a maximum entropy canonical ensemble (the Boltzmann distribution). Simple geometric examples of similar equivalence gives the ‘thin shell’ concentration for balls: the volume of a high-dimensional ball is concentrated near its surface. Moreover, a high-dimensional sphere is concentrated near any equator (waist concentration; the general theory of such phenomena was elaborated by M. Gromov [35]). P. Lévy [36] analysed these effects and proved the first general concentration theorem. Modern measure concentration theory is a mature mathematical discipline with many deep results, comprehensive reviews [32], books [33,37], advanced textbooks [31], and even

elementary geometric introductions [38]. Nevertheless, surprising counterintuitive results continue to appear and push new achievements in machine learning, Artificial Intelligence (AI), and neuroscience.

This mini-review focuses on several novel results: stochastic separation theorems and evaluation of goodness of clustering in high dimensions, and on their applications to corrections of AI errors. Several possible applications to the dynamics of selective memory in the real brain and ‘simplicity revolution in neuroscience’ are also briefly discussed.

2. Stochastic Separation Theorems

2.1. Blessing of Dimensionality Surprises and Correction of AI Mistakes

D. Donoho and J. Tanner [23] formulated several ‘blessing of dimensionality’ surprises. In most cases, they considered M points sampled independently from a standard normal distribution in dimension n . Intuitively, we expect that some of the points will lie on the boundary of the convex hull of these points, and the others will be inside the interior of the hull. However, for large n and M , this expectation is wrong. This is the main surprise. With a high probability $p > 1 - \varepsilon$ all M random points are vertices of their convex hull. It is sufficient that $M < b \exp(an)$ for some a and b that depend on ε only [39,40]. Moreover, with a high probability, each segment connecting a pair of vertices is also an edge of the convex hull, and any simplex with k vertices from the sample is a $k - 1$ -dimensional face of the convex hull for some range of values of k . For uniform distributions in a ball, similar results were proved earlier by I. Bárány and Z. Füredi [41]. According to these results, each point of a random sample can be separated from all other points by a linear functional, even if the set is exponentially large.

Such a separability is important for the solution of a technological problem of fast, robust and non-damaging correction of AI mistakes [30,39,40]. AI systems make mistakes and will make mistakes in the future. If a mistake is detected, then it should be corrected. The complete re-training of the system requires too much resource and is rarely applicable to the correction of a single mistake. We proposed to use additional simple machine learning systems, correctors, for separation of the situations with higher risk of mistake from the situations with normal functioning [19,42] (Figure 1). The decision rules should be changed for situations with higher risk. Inputs for correctors are: the inputs of the original AI systems, the outputs of this system and (some) internal signals of this system [39,40]. The construction of correctors for AI systems is crucial in the development of the future AI ecosystems.

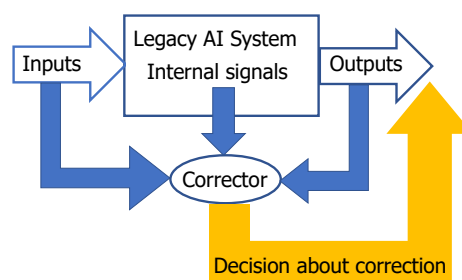


Figure 1. A scheme of corrector. Corrector receives some input, internal, and output signals from the legacy artificial intelligence (AI) system and classifies the situation as ‘high risk’ or ‘normal’ one. For a high-risk situation, it sends the corrected output to users following the correction rules. The high risk/normal situation classifier is prepared by supervised training on situations with diagnosed errors (universal construction). The online training algorithm could be very simple like Fisher’s linear discriminants or their ensembles [30,39,40,43,44]. Correction rules for high-risk situations are specific to a particular problem.

Correctors should [30]:

- be simple;

- not damage the existing skills of the AI system;
- allow fast non-iterative learning;
- correct new mistakes without destroying the previous fixes.

Of course, if an AI system made too many mistakes then their correctors could conflict. In such a case, re-training is needed with the inclusion of new samples.

2.2. Fisher Separability

Linear separation of data points from datasets [23,41] is a good candidate for the development of AI correctors. Nevertheless, from the ‘practical man’ point of view, one particular case, Fisher’s discriminant [45], is much more preferable to the general case because it allows one-shot and explicit creation of the separating functional.

Consider a finite data set Y without any hypothesis about the probability distribution. Let (\cdot, \cdot) be the standard inner product in \mathbb{R}^n . Let us define Fisher’s separability following [39].

Definition 1. A point x is Fisher-separable from a finite set Y with a threshold α ($0 \leq \alpha < 1$) if

$$(\mathbf{x}, \mathbf{y}) \leq \alpha(\mathbf{x}, \mathbf{x}), \quad \text{for all } \mathbf{y} \in Y \quad (1)$$

This definition coincides with the textbook definition of Fisher’s discriminant if the data set Y is whitened, which means that the mean point is in the origin and the sample covariance matrix is the identity matrix. Whitening is often a simple by-product of principal component analysis (PCA) because, on the basis of principal components, whitening is just the normalization of coordinates to unit variance. Again, following the ‘practical’ approach, we stress that the precise PCA and whitening are not necessary but rather a priori bounded condition number is needed: the ratio of the maximal and the minimal eigenvalues of the empirical covariance matrix after whitening should not exceed a given number $\kappa \geq 1$, independently of the dimension.

A finite set is called Fisher-separable, if each point is Fisher-separable from the rest of the set (Definition 3, [39]).

Definition 2. A finite set $Y \subset \mathbb{R}^n$ is called Fisher-separable with threshold $\alpha \in (0, 1)$ if inequality (1) holds for all $\mathbf{x}, \mathbf{y} \in Y$ such that $\mathbf{x} \neq \mathbf{y}$. The set Y is called Fisher-separable if there exists some α ($0 \leq \alpha < 1$) such that Y is Fisher-separable with threshold α .

Inequality (1) holds for vectors \mathbf{x}, \mathbf{y} if and only if \mathbf{x} does not belong to the ball (Figure 2):

$$\left\{ \mathbf{z} \mid \left\| \mathbf{z} - \frac{\mathbf{y}}{2\alpha} \right\| < \frac{\|\mathbf{y}\|}{2\alpha} \right\}. \quad (2)$$

2.3. Stochastic Separation for Distributions with Bounded Support

Let us analyse the separability of a random point from a finite set in the n -dimensional unit ball \mathbb{B}_n . Consider the distributions that can deviate from the equidistribution, and these deviations can grow with dimension n but not faster than the geometric progression with the common ratio $1/r > 1$, and, hence, the maximal density ρ_{\max} satisfies:

$$\rho_{\max} < \frac{C}{r^n V_n(\mathbb{B}_n)}, \quad (3)$$

where constant C does not depend on n .

For such a distribution in the unit ball, the probability ψ to find a random point \mathbf{x} in the excluded volume V_{excl} (Figure 2) tends to 0 as a geometric progression with the common ratio $b/(2r\alpha)$ when $n \rightarrow \infty$.

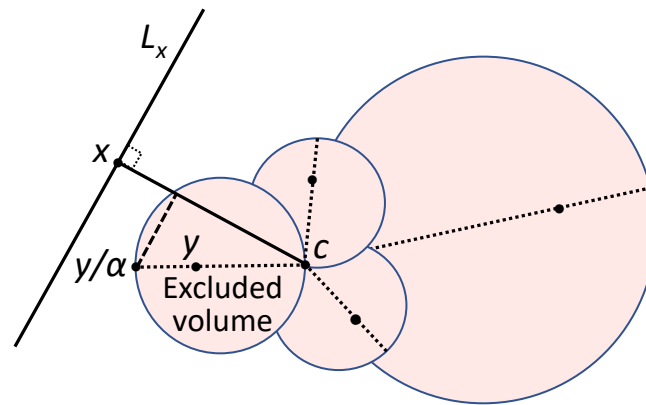


Figure 2. Fisher’s separability of a point x from a set Y . Diameters of the filled balls (excluded volume) are the segments $[c, y/\alpha]$ ($y \in Y$). Point x should not belong to the excluded volume to be separable from $y \in Y$ by the linear discriminant (1) with threshold α . Here, c is the origin (centre), and $L_x = \{z \mid (x, z) = (x, x)\}$ is the hyperplane. A point x should not belong to the union of such balls (filled) for all $y \in Y$ for separability from a set Y . The volume of this union, V_{excl} , does not exceed $V_n(\mathbb{B}_n)|Y|/(2\alpha)^n$.

Theorem 1. (Theorem 1, [39]) Let $Y \subset \mathbb{B}_n$, $|Y| < b^n$, and $2r\alpha > b > 1$. Assume that a probability distribution in the unit ball has a density with maximal value ρ_{max} , which satisfies inequality (3). Then the probability p that a random point from this distribution is Fisher-separable from Y is $p = 1 - \psi$, where the probability of inseparability

$$\psi < C \left(\frac{b}{2r\alpha} \right)^n.$$

Let us evaluate the probability that a random set Y is Fisher-separable. Assume that each point of Y is randomly selected from a distribution that satisfies (3). These distributions could be different for different $y \in Y$.

Theorem 2. (Theorem 2, [39]) Assume that a probability distribution in the unit ball has a density with maximal value ρ_{max} , which satisfies inequality (3). Let $|Y| < b^n$ and $2r\alpha > b^2 > 1$. Then the probability p that Y is Fisher-separable is $p = 1 - \psi$, where the probability of inseparability

$$\psi < |Y|C \left(\frac{b}{2r\alpha} \right)^n < C \left(\frac{b^2}{2r\alpha} \right)^n.$$

The difference in conditions from Theorem 1 is that here $2r\alpha > b^2$ and in Theorem 1 $2r\alpha > b$. Again, $|Y|$ can grow exponentially with the dimension as the geometric progression with the common factor $b > 0$, while $\psi \rightarrow 0$ faster than the geometric progression with the common factor $b^2/2r\alpha < 1$.

For illustration, if Y is an i.i.d. sample from the uniform distribution in the 100-dimensional ball and $|Y| = 2.7 \times 10^6$, then with probability $p > 0.99$ this set is Fisher-separable [42].

2.4. Generalisations

V. Kůrková [46] emphasized that many attractive measure concentration results are formulated for i.i.d. samples from very simple distributions (Gaussian, uniform, etc.), whereas the reality of big data is very different: the data are not i.i.d. samples from simple distributions. The machine learning theory based on the i.i.d. assumption should be revised, indeed [47]. In the theorems above two main restrictions were employed: the probability of a set occupying relatively small volume could not be large (3), and the support of the distribution is bounded. The requirement

of identical distribution of different points is not needed. The independence of the data points can be relaxed [39]. The boundedness of the support of distribution can be transformed to the ‘not-too-heavy-tail’ condition. The condition ‘sets of relatively small volume should not have large probability’ remains in most generalisations. It can be considered as ‘smeared absolute continuity’ because absolute continuity means that the sets of zero volume have zero probability. Theorems 1 and 2 have numerous generalisations [39,40,48,49]. Let us briefly list some of them:

- Product distributions in a unite cube where coordinates X_i are independent random variables with the variances separated from zero, $\text{var}(X_i) > \sigma_0^2 > 0$ (Theorem 2, [42]); significantly improved estimates are obtained by B. Grechuk [48].
- Log-concave distributions (a distribution with density $\rho(x)$ is log-concave if the set $D = \{x | \rho(x) > 0\}$ is convex and $g(x) = -\log \rho(x)$ is a convex function on D). In this case, the possibility of an exponential (non-Gaussian) tail brings a surprise: the upper size bound of the random set $|Y|$, sufficient for Fisher-separability in high dimensions with high probability, grows with dimension n as $\sim \exp(a\sqrt{n})$, i.e. slower than exponential (Theorem 5, [39]).
- Strongly log-concave distributions. A log concave distribution is strongly log-concave if there exists a constant $c > 0$ such that

$$\frac{g(x) + g(y)}{2} - g\left(\frac{x + y}{2}\right) \geq c\|x - y\|^2, \quad \forall x, y \in D.$$

In this case, we return to the exponential estimation of the maximal allowed size of $|Y|$ (Corollary 4, [39]). The comparison theorems [39] allow us to combine different distributions, for example the distribution from Theorem 2 in a ball with the log-concave or strongly log-concave tail outside the ball.

- The kernel versions of the stochastic separation theorem were found, proved and applied to some real-life problems [50].
- There are also various estimations beyond the standard i.i.d. hypothesis [39] but the general theory is yet to be developed.

2.5. Some Applications

The correction methods were tested on various AI applications for videostream processing: detection of faces for security applications and detection of pedestrians [39,44,51], translation of Sign Language into text for communication between deaf-mute people [52], knowledge transfer between AI systems [53], medical image analysis, scanning and classifying archaeological artifacts [54], etc., and even to some industrial systems with relatively high level of errors [43].

Application of the corrector technology to image processing was patented together with industrial partners [55]. A typical test of correctors’ performance is described below. For more detail of this test, we refer to [44]. A convolutional neural network (CNN) was trained to detect pedestrians in images. A set of 114,000 positive pedestrian and 375,000 negative non-pedestrian RGB images, re-sized to 128×128 , were collected and used as a training set. The testing set comprised of 10,000 positives and 10,000 negatives. The training and testing sets did not intersect.

We investigated in the computational experiments if it is possible to take one of cutting edge CNNs and train a one-neuron corrector to eliminate all the false positives produced. We also look at what effect, this corrector had on true positive numbers.

For each positive and false positive we extracted the second to last fully connected layer from CNN. These extracted feature vectors have dimension 4096. We applied PCA to reduce the dimension and analyzed how the effectiveness of the correctors depends on the number of principal components retained. This number varied in our experiments from 50 to 2000. The 25 false positives, taken from the testing set, were chosen at random to model single mistakes of the legacy classifier. Several such samples were chosen. For data projected on more than the first 87 principal components one neuron

with weights selected by the Fisher linear discriminant formula corrected 25 errors without doing any damage to classification capabilities (original skills) of the legacy AI system on the training set. For 50 or less principal components this separation is not perfect.

Single false positives were corrected successfully without any increase of the true positive rates. We removed more than 10 false positives at no cost to true positive detections in the street video data (Nottingham) by the use of a single linear function. Further increasing the number of corrected false positives demonstrated that a single-neuron corrector could result in gradual deterioration of the true positive rates.

3. Clustering in High Dimensions

Producing a special corrector for every single mistake seems to be a non-optimal approach, despite some successes. In practice, happily, often one corrector improves performance and prevents the system from some new mistakes because they are correlated. Moreover, mistakes can be grouped in clusters and we can create correctors for the clusters of situations rather than for single mistakes. Here we meet another measure concentration ‘blessing’. In high dimensions, clusters are good (well-separated) even in the situations when one can expect their strong intersection. For example, consider two clusters and the distance-based clustering. Let r_1^2 and r_2^2 be the mean squared Euclidean distances between the centroids of the clusters and their data points, and ρ be the distance between two centroids. The standard criteria of clusters’ quality [56] compare ρ with $r_1 + r_2$ and assume that for ‘good’ clusters $r_1 + r_2 < \rho$. Assume the opposite, $r_1 + r_2 > \rho$ and evaluate the volume of the intersection of two balls with radii r_1, r_2 . The intersection of the spheres (boundaries of the balls) is a $(n - 1)$ -dimensional sphere with the centre c (Figure 3). Assume $\rho^2 > |r_1^2 - r_2^2|$, which means that c is situated between the centers of the balls (otherwise, the biggest ball includes more than a half of the volume of the smallest one). The intersection of clusters belongs to a ball of radius R :

$$R^2 = \frac{r_1^2 + r_2^2}{2} - \frac{\rho^2}{4} - \frac{(r_1^2 - r_2^2)^2}{4\rho^2}. \quad (4)$$

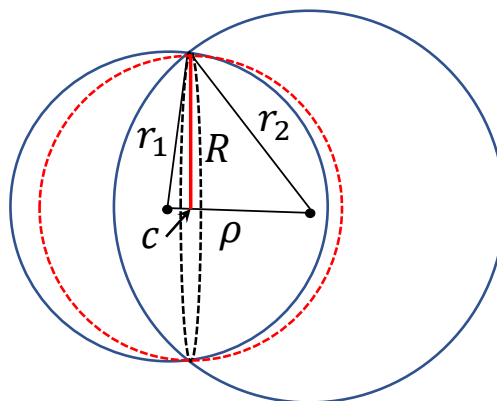


Figure 3. Measure of clustering quality. Intersection of two balls with the radii r_1, r_2 and the distance between centres $\rho < r_1 + r_2$ is included in a ball with radius R (4) and centre c (colored in red).

$R < r_{1,2}$ and the fractions of the volume of the two initial balls in the intersection is less than $(R/r_{1,2})^n$. These fractions evaluate the probability to confuse points between the clusters (for uniform distributions, for the Gaussian distributions the estimates are similar). We can measure the goodness of high-dimensional clusters by

$$\gamma = \left(\frac{R}{r_1}\right)^n + \left(\frac{R}{r_2}\right)^n.$$

Note that γ exponentially tends to zero with n increase. Small γ means ‘good’ clustering.

If $\gamma \ll 1$ then the probability to find a data point in the intersection of the balls (the ‘area of confusion’ between clusters) is negligible for uniform distributions in balls, isotropic Gaussian distributions and always when small volume implies small probability. Therefore, the clustering of mistakes for correction of high-dimensional machine learning systems gives good results even if clusters are not very good in the standard measures, and correction of clustered mistakes requires much fewer correctors for the same or even better accuracy [43].

We implemented the correctors with separation of clustered false-positive mistakes from the set of true positive and tested them on the classical face detection task [43]. The legacy object detector was an OpenCV implementation of the Haar face detector. It has been applied to video footage capturing traffic and pedestrians on the streets of Montreal. The powerful MTCNN face detector was used to generate ground truth data. The total number of true positives was 21,896, and the total number of false positives was 9372. The training set contained randomly chosen 50% of positives and false positives. PCA was used for dimensionality reduction with 200 principal components retained. Single-cluster corrector allows one to filter 90% of all errors at the cost of missing 5% percent of true positives. In dimension 200, a cluster of errors is sufficiently well-separated from the true positives. A significant classification performance gain was observed with more clusters, up to 100.

Further increase of dimension (the number of principal components retained) can even damage the performance because the number of features does not coincide with the dimension of the dataset, and the whitening with retained minor component can lead to ill-posed problems and loss of stability. For more detail, we refer to [43].

4. What Does ‘High Dimensionality’ Mean?

The dimensionality of data should not be naively confused with the number of features. Let us have n objects with p features. The usual data matrix in statistics is a 2D $n \times p$ array with n rows and p columns. The rows give values of features for an individual sample, and the columns give values of a feature for different objects. In classical statistics, we assume that $n \gg p$ and even study asymptotic estimates for $n \rightarrow \infty$ and p fixed. But, the modern ‘post-classical’ world is different [14]: the situation with $n < p$ (and even $n \ll p$) is not anomalous anymore. Moreover, it can be considered in some sense as the generic case: we can measure a very large number of attributes for a relatively small number of individual cases.

In such a situation the default preprocessing method could be recommended [57]: transform the $n \times p$ data matrix with $n < p$ into the square $n \times n$ matrix of inner products (or correlation coefficients) between the individual data vectors. After that, apply PCA and all the standard machinery of machine learning. New data will be presented by projections on the old samples. (Detailed description of this preprocessing and the following steps is presented in [57] with an applied case study for $n = 64$ and $p \approx 5 \times 10^5$.) Such a preprocessing reduces the apparent dimension of the data_space to $p \leq n$.

PCA gives us a tool for estimating the linear dimension of the dataset. Dimensionality reduction is achieved by using only the first few principal components. Several heuristics are used for evaluation of how many principal components should be retained:

- The classical Kaiser rule recommends to retain the principal components corresponding to the eigenvalues of the correlation matrix $\lambda \geq 1$ (or $\lambda \geq \alpha$ where $\alpha < 1$ is a selected threshold; often $\alpha = 0.5$ is selected). This is, perhaps, the most popular choice.
- Control of the fraction of variance unexplained. This approach is also popular, but it can retain too many minor components that can be considered ‘noise’.
- Conditional number control [39] recommends to retain the principal components corresponding to $\lambda \geq \lambda_{\max}/\kappa$, where λ_{\max} is the maximal eigenvalue of the correlation matrix and κ is the upper border of the conditional number (the recommended values are $\kappa \leq 10$ [58]). This recommendation is very useful because it provides direct control of multicollinearity.

After dimensionality reduction, we can perform whitening of data and apply the stochastic separation theorems. This requires a hypothesis about the distribution of data: sets of a relatively small volume should not have a high probability, and there should be no ‘heavy tails’. Unfortunately, this assumption is not always true in the practice of big data analysis. (We are grateful to G. Hinton and V. Kůrková for this comment.)

The separability properties can be affected by various violations of i.i.d. structure of data, inhomogeneity of data, small clusters and fine-grained lumping, and other peculiarities [59]. Therefore, the notion of dimension should be revisited. We proposed to use the Fisher separability of data to estimate the dimension [39]. For regular probability distributions, this estimate will give a standard geometric dimension, whereas, for complex (and often more realistic) cases, it will provide a more useful dimension characteristic. This approach was tested [59] for many bioinformatic datasets.

For analysis of Fisher’s separability and related estimation of dimensionality for general distribution and empirical datasets, an auxiliary random variable is used [39,59]. This is the probability that a randomly chosen point x is not Fisher-separable with threshold α from a given data point y by the discriminant (1):

$$p = p_y(\alpha) = \int_{\|z - \frac{y}{2\alpha}\| \leq \frac{\|y\|}{2\alpha}} \rho(z) dz, \tag{5}$$

where $\rho(z) dz$ is the probability measure for x .

If y is selected at random (not compulsory with the same distribution as x) then $p_y(\alpha)$ is a random variable. For a finite dataset Y the probability $p_Y(\alpha)$ that the data point is not Fisher-separable with threshold α from Y can be evaluated by the sum of $p_y(\alpha)$ for $y \in Y$:

$$p_Y(\alpha) \leq \sum_{y \in Y} p_y(\alpha). \tag{6}$$

Comparison of the empirical distribution of $p_y(\alpha)$ to the distribution evaluated for the high-dimensional sphere $S^{n-1} \subset \mathbb{R}^n$ can be used as information about the ‘effective’ dimension of data. The probability $p_y(\alpha)$ is the same for all $y \in S^{n-1}$ and exponentially decreases for large n . We assume that y is sampled randomly from for the rotationally invariant distribution on the unit sphere $S^{n-1} \subset \mathbb{R}^n$. For large n the asymptotic formula holds [39,59]:

$$p_y(\alpha) \approx \frac{(1 - \alpha^2)^{(n-1)/2}}{\alpha \sqrt{2\pi(n-1)}}. \tag{7}$$

Here $f(n) \approx g(n)$ means that $f(n)/g(n) \rightarrow 1$ when $n \rightarrow \infty$ (the functions here are strictly positive). It was noticed that the asymptotically equivalent formula with the denominator $\alpha \sqrt{2\pi n}$ performs better in small dimensions [59].

The introduced measure of dimension performs competitively with other state-of-the-art measures for simple i.i.d. data situated on manifolds [39,59]. It was shown to perform better in the case of noisy samples and allows estimation of the intrinsic dimension in situations where the intrinsic manifold, regular distribution and i.i.d. assumptions are not valid [59].

After this revision of the definition of data dimension, we can answer the question from the title of this section: What does ‘high dimensionality’ mean? The answer is given by the stochastic separation estimates for the uniform distribution in the unit sphere $S^{n-1} \subset \mathbb{R}^n$. Let $y \in S^{n-1}$. We use notation A_m for the volume (surface) of S^m . The points of S^{n-1} , which are not Fisher-separable from y with a given threshold α , form a spherical cap with the base radius $r = \sqrt{1 - \alpha^2}$ (Figure 4). The area of this cap is estimated from above by the lateral surface of the cone with the same base, which is tangent to the sphere at the base points (see Figure 4). Therefore, the probability ψ_α that a point selected randomly from the rotationally invariant distribution on S^{n-1} is not Fisher-separable from y is estimated from above as

$$p_y(\alpha) < \frac{A_{n-2}}{A_{n-1}} \frac{(1 - \alpha^2)^{(n-1)/2}}{\alpha(n-1)}. \tag{8}$$

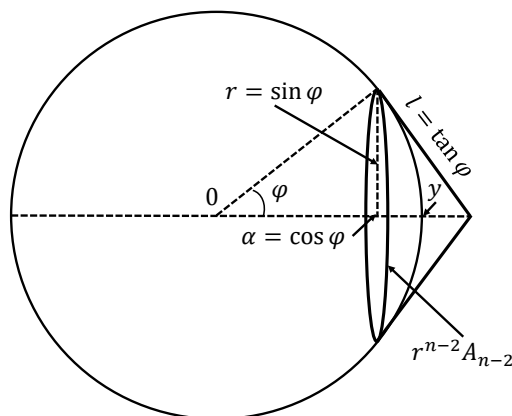


Figure 4. Estimation of the area of the spherical cap. A point of \mathbb{S}^{n-1} is Fisher-separable from $y \in \mathbb{S}^{n-1}$ with the threshold $\alpha = \cos \phi$ if and only if it does not belong to the spherical cap with the base radius $r = \sin \phi$ and the base plane orthogonal to y . The surface of this spherical cap is less than the lateral surface of the cone that is tangent to the base. The $n - 2$ -dimensional surface of the base is $A_B = r^{n-2} A_{n-2}$. The lateral surface of the cone is $l A_B / (n - 1)$.

The surface area of \mathbb{S}^{n-1} is

$$A_{n-1} = \frac{2\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2})}, \tag{9}$$

where Γ is Euler’s gamma-function.

Rewrite the estimate (8) as

$$p_y(\alpha) < \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} \frac{(1 - \alpha^2)^{(n-1)/2}}{\alpha\sqrt{\pi}}. \tag{10}$$

Recall that $\Gamma(x)$ is a monotonically increasing logarithmically convex function for $x \geq 3/2$ [60]. Therefore, for $n \geq 4$

$$\frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} < \frac{\Gamma(\frac{n+1}{2})}{\Gamma(\frac{n}{2})}.$$

Take into account that $\Gamma(\frac{n+1}{2}) = \frac{n-1}{2} \Gamma(\frac{n-1}{2})$ (because $\Gamma(x + 1) = x\Gamma(x)$). After elementary transforms it gives us

$$\frac{\Gamma^2(\frac{n}{2})}{\Gamma^2(\frac{n-1}{2})} < \frac{n-1}{2} \text{ and } \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} < \frac{\sqrt{n-1}}{\sqrt{2}}.$$

Finally, we got an elementary estimate for $p_y(\alpha)$ from above

$$p_y(\alpha) < \frac{(1 - \alpha^2)^{(n-1)/2}}{\alpha\sqrt{2\pi(n-1)}}. \tag{11}$$

Compared to (7), this estimate from above is asymptotically exact.

Estimate from above a probability of a separability violations using (11) and an elementary rule: for any family of events U_1, U_2, \dots, U_m ,

$$\mathbf{P}(U_1 \vee U_2 \vee \dots \vee U_m) \leq \mathbf{P}(U_1) + \mathbf{P}(U_2) + \dots + \mathbf{P}(U_m). \tag{12}$$

According to (11) and (12), if $0 < \psi < 1$, Y is an i.i.d. sample from a rotationally invariant distribution on \mathbb{S}^{n-1} and

$$|Y| \frac{(1 - \alpha^2)^{(n-1)/2}}{\alpha \sqrt{2\pi(n-1)}} < \psi, \tag{13}$$

then all sample points with a probability greater than $1 - \psi$ are Fischer-separable from a given point y with a threshold α . Similarly, if

$$|Y|^2 \frac{(1 - \alpha^2)^{(n-1)/2}}{\alpha \sqrt{2\pi(n-1)}} < \psi, \tag{14}$$

then with probability greater than $1 - \psi$ each sample point is Fisher-separable from the rest of the sample with a threshold α .

Estimates (13) and (14) provide sufficient conditions for separability. The Table 1 illustrates these estimates (the upper borders of $|Y|$ in these estimates are presented in the table with three significant figures). For an illustration of the separability properties, we estimated from above the sample size for which the Fisher-separability is guaranteed with a probability 0.99 and a threshold value $\alpha = 0.8$ (Table 1). These sample sizes grow fast with dimension. From the Fisher-separability point of view, dimensions 30 or 50 are already large. The effects of high-dimensional stochastic separability emerge with increasing dimensionality much earlier than, for example, the appearance of exponentially large quasi-orthogonal bases [28].

Table 1. Stochastic separation on $n - 1$ -dimensional spheres. For a sample size less than $M_{1,99}$, all points of an i.i.d. sample with a probability greater than 0.99 are Fischer-separable from a given point y with a threshold $\alpha = 0.8$. For a sample size less than $M_{2,99}$, with probability greater than 0.99 an i.i.d. sample is Fisher-separable with a threshold $\alpha = 0.8$ (that is, each sample point is Fisher-separable from the rest of the sample with this threshold).

	n = 10	n = 20	n = 30	n = 40	n = 50	n = 60	n = 70	n = 80
$M_{1,99}$	5	1.43×10^3	2.94×10^5	5.91×10^7	1.04×10^{10}	1.89×10^{12}	3.38×10^{14}	5.98×10^{15}
$M_{2,99}$	2	37	542	7.49×10^3	1.02×10^5	1.37×10^6	1.84×10^7	7.73×10^7

5. Discussion: The Heresy of Unheard-of Simplicity and Single Cell Revolution in Neuroscience

V. Kreinovich [61] summarised the impression from the effective AI correctors based on Fisher’s discriminant in high dimensions as “The heresy of unheard-of simplicity” using quotation of the famous Pasternak poetry. Such a simplicity appears also in brain functioning. Despite our expectation that complex intellectual phenomena is a result of a perfectly orchestrated collaboration between many different cells, there is a phenomenon of sparse coding, concept cells, or so-called ‘grandmother cells’ which selectively react to the specific concepts like a grandmother or a well-known actress (‘Jennifer Aniston cells’) [62]. These experimental results continue the single neuron revolution in sensory psychology [63].

The idea of grandmother or concept cells was proposed in the late 1960s. In 1972, Barlow published a manifest about the single neuron revolution in sensory psychology [63]. He suggested: “our perceptions are caused by the activity of a rather small number of neurons selected from a very large population of predominantly silent cells.” Barlow presented many experimental evidences of single-cell perception. In all these examples, neurons reacted selectively to the key patterns (called ‘trigger features’). This reaction was invariant to various changes in conditions.

The modern point of view on the single-cell revolution was briefly summarised recently by R. Quian Quiroga [64]. He mentioned that the ‘grandmother cells’ were invented by Lettvin “to ridicule the idea that single neurons can encode specific concepts”. Later discoveries changed the situation and added more meaning and detail to these ideas. The idea of concept cells was evolved during decades. According to Quian Quiroga, these cells are not involved in identifying a particular stimulus or concept. They are rather involved in creating and retrieving associations and can be

seen as the “building blocks of episodic memory”. Many recent discoveries used data received from intracranial electrodes implanted in the medial temporal lobe (MTL; the hippocampus and surrounding cortex) for patients medications. The activity of dozens of neurons can be recorded while patients perform different tasks. Neurons with high selectivity and invariance were found. In particular, one neuron fired to the presentation of seven different pictures of Jennifer Aniston and her spoken and written name, but not to 80 pictures of other persons. Emergence of associations between images was also discovered.

Some important memory functions are performed by stratified brain structures, such as the hippocampus. The CA1 region of the hippocampus includes a monolayer of morphologically similar pyramidal cells oriented parallel to the main axis (Figure 5). In humans, CA1 region of the hippocampus contains 1.4×10^7 of pyramidal neurons. Excitatory inputs to these neurons come from the CA3 regions (ipsi- and contra-lateral). Each CA3 pyramidal neuron sends an axon that bifurcates and leaves multiple collaterals in the CA1 (Figure 5b). This structural organization allows transmitting multidimensional information from the CA3 region to neurons in the CA1 region. Thus, we have simultaneous convergence and divergence of the information content (Figure 5b, right). A single pyramidal cell can receive around 30,000 excitatory and 1700 inhibitory inputs (data for rats [65]). Moreover, these numbers of synaptic contacts of cells vary greatly between neurons [66]. There are nonuniform and clustered connectivity patterns. Such a variability is considered as a part of the mechanism enhancing neuronal feature selectivity [66]. However, anatomical connectivity is not automatically transferred into functional connectivity and a realistic model should decrease significantly (by several orders of magnitude) the number of functional connections (see, for example, [67]). Nevertheless, even several dozens of effective functional connections are sufficient for the application of stochastic separation theorems (see Table 1).

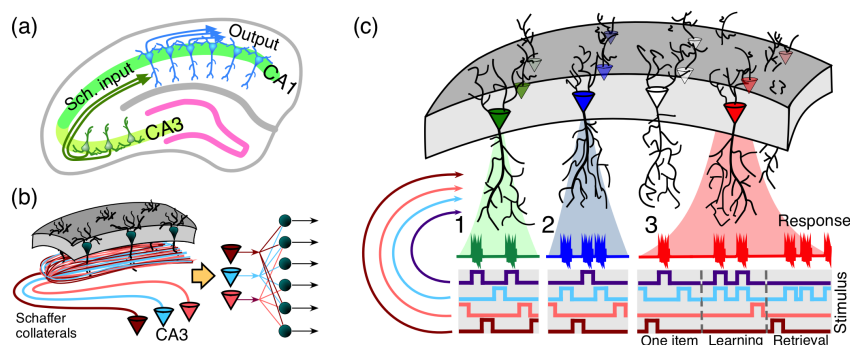


Figure 5. Organisation of encoding memories by single neurons in laminar structures: (a) laminar organization of the CA3 and CA1 areas in the hippocampus facilitates multiple parallel synaptic contacts between neurons in these areas by means of Schaffer collaterals; (b) axons from CA3 pyramidal neurons bifurcate and pass through the CA1 area in parallel (left) giving rise to the convergence–divergence of the information content (right). Multiple CA1 neurons receive multiple synaptic contacts from CA3 neurons; (c) schematic representation of three memory encoding schemes: (1) selectivity. A neuron (shown in green) receives inputs from multiple presynaptic cells that code different information items. It detects (responds to) only one stimulus (purple trace), whereas rejecting the others; (2) clustering. Similar to 1, but now a neuron (shown in blue) detects a group of stimuli (purple and blue traces) and ignores the others; (3) acquiring memories. A neuron (shown in red) learns dynamically a new memory item (blue trace) by associating it with a known one (purple trace). ((Figure 13, [40]), published under CC BY-NC-ND 4.0 license.).

For sufficiently high-dimensional sets of input signals a simple enough functional neuronal model with Hebbian learning (the generalized Oja rule [40,68]) is capable of explaining the following phenomena:

- the extreme selectivity of single neurons to the information content of high-dimensional data (Figure 5(c1)),
- simultaneous separation of several uncorrelated informational items from a large set of stimuli (Figure 5(c2)),
- dynamic learning of new items by associating them with already known ones (Figure 5(c3)).

These results constitute a basis for the organization of complex memories in ensembles of single neurons.

Re-training large ensembles of neurons is extremely time and resources consuming both in the brain and in machine learning. It is, in fact, impossible to realize such a re-training in many real-life situations and applications. “The existence of high discriminative units and a hierarchical organization for error correction are fundamental for effective information encoding, processing and execution, also relevant for fast learning and to optimize memory capacity” [69].

The multidimensional brain is the most puzzling example of the ‘heresy of unheard-of simplicity’, but the same phenomenon has been observed in social sciences and in many other disciplines [61].

There is a fundamental difference and complementarity between analysis of essentially high-dimensional datasets, where simple linear methods are applicable, and reducible datasets for which non-linear methods are needed, both for reduction and analysis [30]. This alternative in neuroscience was described as high-dimensional ‘brainland’ versus low-dimensional ‘flatland’ [70]. The specific multidimensional effects of the ‘blessing of dimensionality’ can be considered as the deepest reason for the discovery of small groups of neurons that control important physiological phenomena. On the other hand, even low dimensional data live often in a higher-dimensional space and the dynamics of low-dimensional models should be naturally embedded into the high-dimensional ‘brainland’. Thus, a “crucial problem nowadays is the ‘game’ of moving from ‘brainland’ to ‘flatland’ and backward” [70].

C. van Leeuwen formulated a radically opposite point of view [71]: neither high-dimensional linear models nor low-dimensional non-linear models have serious relations to the brain.

The devil is in the detail. First of all, the preprocessing is always needed to extract the relevant features. The linear method of choice is PCA. Various versions of non-linear PCA can be also useful [6]. After that, nobody has a guarantee that the dataset is either essentially high-dimensional or reducible. It can be a mixture of both alternatives, therefore both extraction of reducible lower-dimensional subset for nonlinear analysis and linear analysis of the high dimensional residuals could be needed together.

Author Contributions: Conceptualization, A.N.G., V.A.M. and I.Y.T.; Methodology, A.N.G., V.A.M. and I.Y.T.; Writing—Original Draft Preparation, A.N.G.; Writing—Editing, A.N.G., V.A.M. and I.Y.T.; Visualization, A.N.G., V.A.M. and I.Y.T. All authors have read and agreed to the published version of the manuscript.

Funding: The work was supported by the Ministry of Science and Higher Education of the Russian Federation (Project No. 14.Y26.31.0022). Work of A.N.G. and I.Y.T. was also supported by Innovate UK (Knowledge Transfer Partnership grants KTP009890; KTP010522) and University of Leicester. V.A.M acknowledges support from the Spanish Ministry of Economy, Industry, and Competitiveness (grant FIS2017-82900-P).

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Bellman, R. *Dynamic Programming*; Princeton University Press: Princeton, NJ, USA, 1957.
2. Bellman, R. The theory of dynamic programming. *Bull. Am. Math. Soc.* **1954**, *60*, 503–515. [[CrossRef](#)]
3. Bellman, R.; Kalaba, R. Reduction of dimensionality, dynamic programming, and control processes. *J. Basic Eng.* **1961**, *83*, 82–84. [[CrossRef](#)]
4. Gorban, A.N.; Kazantzis, N.; Kevrekidis, I.G.; Öttinger, H.C.; Theodoropoulos, C. (Eds.) *Model Reduction and Coarse-Graining Approaches for Multiscale Phenomena*; Springer: Berlin/Heidelberg, Germany, 2006.
5. Jolliffe, I. *Principal Component Analysis*; Springer: Berlin/Heidelberg, Germany, 1993.
6. Gorban, A.N.; Kégl, B.; Wunsch, D.; Zinovyev, A. (Eds.) *Principal Manifolds for Data Visualisation and Dimension Reduction*; Springer: Berlin/Heidelberg, Germany, 2008. [[CrossRef](#)]

7. Gorban, A.N.; Zinovyev, A. Principal manifolds and graphs in practice: from molecular biology to dynamical systems. *Int. J. Neural Syst.* **2010**, *20*, 219–232. [[CrossRef](#)] [[PubMed](#)]
8. Cichocki, A.; Lee, N.; Oseledets, I.; Phan, A.H.; Zhao, Q.; Mandic, D.P. Tensor networks for dimensionality reduction and large-scale optimization: Part 1 low-rank tensor decompositions. *Found. Trends[®] Mach. Learn.* **2016**, *9*, 249–429. [[CrossRef](#)]
9. Cichocki, A.; Phan, A.H.; Zhao, Q.; Lee, N.; Oseledets, I.; Sugiyama, M.; Mandic, D.P. Tensor networks for dimensionality reduction and large-scale optimization: Part 2 applications and future perspectives. *Found. Trends[®] Mach. Learn.* **2017**, *9*, 431–673. [[CrossRef](#)]
10. Beyer, K.; Goldstein, J.; Ramakrishnan, R.; Shaft, U. When is “nearest neighbor” meaningful? In Proceedings of the 7th International Conference on Database Theory (ICDT), Jerusalem, Israel, 10–12 January 1999; pp. 217–235. [[CrossRef](#)]
11. Pestov, V. Is the k-NN classifier in high dimensions affected by the curse of dimensionality? *Comput. Math. Appl.* **2013**, *65*, 1427–1437. [[CrossRef](#)]
12. Kainen, P.C. Utilizing geometric anomalies of high dimension: when complexity makes computation easier. In *Computer-Intensive Methods in Control and Signal Processing: The Curse of Dimensionality*; Warwick, K., Kárný, M., Eds.; Springer: New York, NY, USA, 1997; pp. 283–294. [[CrossRef](#)]
13. Brown, B.M.; Hall, P.; Young, G.A. On the effect of inliers on the spatial median. *J. Multivar. Anal.* **1997**, *63*, 88–104. [[CrossRef](#)]
14. Donoho, D.L. High-Dimensional Data Analysis: The Curses and Blessings of Dimensionality. Invited Lecture at Mathematical Challenges of the 21st Century. In Proceedings of the AMS National Meeting, Los Angeles, CA, USA, 6–12 August, 2000. Available online: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.329.3392> (accessed on 5 January 2020).
15. Chen, D.; Cao, X.; Wen, F.; Sun, J. Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification. In Proceedings of the 2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, 23–28 June 2013; pp. 3025–3032. [[CrossRef](#)]
16. Liu, G.; Liu, Q.; Li, P. Blessing of dimensionality: Recovering mixture data via dictionary pursuit. *IEEE Trans. Pattern Anal. Mach. Intell.* **2016**, *39*, 47–60. [[CrossRef](#)]
17. Murtagh, F. The remarkable simplicity of very high dimensional data: Application of model-based clustering. *J. Classif.* **2009**, *26*, 249–277. [[CrossRef](#)]
18. Anderson, J.; Belkin, M.; Goyal, N.; Rademacher, L.; Voss, J. The More, the Merrier: The Blessing of Dimensionality for Learning Large Gaussian Mixtures. In Proceedings of the 27th Conference on Learning Theory, Barcelona, Spain, 13–15 June 2014; Balcan, M.F., Feldman, V., Szepesvári, C., Eds.; PMLR: Barcelona, Spain, 2014; Volume 35, pp. 1135–1164.
19. Gorban, A.N.; Tyukin, I.Y.; Romanenko, I. The blessing of dimensionality: Separation theorems in the thermodynamic limit. *IFAC-PapersOnLine* **2016**, *49*, 64–69. [[CrossRef](#)]
20. Li, Q.; Cheng, G.; Fan, J.; Wang, Y. Embracing the blessing of dimensionality in factor models. *J. Am. Stat. Assoc.* **2018**, *113*, 380–389. [[CrossRef](#)] [[PubMed](#)]
21. Landgraf, A.J.; Lee, Y. Generalized principal component analysis: Projection of saturated model parameters. *Technometrics* **2019**. [[CrossRef](#)]
22. Donoho, D.L. Compressed sensing. *IEEE Trans. Inf. Theory* **2006**, *52*, 1289–1306. [[CrossRef](#)]
23. Donoho, D.; Tanner, J. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Phil. Trans. R. Soc. A* **2009**, *367*, 4273–4293. [[CrossRef](#)] [[PubMed](#)]
24. Candes, E.; Rudelson, M.; Tao, T.; Vershynin, R. Error correction via linear programming. In Proceedings of the 46th Annual IEEE Symposium on Foundations of Computer Science (FOCS’05), Pittsburgh, PA, USA, 23–25 October 2005; pp. 668–681. [[CrossRef](#)]
25. Pereda, E.; García-Torres, M.; Melián-Batista, B.; Mañas, S.; Méndez, L.; González, J.J. The blessing of Dimensionality: Feature Selection outperforms functional connectivity-based feature transformation to classify ADHD subjects from EEG patterns of phase synchronisation. *PLoS ONE* **2018**, *13*, e0201660. [[CrossRef](#)] [[PubMed](#)]
26. Kainen, P.; Kůrková, V. Quasiorthogonal dimension of Euclidian spaces. *Appl. Math. Lett.* **1993**, *6*, 7–10. [[CrossRef](#)]

27. Hall, P.; Marron, J.; Neeman, A. Geometric representation of high dimension, low sample size data. *J. R. Stat. Soc. B* **2005**, *67*, 427–444. [[CrossRef](#)]
28. Gorban, A.N.; Tyukin, I.; Prokhorov, D.; Sofeikov, K. Approximation with random bases: Pro et contra. *Inf. Sci.* **2016**, *364–365*, 129–145. [[CrossRef](#)]
29. Dasgupta, S.; Gupta, A. An elementary proof of a theorem of Johnson and Lindenstrauss. *Random Struct. Algorithms* **2003**, *22*, 60–65. [[CrossRef](#)]
30. Gorban, A.N.; Tyukin, I.Y. Blessing of dimensionality: Mathematical foundations of the statistical physics of data. *Philos. Trans. R. Soc. A* **2018**, *376*, 20170237. [[CrossRef](#)]
31. Vershynin, R. *High-Dimensional Probability: An Introduction with Applications in Data Science*; Cambridge Series in Statistical and Probabilistic Mathematics; Cambridge University Press: Cambridge, UK, 2018.
32. Giannopoulos, A.A.; Milman, V.D. Concentration property on probability spaces. *Adv. Math.* **2000**, *156*, 77–106. [[CrossRef](#)]
33. Ledoux, M. *The Concentration of Measure Phenomenon*; Number 89 in Mathematical Surveys & Monographs; AMS: Providence, RI, USA, 2005.
34. Gibbs, J.W. *Elementary Principles in Statistical Mechanics, Developed with Especial Reference to the Rational Foundation of Thermodynamics*; Dover Publications: New York, NY, USA, 1960.
35. Gromov, M. Isoperimetry of waists and concentration of maps. *Geom. Funct. Anal.* **2003**, *13*, 178–215. [[CrossRef](#)]
36. Lévy, P. *Problèmes Concrets D'analyse Fonctionnelle*; Gauthier-Villars: Paris, France, 1951.
37. Dubhashi, D.P.; Panconesi, A. *Concentration of Measure for the Analysis of Randomized Algorithms*; Cambridge University Press: Cambridge, UK, 2009.
38. Ball, K. An Elementary Introduction to Modern Convex Geometry. In *Flavors of Geometry*; Cambridge University Press: Cambridge, UK, 1997; Volume 31.
39. Gorban, A.N.; Golubkov, A.; Grechuk, B.; Mirkes, E.M.; Tyukin, I.Y. Correction of AI systems by linear discriminants: Probabilistic foundations. *Inf. Sci.* **2018**, *466*, 303–322. [[CrossRef](#)]
40. Gorban, A.N.; Makarov, V.A.; Tyukin, I.Y. The unreasonable effectiveness of small neural ensembles in high-dimensional brain. *Phys. Life Rev.* **2019**, *29*, 55–88. [[CrossRef](#)] [[PubMed](#)]
41. Bárány, I.; Füredi, Z. On the shape of the convex hull of random points. *Probab. Theory Relat. Fields* **1988**, *77*, 231–240. [[CrossRef](#)]
42. Gorban, A.N.; Tyukin, I.Y. Stochastic separation theorems. *Neural Netw.* **2017**, *94*, 255–259. [[CrossRef](#)]
43. Tyukin, I.Y.; Gorban, A.N.; McEwan, A.A.; Meshkinfamfard, S. Blessing of dimensionality at the edge. *arXiv* **2019**, arXiv:1910.00445.
44. Gorban, A.N.; Burton, R.; Romanenko, I.; Tyukin, I.Y. One-trial correction of legacy AI systems and stochastic separation theorems. *Inf. Sci.* **2019**, *484*, 237–254. [[CrossRef](#)]
45. Fisher, R.A. The Use of Multiple Measurements in Taxonomic Problems. *Ann. Eugenics* **1936**, *7*, 179–188. [[CrossRef](#)]
46. Kúrková, V. Some insights from high-dimensional spheres: Comment on “The unreasonable effectiveness of small neural ensembles in high-dimensional brain” by Alexander N. Gorban et al. *Phys. Life Rev.* **2019**, *29*, 98–100. [[CrossRef](#)]
47. Gorban, A.N.; Makarov, V.A.; Tyukin, I.Y. Symphony of high-dimensional brain. Reply to comments on “The unreasonable effectiveness of small neural ensembles in high-dimensional brain”. *Phys. Life Rev.* **2019**, *29*, 115–119. [[CrossRef](#)] [[PubMed](#)]
48. Grechuk, B. Practical stochastic separation theorems for product distributions. In Proceedings of the IEEE IJCNN 2019—International Joint Conference on Neural Networks, Budapest, Hungary, 14–19 July 2019. [[CrossRef](#)]
49. Kúrková, V.; Sanguineti, M. Probabilistic Bounds for Binary Classification of Large Data Sets. In Proceedings of the International Neural Networks Society, Genova, Italy, 16–18 April 2019; Oneto, L., Navarin, N., Sperduti, A., Anguita, D., Eds.; Springer: Berlin/Heidelberg, Germany, 2019; Volume 1, pp. 309–319. [[CrossRef](#)]
50. Tyukin, I.Y.; Gorban, A.N.; Grechuk, B. Kernel Stochastic Separation Theorems and Separability Characterizations of Kernel Classifiers. In Proceedings of the IEEE IJCNN 2019—International Joint Conference on Neural Networks, Budapest, Hungary, 14–19 July 2019. [[CrossRef](#)]

51. Meshkinfamfard, S.; Gorban, A.N.; Tyukin, I.V. Tackling Rare False-Positives in Face Recognition: A Case Study. In *Proceedings of the 2018 IEEE 20th International Conference on High Performance Computing and Communications; IEEE 16th International Conference on Smart City; IEEE 4th International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*; IEEE: Exeter, UK, 2018; pp. 1592–1598. [[CrossRef](#)]
52. Tyukin, I.Y.; Gorban, A.N.; Green, S.; Prokhorov, D. Fast construction of correcting ensembles for legacy artificial intelligence systems: Algorithms and a case study. *Inf. Sci.* **2019**, *485*, 230–247. [[CrossRef](#)]
53. Tyukin, I.Y.; Gorban, A.N.; Sofeikov, K.; Romanenko, I. Knowledge transfer between artificial intelligence systems. *Front. Neurobot.* **2018**, *12*. [[CrossRef](#)] [[PubMed](#)]
54. Allison, P.M.; Sofeikov, K.; Levesley, J.; Gorban, A.N.; Tyukin, I.; Cooper, N.J. Exploring automated pottery identification [Arch-I-Scan]. *Internet Archaeol.* **2018**, *50*. [[CrossRef](#)]
55. Romanenko, I.; Gorban, A.; Tyukin, I. Image Processing. U.S. Patent 10,489,634 B2, 26 November 2019. Available online: <https://patents.google.com/patent/US10489634B2/en> (accessed on 5 January 2020).
56. Xu, R.; Wunsch, D. *Clustering*; Wiley: Hoboken, NJ, USA, 2008.
57. Moczko, E.; Mirkes, E.M.; Cáceres, C.; Gorban, A.N.; Piletsky, S. Fluorescence-based assay as a new screening tool for toxic chemicals. *Sci. Rep.* **2016**, *6*, 33922. [[CrossRef](#)] [[PubMed](#)]
58. Dormann, C.F.; Elith, J.; Bacher, S.; Buchmann, C.; Carl, G.; Carré, G.; Marquéz, J.R.; Gruber, B.; Lafourcade, B.; Leitão, P.J.; et al. Collinearity: A review of methods to deal with it and a simulation study evaluating their performance. *Ecography* **2013**, *36*, 27–46. [[CrossRef](#)]
59. Albergante, L.; Bac, J.; Zinovyev, A. Estimating the effective dimension of large biological datasets using Fisher separability analysis. In *Proceedings of the IEEE IJCNN 2019—International Joint Conference on Neural Networks*, Budapest, Hungary, 14–19 July 2019. [[CrossRef](#)]
60. Artin, E. *The Gamma Function*; Courier Dover Publications: Mineola, NY, USA, 2015.
61. Kreinovich, V. The heresy of unheard-of simplicity: Comment on “The unreasonable effectiveness of small neural ensembles in high-dimensional brain” by A.N. Gorban, V.A. Makarov, and I.Y. Tyukin. *Phys. Life Rev.* **2019**, *29*, 93–95. [[CrossRef](#)]
62. Quian Quiroga, R.; Reddy, L.; Kreiman, G.; Koch, C.; Fried, I. Invariant visual representation by single neurons in the human brain. *Nature* **2005**, *435*, 1102–1107. [[CrossRef](#)]
63. Barlow, H.B. Single units and sensation: A neuron doctrine for perceptual psychology? *Perception* **1972**, *1*, 371–394. [[CrossRef](#)]
64. Quian Quiroga, R. Akakhievitch revisited: Comment on “The unreasonable effectiveness of small neural ensembles in high-dimensional brain” by Alexander N. Gorban et al. *Phys. Life Rev.* **2019**, *29*, 111–114. [[CrossRef](#)]
65. Megias, M.; Emri, Z.S.; Freund, T.F.; Gulyás, A.I. Total number and distribution of inhibitory and excitatory synapses on hippocampal CA1 pyramidal cells. *Neuroscience* **2001**, *102*, 527–540. [[CrossRef](#)]
66. Druckmann, S.; Feng, L.; Lee, B.; Yook, C.; Zhao, T.; Magee, J.C.; Kim, J. Structured synaptic connectivity between hippocampal regions. *Neuron* **2014**, *81*, 629–640. [[CrossRef](#)] [[PubMed](#)]
67. Brivanlou, I.H.; Dantzer, J.L.; Stevens, C.F.; Callaway, E.M. Topographic specificity of functional connections from hippocampal CA3 to CA1. *Proc. Natl. Acad. Sci. USA* **2004**, *101*, 2560–2565. [[CrossRef](#)] [[PubMed](#)]
68. Tyukin, I.; Gorban, A.N.; Calvo, C.; Makarova, J.; Makarov, V.A. High-dimensional brain: A tool for encoding and rapid learning of memories by single neurons. *Bull. Math. Biol.* **2019**, *81*, 4856–4888. [[CrossRef](#)] [[PubMed](#)]
69. Varona, P. High and low dimensionality in neuroscience and artificial intelligence: Comment on “The unreasonable effectiveness of small neural ensembles in high-dimensional brain” by A.N. Gorban et al. *Phys. Life Rev.* **2019**, *29*, 106–107. [[CrossRef](#)]
70. Barrio, R. “Brainland” vs. “flatland”: how many dimensions do we need in brain dynamics? Comment on the paper “The unreasonable effectiveness of small neural ensembles in high-dimensional brain” by Alexander N. Gorban et al. *Phys. Life Rev.* **2019**, *29*, 108–110. [[CrossRef](#)]
71. van Leeuwen, C. The reasonable ineffectiveness of biological brains in applying the principles of high-dimensional cybernetics: Comment on “The unreasonable effectiveness of small neural ensembles in high-dimensional brain” by Alexander N. Gorban et al. *Phys. Life Rev.* **2019**, *29*, 104–105. [[CrossRef](#)]

