

Augmented Artificial Intelligence: a Conceptual Framework

1st Alexander N Gorban
University of Leicester

and Lobachevsky University
Leicester, UK, and Nizhni Novgorod, Russia
a.n.gorban@le.ac.uk

2nd Bogdan Grechuk
University of Leicester

Leicester, UK
bg83@le.ac.uk

3rd Ivan Y Tyukin

University of Leicester
and Lobachevsky University
Leicester, UK, and Nizhni Novgorod, Russia
i.tyukin@le.ac.uk

Abstract—All artificial Intelligence (AI) systems make errors. These errors are unexpected, and differ often from the typical human mistakes (“non-human” errors). The AI errors should be corrected without damage of existing skills and, hopefully, avoiding direct human expertise. This paper presents an initial summary report of project taking new and systematic approach to improving the intellectual effectiveness of the individual AI by communities of AIs. We combine some ideas of learning in heterogeneous multiagent systems with new and original mathematical approaches for non-iterative corrections of errors of legacy AI systems. The mathematical foundations of AI non-destructive correction are presented and a series of new stochastic separation theorems is proven. These theorems provide a new instrument for the development, analysis, and assessment of machine learning methods and algorithms in high dimension. They demonstrate that in high dimensions and even for exponentially large samples, linear classifiers in their classical Fisher’s form are powerful enough to separate errors from correct responses with high probability and to provide efficient solution to the non-destructive corrector problem. In particular, we prove some hypotheses formulated in our paper ‘Stochastic Separation Theorems’ (Neural Networks, 94, 255–259, 2017), and answer one general problem published by Donoho and Tanner in 2009.

Index Terms—multiscale experts, knowledge transfer, non-iterative learning, error correction, measure concentration, blessing of dimensionality

I. INTRODUCTION

The history of neural networks research can be represented as a series of explosions or waves of inventions and expectations. This history ensures us that the popular Gartner’s hype cycle for emerging technologies presented by the solid curve on Fig. 1 (see, for example [1]) should be supplemented by the new peak of expectation explosion (dashed line). Some expectations from the previous peak are realized and move to the “Plateau of Productivity” but the majority of them jump to the next “Peak of Inflated Expectations”. This observation relates not only to neural technologies but perhaps to majority of IT innovations. It is surprising to see, how expectations reappear in the new wave from the previous peak often without modifications, just changing the human carriers.

ANG and IYT were Supported by Innovate UK (KTP009890 and KTP010522) and Ministry of science and education, Russia (14.Y26.31.0022). BG thanks the University of Leicester for granting him academic study leave to do this research.

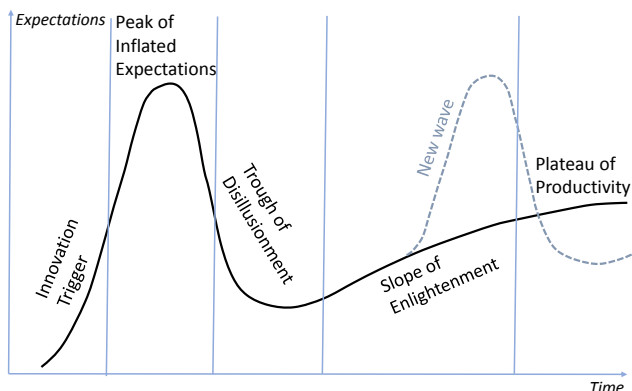


Figure 1. Gartner’s Hype Cycle for emerging technologies supplemented by a new peak.

Computers and networks have been expected to augment the human intelligence [2]. In 1998 one of the authors had been inspired by 8 years of success of knowledge discovery by deep learning neural network and by the transformation of their hidden knowledge into explicit knowledge in the form of “logically transparent networks” [3] by means of pruning, binarization and other simplification procedures [4], [5], and wrote: “I am sure that the neural network technology of knowledge discovery is a “point of growth”, which will remodel neuroinformatics, transform many areas of information technologies and create new approaches” [6]. Now it seems that this prediction will not be fulfilled: most customers do not care about gaining knowledge but prefer the “one button solutions”, which exclude humans from the process as far as it is possible. This is not a new situation in history. New intellectual technologies increase intellectual abilities of mankind, but not the knowledge of individual humans. Here, we can refer to Plato “There is an old Egyptian tale of Theuth, the inventor of writing, showing his invention to the god Thamus, who told him that he would only spoil men’s memories and take away their understandings” [7]. The adequate model of future Artificial Intelligence (AI) usage should include large communities of AI systems. Knowledge should circulate and grow in these communities. Participation

of humans in these processes should be minimized. In the course of this technical revolution not the “Augmented human intellect” but the continuously augmenting AI will be created.

In this work, we propose the conceptual framework for augmenting AI in communities or “social networks” of AIs. For construction of such social networks, we employ several ideas in addition to the classical machine learning. The first of them is separation of the problem areas between small local (neural) experts, their competitive and collaborative learning, and conflict resolution. In 1991, two first papers with this idea were published simultaneously [8], [9]. The techniques for distribution of tasks between small local experts were developed. In our version of this technology [8] and in all our applied software [3], [10]–[12] the neural network answers were always complemented by the evaluation of the network self-confidence. This self-confidence level is an important instrument for community learning.

The second idea is the blessing of dimensionality [13]–[17] and the AI correction method [21] based on stochastic separation theorems [20]. The “sparsity” of high-dimensional spaces and concentration of measure phenomena make some low-dimensional approaches impossible in high dimensions. This problem is widely known as the “curse of dimensionality”. Surprisingly, the same phenomena can be efficiently employed for creation of new, high-dimensional methods, which seem to be much simpler than in low dimensions. This is the “blessing of dimensionality”.

The classical theorems about concentration of measure state that random points in a highly-dimensional data distribution are concentrated in a thin layer near an average or median level set of a Lipschitz function (for introduction into this area we refer to [18]). The newly discovered stochastic separation theorems [20] revealed the fine structure of these thin layers: the random points are all linearly separable from the rest of the set even for exponentially large random sets. Of course, the probability distribution should be ‘genuinely’ high-dimensional for all these concentration and separation theorems.

Linear separability of exponentially large random subsets in high dimension allows us to solve the problem of nondestructive correction of legacy AI systems: the linear classifiers in their simplest Fisher’s form can separate errors from correct responses with high probability [21]. It is possible to avoid the standard assumption about independence and identical distribution of data points (i.i.d.). The non-iterative and non-destructive correctors can be employed for skills transfer in communities of AI systems [22].

These two ideas are joined in a special organisational environment of community learning which is organized in several phases:

- Initial supervising learning where community of newborn experts assimilate the knowledge hidden in labeled tasks from a problem-book (the problem-book is a continuously growing and transforming collection of samples);
- Non-iterative learning of community with self-labeling of real-life or additional training samples on the basis

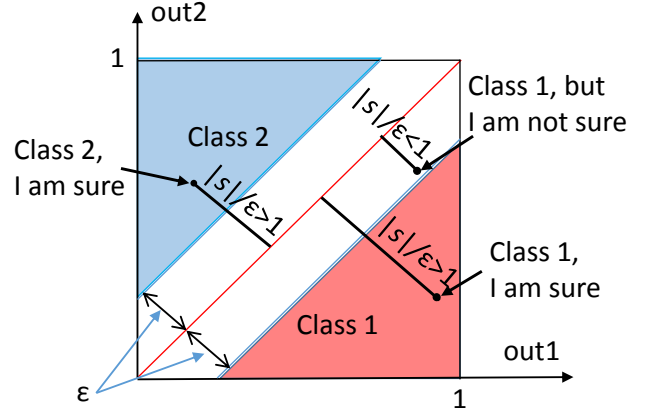


Figure 2. Answers and assurance; s is the deviation from the diagonal.

of separation of expertise between local experts, their continuous adaptation and mutual correction for the assimilation of gradual changes in reality.

- Interiorisation of the results of the self-supervising learning of community in the internal skills of experts.
- Development and learning of special network manager that evaluates the level of expertise of the local experts for a problem and distributes the incoming task flow between them.
- Using an “ultimate auditor” to assimilate qualitative changes in the environment and correct collective errors; it may be human inspection, a feedback from real life, or another system of interference into the self-labeling process.

We describe the main constructions of this approach using the example of classification problems and simple linear correctors. The correctors with higher abilities can be constructed on the basis of small neural networks with uncorrelated neurons [21] but already single-neuron correctors (Fisher’s discriminants) can help in explanation of a wealth of empirical evidence related to in-vivo recordings of “Grandmother” cells and “concept” cells [17], [23]. We pay special attention to the mathematical backgrounds of the technology and prove a series of new stochastic separation theorems. In particular, we find that the hypothesis about stochastic separability for general log-concave distribution [22] is true and describe a general class of probability measures with linear stochastic separability, the question was asked in 2009 by Donoho and Tanner [19].

II. SUPERVISING STAGE: PROBLEM OWNERS, MARGINS, SELF-CONFIDENCE, AND ERROR FUNCTIONS

Consider binary classification problems. The neural experts with arbitrary internal structure have two outputs, out1 and out2, with interpretation: the sample belongs to class 1 if $out1 \geq out2$ and it belongs to class 2 if $out1 < out2$. For any given $\varepsilon > 0$ we can define the level of (self-)confidence in the

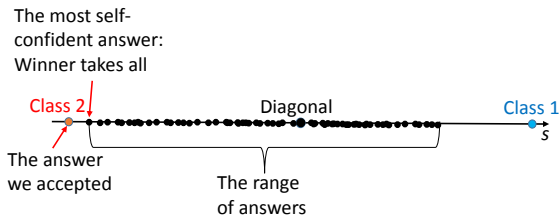


Figure 3. Interpretation of community answer: Most self-confident winner takes all. Dots correspond to the various agents' answers, s is defined in Fig. 2.

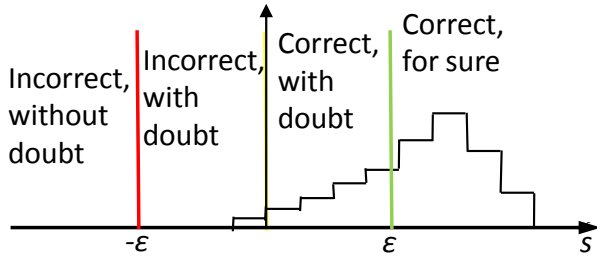


Figure 4. Histogram of answers for trained community of agents: they should give a correct self-confident answer to the samples they own, and do not make large mistakes on all other examples they never met before.

classification answer as it demonstrated in Fig. 2. The *owner of a sample* is an expert that gives the best (correct and most confident) answer for this sample. If we assume the single owner for every sample then in the community functioning for problem solving this single owner gives the final result (Fig. 3).

We aim to train the community of agents in such a way that they will give correct self-confident answers to the samples they own, and do not make large mistakes on all other examples they never met before. The desired histogram of answers is presented in Fig. 3.

Learning is minimisation of error functionals, which is defined for any selected sample and any local expert. This error function should be different for owners and non-owners of the sample. If we assume that each sample has a single owner then the error function presented in Fig. 5 can be used.

Voting of k most self-confident experts (Fig. 6) can make the decision more stable. This voting may be organised with weights of votes, which depend on the individual experts' level of confidence, or without weights, just as a simple voting. The modified error function for system with collective ownership (each sample has k owners) is needed (Fig. 7). This function is constructed to provide proper answers of all k owners.

III. SELF-LEARNING STAGE: COMMUNITIES AND RECOMMENDER SYSTEMS

After the stage of supervising learning, community of local experts can start working with new, previously unlabeled data.

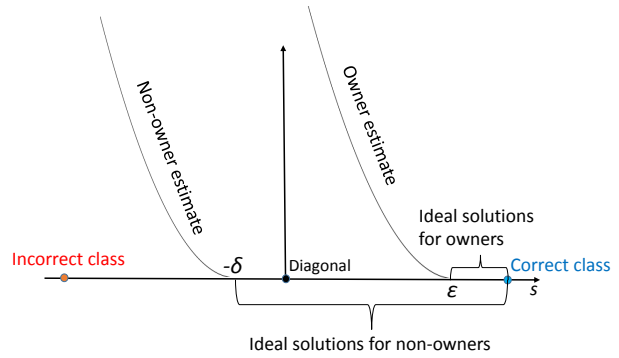


Figure 5. Soft margin error function for owners and non-owners (one owner).

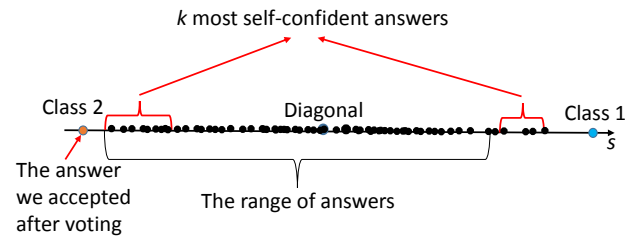


Figure 6. Interpretation of community answer with collective ownership: Voting of k most self-confident winners.

For a new example, the owners will be identified and the task will be solved by the owners following decision from Figs. 3, 6 or similar rules with distribution of responsibility between the most self-confident experts. After such labeling steps the learning cycles should follow with improvement of experts' performance (they should give the correct self-confident answers to the samples they own, and do not make large mistakes for all other examples).

This regular alternation, solving new tasks – learning – solving new task – ..., provides adaptation to the gradual change in reality and assimilation of growing data. It is not compulsory that all local experts are answering the same tasks. A sort of soft *biclustering systems* of experts and problems should be implemented to link a problem to potential experts

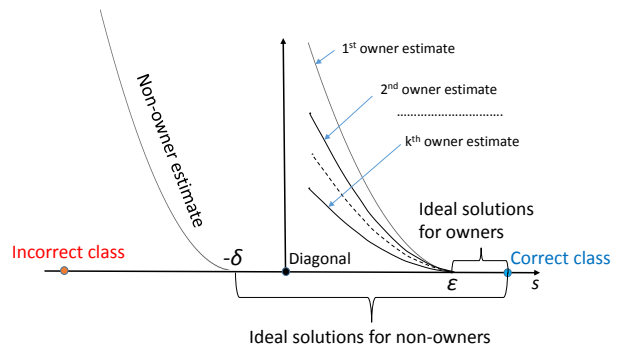


Figure 7. Soft margin error function for owners and non-owners (k owners).

and an expert to tasks it can own. Selection of experts should be done with some excess to guarantee sufficient number of selected skilled experts for correct solution. Originally [8], a version of neural network associative memory was proposed to calculate the relative weight of an expert for solution of a problem (we can call it “affinity of an expert to a problem”). A well-developed technology of recommender systems [24] includes many solutions potentially usable for recommendations of local experts to problems and problems to local experts. Implementation of a recommender system for the assignment of local experts to solve problems transforms the community of agents into hierarchical “social network” with various nodes and groups.

IV. CORRECTORS, KNOWLEDGE TRANSFER, AND INTERIORISATION

Objectives of the community self-learning are:

- Assimilation of incrementally growing data;
- Adaptation to graduate change in reality;
- Non-iterative knowledge transfer from the locally best experts to other agents;

In the community self-learning process for each sample the locally best experts (owners) find the label. After the labeling, the skills should be improved. The supervised learning of large multiagent system requires large resources. It should not destroy the previous skills and, therefore, the large labeled data base of previous tasks should be used. It can require large memory and many iterations, which involve all the local experts. It is desirable to correct the errors (or increase the level of confidence, if it is too low) without destroying of previously learned skills. It is also very desirable to avoid usage of large database and long iterative process.

Communities of AI systems in real world will work on the basis of heterogeneous networks of computational devices and in heterogeneous infrastructure. Real-time correction of mistakes in such heterogeneous systems by re-training is not always viable due to the resources involved. We can, therefore, formulate the technical requirements for the correction procedures [17]. *Corrector* should:

- be simple;
- not destroy the existing skills of the AI systems;
- allow fast non-iterative learning;
- allow correction of new mistakes without destroying of previous corrections.

Surprisingly, all these requirements can be met in sufficiently high dimensions. For this purpose, we propose to employ the concept of corrector of legacy AI systems, developed recently [16], [21] on the basis of stochastic separation theorems [20]. For high-dimensional distributions in n -dimensional space every point from a large finite set can be separated from all other points by a simple linear discriminant. The size of this finite set can grow exponentially with n . For example, for the equidistribution in an 100-dimensional ball, with probability > 0.99 every point in $2.7 \cdot 10^6$ independently chosen random points is linearly separable from the set of all other points.

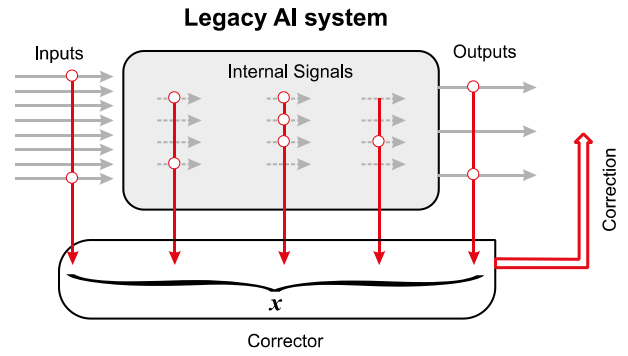


Figure 8. Corrector of AI errors.

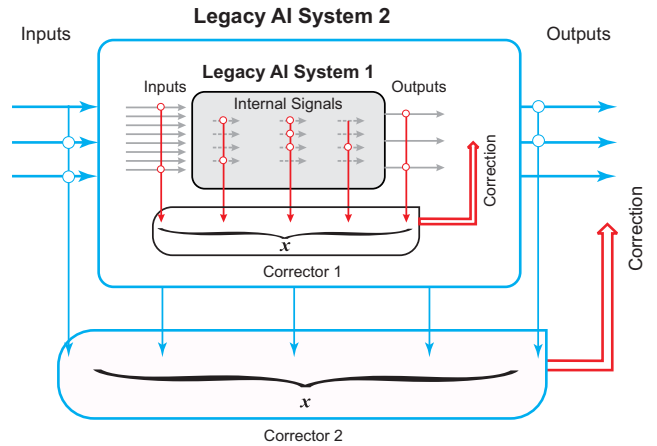


Figure 9. Cascade of AI correctors

The idea of a corrector is simple. It corrects an error of a single local expert. Separate the sample with error from all other samples by a simple discriminant. This discriminant splits the space of samples into two subsets: the sample with errors belongs to one of them, and all other samples belong to “another half”. Modify the decision rule for the set, which includes the erroneous sample. This is corrector of a legacy AI system (Fig. 8). Inputs for this corrector may include input signals, and any internal or output signal of the AI system.

One corrector can correct several errors (it is useful to cluster them before corrections). For correction of many errors, cascades of correctors are employed [17]: the AI system with the first corrector is a new legacy AI system and can be corrected further, as presented in Fig. 9. In this diagram, the original legacy AI system (shown as Legacy AI System 1) is supplied with a corrector altering its responses. The combined new AI system can in turn be augmented by another corrector, leading to a cascade of AI correctors.

Fast *knowledge transfer* between AI systems can be organised using correctors [22]. The “teacher” AI labels the samples, and a “student” AI also attempts to label them. If their decisions coincide (with the desired level of confidence) then nothing happens. If they do not coincide (or the level of confidence of a student is too low) then a corrector is created

for the student. From the technological point of view it is more efficient to collect samples with student’s errors, then cluster these samples and create correctors for the clusters, not for the individual mistakes. Moreover, new real-world samples are not compulsory needed in the knowledge transfer process. Just a large set of randomly generated (simulated) samples labeled by the teacher AI and the student AI can be used for correction of the student AI with skill transfer from the teacher AI.

Correctors assimilate new knowledge in the course of the community self-learning process (Fig. 10). After collection of a sufficiently large cascade of correctors, a local expert needs to assimilate this knowledge in its internal structure. The main reason for such *interiorisation* is restoring of the regular essentially high-dimensional structure of the distribution of preprocessed samples with preservation of skills. This process can be iterative but it is much simpler than the initial supervising learning. The local expert with the cascade of correctors becomes the teacher AI, and the same expert without correctors becomes the student AI (see Fig. 10). Available real dataset can be supplemented by the randomly simulated samples and, after iterative learning the skills from the teacher are transferred to the student (if the capacity of the student is sufficient). The student with updated skills returns to the community of local experts.

Two important subsystems are not present in Fig. 10): the manager – recommender and the ultimate auditor. The *manager – recommender* distributes tasks to local experts and local experts to tasks. It takes decisions on the basis of the previous experience of problem solving and assigns experts to problems with an adequate surplus, for reliability, and with some stochastisation, for the training of various experts and for the extension of experts’ pool.

In practice, the self-learning and self-labeling of samples performed by the selected local experts is supplemented by the labeling of samples and critics of decisions by an *ultimate auditor*. First of all, this auditor is the real practice itself: the real consequences of the decisions return to the systems. Secondly, the ultimate audit may include inspection by a qualified human or by a special AI audit system with additional skills.

V. MATHEMATICAL FOUNDATIONS OF NON-DESTRUCTIVE AI CORRECTION

A. General stochastic separation theorem

Bárány and Zoltán [25] studied properties of high-dimensional polytopes deriving from uniform distribution in the n -dimensional unit ball. They found that in the envelope of M random points *all* of the points are on the boundary of their convex hull and none belong to the interior (with probability greater than $1 - c^2$, provided that $M \leq c2^{n/2}$, where $c > 0$ in an arbitrary constant). They also show that the bound on M is nearly tight, up to polynomial factor in n . Donoho and Tanner [19] derived a similar result for i.i.d. points from the Gaussian distribution. They also mentioned that in applications it often seems that Gaussianity is not required and stated the problem of characterisation of ensembles leading to the same

qualitative effects (‘phase transitions’), which are found for Gaussian polytopes.

Recently, we noticed that these results could be proven for many other distributions, indeed, and one more important (and surprising) property is also typical: *every point in this M -point random set can be separated from all other points of this set by the simplest linear Fisher discriminant* [16], [20]. This observation allowed us to create the corrector technology for legacy AI systems [21]. We used the ‘thin shell’ measure concentration inequalities to prove these results [17], [20]. Separation by linear Fisher’s discriminant is practically most important *Surprise 4* in addition to three surprises mentioned in [19].

The standard approach assumes that the random set consists of independent identically distributed (i.i.d.) random vectors. The new stochastic separation theorem presented below does not assume that the points are identically distributed. It can be very important: in the real practice the new datapoints are not compulsory taken from the same distribution that the previous points. In that sense the typical situation with the real data flow is far from an i.i.d. sample (we are grateful to G. Hinton for this important remark). This new theorem gives also an answer to the *open problem* from [19]: it gives the general characterisation of the wide class of distributions with stochastic separation theorems (the SmAC condition below). Roughly speaking, this class consists of distributions without sharp peaks in sets with exponentially small volume (the precise formulation is below). We call this property “Smear Absolute Continuity” (or SmAC for short) with respect to the Lebesgue measure: the absolute continuity means that the sets of zero measure have zero probability, and the SmAC condition below requires that the sets with exponentially small volume should not have high probability. Below \mathbb{B}_n is a unit ball in \mathbb{R}^n and V_n denotes the n -dimensional Lebesgue measure.

Consider a *family* of distributions, one for each pair of positive integers M and n . The general SmAC condition is

Definition 1. *The joint distribution of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ has SmAC property if there exist constants $A > 0$, $B \in (0, 1)$, and $C > 0$, such that for every positive integer n , any convex set $S \in \mathbb{R}^n$ such that*

$$\frac{V_n(S)}{V_n(\mathbb{B}_n)} \leq A^n,$$

any index $i \in \{1, 2, \dots, M\}$, and any points $\mathbf{y}_1, \dots, \mathbf{y}_{i-1}, \mathbf{y}_{i+1}, \dots, \mathbf{y}_M$ in \mathbb{R}^n , we have

$$\mathbb{P}(\mathbf{x}_i \in \mathbb{B}_n \setminus S \mid \mathbf{x}_j = \mathbf{y}_j, \forall j \neq i) \geq 1 - CB^n. \quad (1)$$

We remark that

- We do not require for SmAC condition to hold for *all* $A < 1$, just for *some* $A > 0$. However, constants A , B , and C should be independent from M and n .
- We do not require that \mathbf{x}_i are independent. If they are, (1) simplifies to

$$\mathbb{P}(\mathbf{x}_i \in \mathbb{B}_n \setminus S) \geq 1 - CB^n.$$

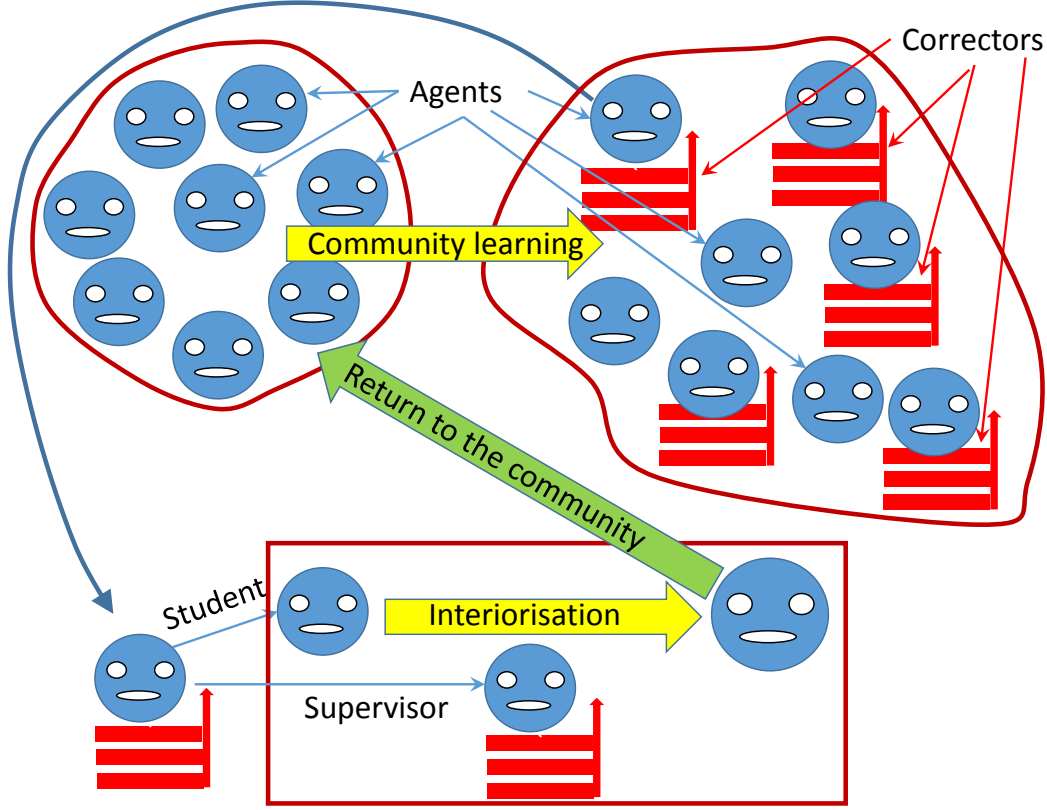


Figure 10. Community learning, self-learning, and interiorisation of knowledge.

- We do not require that \mathbf{x}_i are identically distributed. □
- The unit ball \mathbb{B}_n in SmAC condition can be replaced by an arbitrary ball, due to rescaling.
- We do not require the distribution to have a bounded support - points \mathbf{x}_i are allowed to be outside the ball, but with exponentially small probability.

The following proposition establishes a sufficient condition for SmAC condition to hold.

Proposition 1. Assume that $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ are continuously distributed in \mathbb{B}_n with conditional density satisfying

$$\rho_n(\mathbf{x}_i | \mathbf{x}_j = \mathbf{y}_j, \forall j \neq i) \leq \frac{C}{r^n V_n(\mathbb{B}_n)} \quad (2)$$

for any n , any index $i \in \{1, 2, \dots, M\}$, and any points $\mathbf{y}_1, \dots, \mathbf{y}_{i-1}, \mathbf{y}_{i+1}, \dots, \mathbf{y}_M$ in \mathbb{R}^n , where $C > 0$ and $r > 0$ are some constants. Then SmAC condition holds with the same C , any $B \in (0, 1)$, and $A = Br$.

Proof.

$$\begin{aligned} \mathbb{P}(\mathbf{x}_i \in S | \mathbf{x}_j = \mathbf{y}_j, \forall j \neq i) &= \int_S \rho_n(\mathbf{x}_i | \mathbf{x}_j = \mathbf{y}_j, \forall j \neq i) dV \\ &\leq \int_S \frac{C}{r^n V_n(\mathbb{B}_n)} dV = V_n(S) \frac{C}{r^n V_n(\mathbb{B}_n)} \\ &\leq A^n V_n(\mathbb{B}_n) \frac{C}{r^n V_n(\mathbb{B}_n)} = CB^n. \end{aligned}$$

If $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ are independent with \mathbf{x}_i having density $\rho_{i,n} : \mathbb{B}_n \rightarrow [0, \infty)$, (2) simplifies to

$$\rho_{i,n}(\mathbf{x}) \leq \frac{C}{r^n V_n(\mathbb{B}_n)}, \quad \forall n, \forall i, \forall \mathbf{x} \in \mathbb{B}_n, \quad (3)$$

where $C > 0$ and $r > 0$ are some constants.

With $r = 1$, (3) implies that SmAC condition holds for probability distributions whose density is bounded by a constant times density $\rho_n^{uni} := \frac{1}{V_n(\mathbb{B}_n)}$ of uniform distribution in the unit ball. With arbitrary $r > 0$, (3) implies that SmAC condition holds whenever ratio $\rho_{i,n}/\rho_n^{uni}$ grows at most exponentially in n . This condition is general enough to hold for many distributions of practical interest.

Example 1. (Unit ball) If $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ are i.i.d random points from the equidistribution in the unit ball, then (3) holds with $C = r = 1$.

Example 2. (Randomly perturbed data) Fix parameter $\epsilon \in (0, 1)$ (random perturbation parameter). Let $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M$ be the set of M arbitrary (non-random) points inside the ball with radius $1 - \epsilon$ in \mathbb{R}^n . They might be clustered in arbitrary way, all belong to a subspace of very low dimension, etc. Let $\mathbf{x}_i, i = 1, 2, \dots, M$ be a point, selected uniformly at random from a ball with center \mathbf{y}_i and radius ϵ . We think about \mathbf{x}_i

as “perturbed” version of \mathbf{y}_i . In this model, (3) holds with $C = 1$, $r = \epsilon$.

Example 3. (Uniform distribution in a cube) Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ be i.i.d random points from the equidistribution in the unit cube. Without loss of generality, we can scale the cube to have side length $s = \sqrt{4/n}$. Then (3) holds with $r < \sqrt{\frac{2}{\pi e}}$.

Remark 1. In this case,

$$\begin{aligned} V_n(\mathbb{B}_n)\rho_{i,n}(\mathbf{x}) &= \frac{V_n(\mathbb{B}_n)}{(\sqrt{4/n})^n} = \frac{\pi^{n/2}/\Gamma(n/2+1)}{(4/n)^{n/2}} \\ &< \frac{(\pi/4)^{n/2}n^{n/2}}{\Gamma(n/2)} \approx \frac{(\pi/4)^{n/2}n^{n/2}}{\sqrt{4\pi/n}(n/2e)^{n/2}} \leq \frac{1}{2\sqrt{\pi}} \left(\sqrt{\frac{\pi e}{2}}\right)^n, \end{aligned}$$

where \approx means Stirling’s approximation for gamma function Γ .

Example 4. (Product distribution in unit cube) Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ be independent random points from the product distribution in the unit cube, with component j of point \mathbf{x}_i having a continuous distribution with density $\rho_{i,j}$. Assume that all $\rho_{i,j}$ are bounded from above by some absolute constant K . Then (3) holds with $r < \frac{1}{K}\sqrt{\frac{2}{\pi e}}$ (after appropriate scaling of the cube).

A finite set $F \subset \mathbb{R}^n$ is called *linearly separable* if the following equivalent conditions hold.

- For each $\mathbf{x} \in F$ there exists a linear functional l such that $l(\mathbf{x}) > l(\mathbf{y})$ for all $\mathbf{y} \in F$, $\mathbf{y} \neq \mathbf{x}$;
- Each $\mathbf{x} \in F$ is an extreme point (vertex) of convex hull of F .

Below we prove the separation theorem for distributions satisfying SmAC condition. The proof is based on the following result, see [26]

Proposition 2. Let

$$V(n, M) = \frac{1}{V_n(\mathbb{B}_n)} \max_{\mathbf{x}_1, \dots, \mathbf{x}_M \in \mathbb{B}_n} V_n(\text{conv}\{\mathbf{x}_1, \dots, \mathbf{x}_M\}),$$

where *conv* denotes the convex hull. Then

$$V(n, c^n)^{1/n} < (2e \log c)^{1/2}(1 + o(1)), \quad 1 < c < 1.05.$$

Proposition 2 implies that for every $c \in (1, 1.05)$, there exists a constant $N(c)$, such that

$$V(n, c^n) < (3\sqrt{\log c})^n, \quad n > N(c). \quad (4)$$

Theorem 1. Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be a set of random points in \mathbb{R}^n from distribution satisfying SmAC condition. Then $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is linearly separable with probability greater than $1 - \delta$, $\delta > 0$, provided that

$$M \leq ab^n,$$

where

$$b = \min\{1.05, 1/B, \exp((A/3)^2)\}, \quad a = \min\{1, \delta/2C, b^{-N(b)}\}.$$

Proof. If $n < N(b)$, then $M \leq ab^n \leq b^{-N(b)}b^n < 1$, a contradiction. Let $n \geq N(b)$, and let $F = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M\}$. Then

$$\begin{aligned} \mathbb{P}(F \subset \mathbb{B}_n) &\geq 1 - \sum_{i=1}^M \mathbb{P}(\mathbf{x}_i \notin \mathbb{B}_n) \\ &\geq 1 - \sum_{i=1}^M CB^n = 1 - MCB^n, \end{aligned}$$

where the second inequality follows from (1). Next,

$$\begin{aligned} \mathbb{P}(F \text{ is linearly separable} \mid F \subset \mathbb{B}_n) \\ \geq 1 - \sum_{i=1}^M \mathbb{P}(\mathbf{x}_i \in \text{conv}(F \setminus \{\mathbf{x}_i\}) \mid F \subset \mathbb{B}_n). \end{aligned}$$

For set $S = \text{conv}(F \setminus \{\mathbf{x}_i\})$

$$\begin{aligned} \frac{V_n(S)}{V_n(\mathbb{B}_n)} &\leq V(n, M-1) \leq V(n, b^n) \\ &< \left(3\sqrt{\log(b)}\right)^n \leq A^n, \end{aligned}$$

where we have used (4) and inequalities $a \leq 1$, $b \leq \exp((A/3)^2)$. Then SmAC condition implies that

$$\begin{aligned} \mathbb{P}(\mathbf{x}_i \in \text{conv}(F \setminus \{\mathbf{x}_i\}) \mid F \subset \mathbb{B}_n) \\ = \mathbb{P}(\mathbf{x}_i \in S \mid F \subset \mathbb{B}_n) \leq CB^n. \end{aligned}$$

Hence,

$$\mathbb{P}(F \text{ is linearly separable} \mid F \subset \mathbb{B}_n) \geq 1 - MCB^n,$$

and

$$\begin{aligned} \mathbb{P}(F \text{ is linearly separable}) &\geq (1 - MCB^n)^2 \\ &\geq 1 - 2MCB^n \geq 1 - 2ab^nCB^n \geq 1 - \delta, \end{aligned}$$

where the last inequality follows from $a \leq \delta/2C$, $b \leq 1/B$. \square

B. Stochastic separation by Fisher’s linear discriminant

According to the general stochastic separation theorems there exist linear functionals, which separate points in a random set (with high probability and under some conditions). Such a linear functional can be found by various iterative methods, from the Rosenblatt perceptron learning rule to support vector machines. This existence is nice but for applications we need the non-iterative learning. It would be very desirable to have an explicit expression for separating functionals.

There exists a general scheme for creation of linear discriminants [17], [22]. For separation of single points from a data cloud it is necessary:

- 1) Centralise the cloud (subtract the mean point from all data vectors).
- 2) Escape strong multicollinearity, for example, by principal component analysis and deleting minor components, which correspond to the small eigenvalues of empiric covariance matrix.

- 3) Perform whitening (or spheric transformation), that is a linear transformation, after that the covariance matrix becomes a unit matrix. In principal components whitening is simply the normalisation of coordinates to unit variance.
- 4) The linear inequality for separation of a point \mathbf{x} from the cloud Y in new coordinates is

$$(\mathbf{x}, \mathbf{y}) \leq \alpha(\mathbf{x}, \mathbf{x}), \text{ for all } \mathbf{y} \in Y. \quad (5)$$

where $\alpha \in (0, 1)$ is a threshold, and (\bullet, \bullet) is the standard Euclidean inner product in new coordinates.

In real applied problems, it could be difficult to perform the precise whitening but a rough approximation to this transformation could also create useful discriminants (5). We will call ‘Fisher’s discriminants’ all the discriminants created non-iteratively by inner products (5), with some extension of meaning.

Formally, we say that finite set $F \subset \mathbb{R}^n$ is *Fisher-separable* if

$$(\mathbf{x}, \mathbf{x}) > (\mathbf{x}, \mathbf{y}), \quad (6)$$

holds for all $\mathbf{x}, \mathbf{y} \in F$ such that $\mathbf{x} \neq \mathbf{y}$.

Two following theorems demonstrate that Fisher’s discriminants are powerful in high dimensions.

Theorem 2 (Equidistribution in \mathbb{B}_n [16], [20]). *Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be a set of M i.i.d. random points from the equidistribution in the unit ball \mathbb{B}_n . Let $0 < r < 1$, and $\rho = \sqrt{1 - r^2}$. Then*

$$\mathbf{P} \left(\|\mathbf{x}_M\| > r \text{ and } \left(\mathbf{x}_i, \frac{\mathbf{x}_M}{\|\mathbf{x}_M\|} \right) < r \text{ for all } i \neq M \right) \geq 1 - r^n - 0.5(M - 1)\rho^n; \quad (7)$$

$$\mathbf{P} \left(\|\mathbf{x}_j\| > r \text{ and } \left(\mathbf{x}_i, \frac{\mathbf{x}_j}{\|\mathbf{x}_j\|} \right) < r \text{ for all } i, j, i \neq j \right) \geq 1 - Mr^n - 0.5M(M - 1)\rho^n; \quad (8)$$

$$\mathbf{P} \left(\|\mathbf{x}_j\| > r \text{ and } \left(\frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}, \frac{\mathbf{x}_j}{\|\mathbf{x}_j\|} \right) < r \text{ for all } i, j, i \neq j \right) \geq 1 - Mr^n - M(M - 1)\rho^n. \quad (9)$$

According to Theorem 2, the probability that a single element \mathbf{x}_M from the sample $\mathcal{S} = \{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is linearly separated from the set $\mathcal{S} \setminus \{\mathbf{x}_M\}$ by the hyperplane $l(x) = r$ is at least

$$1 - r^n - 0.5(M - 1)(1 - r^2)^{\frac{n}{2}}.$$

This probability estimate depends on both $M = |\mathcal{S}|$ and dimensionality n . An interesting consequence of the theorem is that if one picks a probability value, say $1 - \vartheta$, then the maximal possible values of M for which the set \mathcal{S} remains linearly separable with probability that is no less than $1 - \vartheta$ grows at least exponentially with n . In particular, the following holds

Corollary 1. *Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be a set of M i.i.d. random points from the equidistribution in the unit ball \mathbb{B}_n . Let $0 < r, \vartheta < 1$, and $\rho = \sqrt{1 - r^2}$. If*

$$M < 2(\vartheta - r^n)/\rho^n, \quad (10)$$

then $\mathbf{P}((\mathbf{x}_i, \mathbf{x}_M) < r \|\mathbf{x}_M\| \text{ for all } i = 1, \dots, M-1) > 1 - \vartheta$. If

$$M < (r/\rho)^n \left(-1 + \sqrt{1 + 2\vartheta\rho^n/r^{2n}} \right), \quad (11)$$

then $\mathbf{P}((\mathbf{x}_i, \mathbf{x}_j) < r \|\mathbf{x}_i\| \text{ for all } i, j = 1, \dots, M, i \neq j) \geq 1 - \vartheta$.

In particular, if inequality (11) holds then the set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is Fisher-separable with probability $p > 1 - \vartheta$.

Note that (9) implies that elements of the set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ are pair-wise almost or ε -orthogonal, i.e. $|\cos(\mathbf{x}_i, \mathbf{x}_j)| \leq \varepsilon$ for all $i \neq j$, $1 \leq i, j \leq M$, with probability larger or equal than $1 - 2Mr^n - 2M(M - 1)\rho^n$. Similar to Corollary 1, one can conclude that the cardinality M of samples with such properties grows at least exponentially with n . Existence of the phenomenon has been demonstrated in [27]. Theorem 2, Eq. (9), shows that the phenomenon is typical in some sense (cf. [28], [29]).

The linear separability property of finite but exponentially large samples of random i.i.d. elements is not restricted to equidistributions in \mathbb{B}_n . As has been noted in [21], it holds for equidistributions in ellipsoids as well as for the Gaussian distributions. Moreover, it can be generalized to product distributions in a unit cube. Consider, e.g. the case when coordinates of the vectors $\mathbf{x} = (X_1, \dots, X_n)$ in the set \mathcal{S} are independent random variables X_i , $i = 1, \dots, n$ with expectations \bar{X}_i and variances $\sigma_i^2 > \sigma_0^2 > 0$. Let $0 \leq X_i \leq 1$ for all $i = 1, \dots, n$. The following analogue of Theorem 2 can now be stated.

Theorem 3 (Product distribution in a cube [20]). *Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be i.i.d. random points from the product distribution in a unit cube. Let*

$$R_0^2 = \sum_i \sigma_i^2 \geq n\sigma_0^2.$$

Assume that data are centralised and $0 < \delta < 2/3$. Then

$$\mathbf{P} \left(1 - \delta \leq \frac{\|\mathbf{x}_j\|^2}{R_0^2} \leq 1 + \delta \text{ and } \frac{(\mathbf{x}_i, \mathbf{x}_M)}{R_0 \|\mathbf{x}_M\|} < \sqrt{1 - \delta} \right. \\ \left. \text{for all } i, j, i \neq M \right) \geq 1 - 2M \exp(-2\delta^2 R_0^4/n) \\ - (M - 1) \exp(-2R_0^4(2 - 3\delta)^2/n); \quad (12)$$

$$\mathbf{P} \left(1 - \delta \leq \frac{\|\mathbf{x}_j\|^2}{R_0^2} \leq 1 + \delta \text{ and } \frac{(\mathbf{x}_i, \mathbf{x}_j)}{R_0 \|\mathbf{x}_j\|} < \sqrt{1 - \delta} \right. \\ \left. \text{for all } i, j, i \neq j \right) \geq 1 - 2M \exp(-2\delta^2 R_0^4/n) \\ - M(M - 1) \exp(-2R_0^4(2 - 3\delta)^2/n). \quad (13)$$

In particular, under the conditions of Theorem 3, set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is Fisher-separable with probability $p > 1 - \vartheta$, provided that $M \leq ab^n$, where $a > 0$ and $b > 1$ are some constants depending only on ϑ and σ_0 .

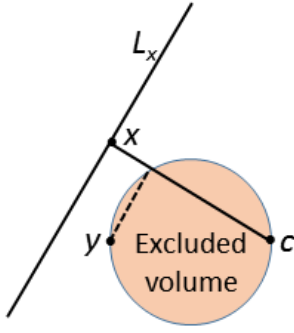


Figure 11. Point x should not belong to the filled ball (excluded volume) to be separable from y by the linear discriminant (5). Here, c is the centre of data cloud, and L_x is the hyperplane such that $(x, z) = (x, x)$ for $z \in L_x$.

The proof of Theorem 3 is based on concentration inequalities in product spaces [30]. Numerous generalisations of Theorems 2, 3 are possible for different classes of distributions, for example, for weakly dependent variables, etc.

We can see from Theorem 3 that the discriminant (5) works without precise whitening. Just the absence of strong degeneration is required: the support of the distribution contains in the unit cube (that is bounding from above) and, at the same time, the variance of each coordinate is bounded from below by $\sigma_0 > 0$.

Linear separability, as an inherent property of data sets in high dimension, is not necessarily confined to cases whereby a linear functional separates a single element of a set from the rest. Theorems 2, 3 be generalized to account for m -tuples, $m > 1$ too [17], [22].

Let us make several remarks for the general distributions. For each data point y the probability that a randomly chosen point x is *not* separated from y by the discriminant (6) (i.e. that the inequality (6) is false) is (Fig. 11)

$$p = p_y = \int_{\|z - \frac{y}{2}\| \leq \frac{\|y\|}{2}} \rho(z) dz, \quad (14)$$

where $\rho(z) dz$ is the probability measure. We need to evaluate the probability of finding a random point outside the union of N such excluded volumes. For example, for the equidistribution in a ball \mathbb{B}_n , $p_y < 1/2^n$ for all y from the ball. The probability to select a point inside the union of N ‘forbidden balls’ is less than $N/2^n$ for any position of N points y in \mathbb{B}_n . Inequalities (7), (8), and (9) are closely related to that fact.

Instead of equidistribution in a ball \mathbb{B}_n , we can take probability distributions with bounded density ρ in a ball \mathbb{B}_n

$$\rho(y) < \frac{C}{r^n V_n(\mathbb{B}_n)}, \quad (15)$$

where $C > 0$ is an arbitrary constant, $V_n(\mathbb{B}_n)$ is the volume of the ball, and radius $r > 1/2$. This inequality guarantees that the probability of each ball with radius less or equal than $1/2$ exponentially decays for $n \rightarrow \infty$. It should be stressed that in asymptotic analysis for large n the constant $C > 0$ is arbitrary but does not depend on n .

For the bounded distributions (15) the separability by linear discriminants (5) is similar to separability for the equidistributions. The proof through the estimation of the probability to avoid the excluded volume is straightforward.

For the practical needs, the number of points N is large. We have to evaluate the sum p_y for N points y . For most estimates, evaluation of the expectations and variances of p_y (14) will be sufficient, if points y are independently chosen.

If the distribution is unknown and exists just in the form of the sample of empirical points, we can evaluate $\mathbf{E}(p)$ and $\text{var}(p)$ (14) from the sample directly, without knowledge of theoretical probabilities.

Bound (15) is the special case of (3) with $r > 1/2$, and it is more restrictive: in Examples 2, 3, and 4, the distributions satisfy (3) with some $r < 1$, but fail to satisfy (15) if $r < 1/2$. Such distributions has the SmAC property and the corresponding set of points $\{x_1, \dots, x_M\}$ is linearly separable by Theorem 1, but different technique is needed to establish its Fisher-separability. One option is to estimate the distribution of p in (14). Another technique is based on concentration inequalities. For some distributions, one can prove that, with exponentially high probability, random point x satisfies

$$r_1(n) \leq \|x\| \leq r_2(n), \quad (16)$$

where $\|\bullet\|$ denotes the Euclidean norm in \mathbb{R}^n , and $r_1(n)$ and $r_2(n)$ are some lower and upper bounds, depending on n . If $r_2(n) - r_1(n)$ is small comparing to $r_1(n)$, it means that the distribution is concentrated in a thin shell between the two spheres. If x and y satisfy (16), inequality (6) may fail only if y belongs to a ball with radius $R = \sqrt{r_2^2(n) - r_1^2(n)}$. If R is much lower than $r_1(n)/2$, this method may provide much better probability estimate than (14). This is how Theorem 3 was proved in [20].

VI. SEPARATION THEOREM FOR LOG-CONCAVE DISTRIBUTIONS

A. Log-concave distributions

In [20] we proposed several possible generalisations of Theorems 2, 3. One of them is the hypothesis that for the uniformly log-concave distributions the similar result can be formulated and proved. Below we demonstrate that this hypothesis is true, formulate and prove the stochastic separation theorems for several classes of log-concave distributions. Additionally, we prove the comparison (domination) Theorem 5 that allows to extend the proven theorems to wider classes of distributions.

In this subsection, we introduce several classes of log-concave distributions and prove some useful properties of these distributions.

Let $\mathcal{P} = \{\mathbb{P}_n, n = 1, 2, \dots\}$ be a family of probability measures with densities $\rho_n : \mathbb{R}^n \rightarrow [0, \infty)$, $n = 1, 2, \dots$. Below, x is a random variable (r.v) with density ρ_n , and $\mathbb{E}_n[f(x)] := \int_{\mathbb{R}^n} f(z) \rho_n(z) dz$ is the expectation of $f(x)$.

We say that density $\rho_n : \mathbb{R}^n \rightarrow [0, \infty)$ (and the corresponding probability measure \mathbb{P}_n):

- is whitened, or *isotropic*, if $\mathbb{E}_n[\mathbf{x}] = 0$, and

$$\mathbb{E}_n[(\mathbf{x}, \theta)^2] = 1 \quad \forall \theta \in S^{n-1}, \quad (17)$$

where S^{n-1} is the unit sphere in \mathbb{R}^n , and (\bullet, \bullet) is the standard Euclidean inner product in \mathbb{R}^n . The last condition is equivalent to the fact that the covariance matrix of the components of \mathbf{x} is the identity matrix, see [32].

- is *log-concave*, if set $D_n = \{z \in \mathbb{R}^n \mid \rho_n(z) > 0\}$ is convex and $g(z) = -\log(\rho_n(z))$ is a convex function on D_n .
- is *strongly log-concave* (SLC), if $g(z) = -\log(\rho_n(z))$ is strongly convex, that is, there exists a constant $c > 0$ such that

$$\frac{g(u) + g(v)}{2} - g\left(\frac{u+v}{2}\right) \geq c\|u-v\|^2, \quad \forall u, v \in D_n.$$

For example, density $\rho_G(z) = \frac{1}{\sqrt{(2\pi)^n}} \exp(-\frac{1}{2}\|z\|^2)$ of n -dimensional standard normal distribution is strongly log-concave with $c = \frac{1}{8}$.

- has sub-Gaussian decay for the norm (SGDN), if there exists a constant $\epsilon > 0$ such that

$$\mathbb{E}_n[\exp(\epsilon\|\mathbf{x}\|^2)] < +\infty. \quad (18)$$

In particular, (18) holds for ρ_G with any $\epsilon < \frac{1}{2}$. However, unlike SLC, (18) is an asymptotic property, and is not affected by local modifications of the underlying density. For example, density $\rho(z) = \frac{1}{C} \exp(-g(\|z\|))$, $z \in \mathbb{R}^n$, where $g(t) = \frac{1}{2} \max\{1, t^2\}$, $t \in \mathbb{R}$ and $C = \int_{\mathbb{R}^n} \exp(-g(\|z\|)) dz$ has SGDN with any $\epsilon < \frac{1}{2}$, but it is not strongly log-concave.

- has sub-Gaussian decay in every direction (SGDD), if there exists a constant $B > 0$ such that inequality

$$\mathbb{P}_n[(\mathbf{x}, \theta) \geq t] \leq 2 \exp\left(-\frac{t}{B}\right)^2$$

holds for every $\theta \in S^{n-1}$ and $t > 0$.

- is ψ_α with constant $B_\alpha > 0$, $\alpha \in [1, 2]$, if

$$(\mathbb{E}_n|(\mathbf{x}, \theta)|^p)^{1/p} \leq B_\alpha p^{1/\alpha} (\mathbb{E}_n|(\mathbf{x}, \theta)|^2)^{1/2} \quad (19)$$

holds for every $\theta \in S^{n-1}$ and all $p \geq 2$.

Proposition 3. Let $\rho_n : \mathbb{R}^n \rightarrow [0, \infty)$ be an isotropic log-concave density, and let $\alpha \in [1, 2]$. The following implications hold.

$$\begin{aligned} \boxed{\rho_n \text{ is SLC}} &\Rightarrow \boxed{\rho_n \text{ has SGDN}} \Rightarrow \boxed{\rho_n \text{ has SGDD}} \Leftrightarrow \\ &\Leftrightarrow \boxed{\rho_n \text{ is } \psi_2} \Rightarrow \boxed{\rho_n \text{ is } \psi_\alpha} \Rightarrow \boxed{\rho_n \text{ is } \psi_1} \Leftrightarrow \boxed{\text{ALL}}, \end{aligned}$$

where the last \Leftrightarrow means the class of isotropic log-concave densities which are ψ_1 actually coincides with the class of all isotropic log-concave densities.

Proof. Proposition 3.1 in [33] states that if there exists $c_1 > 0$ such that $g(\mathbf{x}) = -\log(\rho_n(\mathbf{x}))$ satisfies

$$tg(u) + sg(v) - g(tu + sv) \geq \frac{c_1 ts}{2} \|u - v\|^2, \quad \forall u, v \in D_n. \quad (20)$$

for all $t, s > 0$ such that $t + s = 1$, then inequality

$$\begin{aligned} \mathbb{E}_n[f^2(\mathbf{x}) \log f^2(\mathbf{x})] - \mathbb{E}_n[f^2(\mathbf{x})] \mathbb{E}_n[\log f^2(\mathbf{x})] &\leq \\ &\leq \frac{2}{c_1} \mathbb{E}_n[\|\nabla f(\mathbf{x})\|^2] \end{aligned} \quad (21)$$

holds for every smooth function f on \mathbb{R}^n . As remarked in [33, p. 1035], “it is actually enough that (20) holds for some $t, s > 0, t + s = 1$ ”. With $t = s = 1/2$, this implies that (21) holds for every strongly log-concave distribution, with $c_1 = 8c$. By [34, Theorem 3.1], (21) holds for ρ_n if and only if it has sub-Gaussian decay for the norm, and the implication $\boxed{\rho_n \text{ is SLC}} \Rightarrow \boxed{\rho_n \text{ has SGDN}}$ follows. Also, by [37, Theorem 1(i)], if (21) holds for ρ_n , then it is ψ_2 with constant $B_2 = d/\sqrt{c}$, where d is a universal constant, hence $\boxed{\rho_n \text{ has SGDN}} \Rightarrow \boxed{\rho_n \text{ is } \psi_2}$. The equivalence $\boxed{\rho_n \text{ has SGDD}} \Leftrightarrow \boxed{\rho_n \text{ is } \psi_2}$ follows from (17) and [35, Lemma 2.2.4]. The implications $\boxed{\rho_n \text{ is } \psi_2} \Rightarrow \boxed{\rho_n \text{ is } \psi_\alpha} \Rightarrow \boxed{\rho_n \text{ is } \psi_1}$ follow from (19). Finally, [35, Theorem 2.4.6] implies that every log-concave density ρ_n is ψ_1 with some universal constant. \square

B. Fisher-separability for log-concave distributions

Below we prove Fisher-separability for i.i.d samples from isotropic log-concave ψ_α distributions, using the technique based on concentration inequalities.

Theorem 4. Let $\alpha \in [1, 2]$, and let $\mathcal{P} = \{\mathbb{P}_n, n = 1, 2, \dots\}$ be a family of probability measures with densities $\rho_n : \mathbb{R}^n \rightarrow [0, \infty)$, $n = 1, 2, \dots$, which are ψ_α with constant $B_\alpha > 0$, independent from n . Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be a set of M i.i.d. random points from ρ_n . Then there exist constants $a > 0$ and $b > 0$, which depends only on α and B_α , such that, for any $i, j \in \{1, 2, \dots, M\}$, inequality

$$(\mathbf{x}_i, \mathbf{x}_i) > (\mathbf{x}_i, \mathbf{x}_j)$$

holds with probability at least $1 - a \exp(-bn^{\alpha/2})$. Hence, for any $\delta > 0$, set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is Fisher-separable with probability greater than $1 - \delta$, provided that

$$M \leq \sqrt{\frac{2\delta}{a}} \exp\left(\frac{b}{2} n^{\alpha/2}\right). \quad (22)$$

Proof. Let \mathbf{x} and \mathbf{y} be two points, selected independently at random from the distribution with density ρ_n . [31, Theorem 1.1], (applied with $A = I_n$, where I_n is $n \times n$ identity matrix) states that, for any $t \in (0, 1)$, (16) holds with $r_1(n) = (1-t)\sqrt{n}$, $r_2(n) = (1+t)\sqrt{n}$, and with probability at least $1 - A \exp(-Bt^{2+\alpha}n^{\alpha/2})$, where $A, B > 0$ are constants depending only on α . If (16) holds for \mathbf{x} and \mathbf{y} , inequality (6) may fail only if \mathbf{y} belongs to a ball with radius $R_n = \sqrt{r_2^2(n) - r_1^2(n)} = \sqrt{4tn}$. Theorem 6.2 in [36], applied with $A = I_n$, states that, for any $\epsilon \in (0, \epsilon_0)$, \mathbf{y} does not belong to a ball with any center and radius $\epsilon\sqrt{n}$, with probability at least $1 - \epsilon^{Cn^{\alpha/2}}$ for some constants $\epsilon_0 > 0$ and $C > 0$. By selecting $t = \epsilon_0^2/8$, and $\epsilon = \sqrt{4t} = \epsilon_0/2$, we conclude that (6) holds with probability at least

$1 - 2A \exp(-Bt^{2+\alpha} n^{\alpha/2}) - (\sqrt{4t})^{Cn^{\alpha/2}}$. This is greater than $1 - a \exp(-bn^{\alpha/2})$ for some constants $a > 0$ and $b > 0$. Hence, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M$ are Fisher-separable with probability greater than $1 - \frac{M(M-1)}{2} a \exp(-bn^{\alpha/2})$. This is greater than $1 - \delta$ provided that M satisfies (22). \square

Corollary 2. *Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be a set of M i.i.d. random points from an isotropic log-concave distribution in \mathbb{R}^n . Then set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is Fisher-separable with probability greater than $1 - \delta$, $\delta > 0$, provided that*

$$M \leq ac\sqrt{n},$$

where $a > 0$ and $c > 1$ are constants, depending only on δ .

Proof. This follows from Theorem 4 with $\alpha = 1$ and the fact that all log-concave densities are ψ_1 with some universal constant, see Proposition 3. \square

We say that family $\mathcal{P} = \{\mathbb{P}_n, n = 1, 2, \dots\}$ of probability measure has *exponential Fisher separability* if there exist constants $a > 0$ and $b \in (0, 1)$ such that, for all n , inequality (6) holds with probability at least $1 - ab^n$, where \mathbf{x} and \mathbf{y} are i.i.d vectors in \mathbb{R}^n selected with respect to \mathbb{P}_n . In this case, for any $\delta > 0$, M i.i.d vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ are Fisher-separable with probability at least $1 - \delta$ provided that

$$M \leq \sqrt{\frac{2\delta}{a}} \left(\frac{1}{\sqrt{b}} \right)^n.$$

Corollary 3. *Let $\mathcal{P} = \{\mathbb{P}_n, n = 1, 2, \dots\}$ be a family of isotropic log-concave probability measures which are all ψ_2 with the same constant $B_2 > 0$. Then \mathcal{P} has exponential Fisher separability.*

Proof. This follows from Theorem 4 with $\alpha = 2$. \square

Corollary 4. *Let $\mathcal{P} = \{\mathbb{P}_n, n = 1, 2, \dots\}$ be a family of isotropic probability measures which are all strongly log-concave with the same constant $c > 0$. Then \mathcal{P} has exponential Fisher separability.*

Proof. The proof of Proposition 3 implies that \mathbb{P}_n are all ψ_2 with the same constant $B_2 = d/\sqrt{c}$, where d is a universal constant. The statement then follows from Corollary 3. \square

Example 5. *Because standard normal distribution in \mathbb{R}_n is strongly log-concave with $c = \frac{1}{8}$, Corollary 4 implies that the family of standard normal distributions has exponential Fisher separability.*

C. Domination

We say that family $\mathcal{P}' = \{\mathbb{P}'_n, n = 1, 2, \dots\}$ dominates family $\mathcal{P} = \{\mathbb{P}_n, n = 1, 2, \dots\}$ if there exists a constant C such that

$$\mathbb{P}_n(S) \leq C \cdot \mathbb{P}'_n(S) \quad (23)$$

holds for all n and all measurable subsets $S \subset \mathbb{R}^n$. In particular, if \mathbb{P}'_n and \mathbb{P}_n have densities $\rho'_n : \mathbb{R}^n \rightarrow [0, \infty)$ and $\rho_n : \mathbb{R}^n \rightarrow [0, \infty)$, respectively, then (23) is equivalent to

$$\rho_n(\mathbf{x}) \leq C \cdot \rho'_n(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (24)$$

Theorem 5. *If family \mathcal{P}' has exponential Fisher separability, and \mathcal{P}' dominates \mathcal{P} , then \mathcal{P} has exponential Fisher separability.*

Proof. For every $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ and $\mathbf{y} = (y_1, \dots, y_n) \in \mathbb{R}^n$, let $\mathbf{x} \times \mathbf{y}$ be a point in \mathbb{R}^{2n} with coordinates $(x_1, \dots, x_n, y_1, \dots, y_n)$. Let \mathbb{Q}_n be the product measure of \mathbb{P}_n with itself, that is, for every measurable set $S \subset \mathbb{R}^{2n}$, $\mathbb{Q}_n(S)$ denotes the probability that $\mathbf{x} \times \mathbf{y}$ belongs to S , where vectors \mathbf{x} and \mathbf{y} are i.i.d vectors selected with respect to \mathbb{P}_n . Similarly, let \mathbb{Q}'_n be the product measure of \mathbb{P}'_n with itself. Inequality (23) implies that

$$\mathbb{Q}_n(S) \leq C^2 \cdot \mathbb{Q}'_n(S), \quad \forall S \subset \mathbb{R}^{2n}.$$

Let $A_n \subset \mathbb{R}^{2n}$ be the set of all $\mathbf{x} \times \mathbf{y}$ such that $(\mathbf{x}, \mathbf{x}) \leq (\mathbf{x}, \mathbf{y})$. Because \mathcal{P}' has exponential Fisher separability, $\mathbb{Q}'_n(A_n) \leq ab^n$ for some $a > 0$, $b \in (0, 1)$. Hence,

$$\mathbb{Q}_n(A_n) \leq C^2 \cdot \mathbb{Q}'_n(A_n) \leq (aC^2)b^n,$$

and exponential Fisher separability of \mathcal{P} follows. \square

Corollary 5. *Let $\mathcal{P} = \{\mathbb{P}_n, n = 1, 2, \dots\}$ be a family of distributions which is dominated by a family of (possibly scaled) standard normal distributions. Then \mathcal{P} has exponential Fisher separability.*

Proof. This follows from Example 5, Theorem 5, and the fact that scaling does not change Fisher separability. \square

VII. QUASIORTHOGONAL SETS AND FISHER SEPARABILITY OF NOT I.I.D. DATA

The technique based on concentration inequalities usually fails if the data are not identically distributed, because, in this case, each \mathbf{x}_i may be concentrated in its own spherical shell. An alternative approach to prove separation theorems is to use the fact that, in high dimension, almost all vectors are almost orthogonal [28], which implies that (\mathbf{x}, \mathbf{y}) in (6) is typically “small”. Below we apply this idea to prove Fisher separability of exponentially large families in the “randomly perturbed” model described in Example 2.

Consider the “randomly perturbed” model from Example 2. In this model, Fisher’s hyperplane for separation each point \mathbf{x}_i will be calculated assuming that coordinate center is the corresponding cluster centre \mathbf{y}_i .

Theorem 6. *Let $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ be a set of M random points in the “randomly perturbed” model (see Example 2) with random perturbation parameter $\epsilon > 0$. For any $\frac{1}{\sqrt{n}} < \delta < 1$, set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is Fisher-separable with probability at least*

$$1 - \frac{2M^2}{\delta\sqrt{n}} \left(\sqrt{1 - \delta^2} \right)^{n+1} - M \left(\frac{2\delta}{\epsilon} \right)^n.$$

In particular, set $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ is Fisher-separable with probability at least $1 - v$, $v > 0$, provided that $M < ab^n$, where a, b are constants depending only on v and ϵ .

Proof. Let $\mathbf{x} \in \mathbb{R}^n$ be an arbitrary non-zero vector, and let \mathbf{u} be a vector selected uniformly at random from a unit ball. Then, for any $\frac{1}{\sqrt{n}} < \delta < 1$,

$$P\left(\left|\left(\frac{\mathbf{x}}{\|\mathbf{x}\|}, \frac{\mathbf{u}}{\|\mathbf{u}\|}\right)\right| \geq \delta\right) \leq \frac{2}{\delta\sqrt{n}} \left(\sqrt{1-\delta^2}\right)^{n+1}, \quad (25)$$

see [32, Lemma 4.1].

Applying (25) to $\mathbf{u} = \mathbf{x}_i - \mathbf{y}_i$, we get

$$P\left(\left|\left(\frac{\mathbf{x}_j}{\|\mathbf{x}_j\|}, \frac{\mathbf{u}}{\|\mathbf{u}\|}\right)\right| \geq \delta\right) \leq \frac{2}{\delta\sqrt{n}} \left(\sqrt{1-\delta^2}\right)^{n+1}, \quad j \neq i,$$

and also

$$P\left(\left|\left(\frac{\mathbf{y}_i}{\|\mathbf{y}_i\|}, \frac{\mathbf{u}}{\|\mathbf{u}\|}\right)\right| \geq \delta\right) \leq \frac{2}{\delta\sqrt{n}} \left(\sqrt{1-\delta^2}\right)^{n+1}.$$

On the other hand

$$P(\|\mathbf{x}_i - \mathbf{y}_i\| \leq 2\delta) = \left(\frac{2\delta}{\epsilon}\right)^n.$$

If none of the listed events happen, then projections of all points \mathbf{x}_j , $j \neq i$, on \mathbf{u} have length at most δ (because $\|\mathbf{x}_j\| \leq 1, \forall j$), while the length of projection of \mathbf{x}_i on \mathbf{u} is greater than δ , hence \mathbf{x}_i is separable from other points by Fisher discriminant (with center \mathbf{y}_i). Hence, the probability that \mathbf{x}_i is not separable is at most

$$\frac{2M}{\delta\sqrt{n}} \left(\sqrt{1-\delta^2}\right)^{n+1} + \left(\frac{2\delta}{\epsilon}\right)^n$$

The probability that there exist some index i such that \mathbf{x}_i is not separable is at most the same expression multiplied by M . \square

Theorem 6 is yet another illustration of why randomization and randomized approaches to learning may improve performance of AI systems (see e.g. [40], [41] for more detailed discussion on the randomized approaches and supervisory mechanisms for random parameter assignment).

Moreover, Theorem 6 shows that the cluster structure of data is not an insurmountable obstacle for separation theorems. The practical experience ensures us that combination of cluster analysis with stochastic separation theorems works much better than the stochastic separation theorems directly, if there exists a pronounced cluster structure in data. The preferable way of action is:

- Find clusters in data clouds;
- Create classifiers for distribution of newly coming data between clusters;
- Apply stochastic separation theorems with discriminant (5) for each cluster separately.

This is a particular case of the general rule about complementarity between low-dimensional non-linear structures and high-dimensional stochastic separation [17].

VIII. CONCLUSION

The continuous development of numerous automated AI systems for data mining is inevitable. Well-known AI products capable of responding, at least partially, to elements of the Big Data Challenge have already been developed by technological giants such as Amazon, IBM, Google, Facebook, SoftBank and many others. State-of-the art AI systems for data mining consume huge and fast-growing collections of heterogeneous data. Multiple versions of these huge-size systems have been deployed to date on millions of computers and gadgets across many various platforms. Inherent uncertainties in data result in unavoidable mistakes (e.g. mislabelling, false alarms, mis-detections, wrong predictions etc.) of the AI data mining systems, which require judicious use. The successful operation of any AI system dictates that mistakes must be detected and corrected immediately and locally in the networks. However, it is prohibitively expensive and even dangerous to reconfigure big AI systems in real time.

The future development of sustainable large AI systems for mining of big data requires creation of technology and methods for fast non-iterative, non-destructive, and reversible corrections of Big Data analytics systems and for fast assimilation of new skills by the network of AI. This process should exclude human expertise as far as it is possible. In this paper we presented a brief outline of an approach to Augmented AI. This approach uses communities of interacting AI systems, fast non-destructive and non-iterative correctors of mistakes, knowledge transfer between AIs, a recommender system for distribution of problems to experts and experts to problems, and various types of audit systems. Some parts and versions of this technology have been implemented and tested.

Linear Fisher's discriminant is very convenient and efficiently separates data for many distributions in high dimensions. The cascades of independent linear discriminants are also very simple and more efficient [16], [21]. We have systematically tested linear and cascade correctors with simulated data and on the processing of real videostream data [21]. The combination of low-dimensional non-linear decision rules with the high-dimensional simple linear discriminants is a promising direction of the future development of algorithms.

New stochastic separation theorems demonstrate that the corrector technology can be used to handle errors in data flows with very general probability distributions and far away from the classical i.i.d. hypothesis.

The technology of Augmented AI is necessary for prevention of the deep failure from the current peak of inflated interest to intellectual solutions (Fig. 1) into the gorge of deceived expectations.

ACKNOWLEDGMENT

The proposed corrector methodology was implemented and successfully tested with videostream data and security tasks in collaboration with industrial partners: Apical, ARM, and VMS under support of InnovateUK. We are grateful to them and personally to I. Romanenko, R. Burton, and K. Sofeikov. We are grateful to M. Gromov, who attracted our attention to

the seminal question about product distributions in a multidimensional cube, and to G. Hinton for the important remark that the typical situation with the real data flow is far from an i.i.d. sample (the points we care about seem to be from different distributions).

REFERENCES

- [1] L. Columbus, “Gartner’s hype cycle for emerging technologies, 2017 adds 5G and deep learning for first time,” *Forbes / Tech / #CuttingEdge*, August 15, 2017.
- [2] D.C. Engelbart, “Augmenting human intellect: a conceptual framework,” Stanford Research Institute, Summary Report, AFOSR-3223, Contract AF 49(638)-1024. 144 pp. 1962.
- [3] A.N. Gorban, E.M. Mirkes, V.G. Tsaregorodtsev, “Generation of explicit knowledge from empirical data through pruning of trainable neural networks.” In *Neural Networks, 1999. IJCNN’99. International Joint Conference on 1999 Jul (Vol. 6, pp. 4393-4398)*. IEEE.
- [4] A.N. Gorban, *Training Neural Networks*, Moscow: USSR-USA JV “ParaGraph”. 1990.
- [5] Y. LeCun, J.S. Denker, S.A. Solla, “Optimal brain damage,” In *Advances in Neural Information Processing Systems*, pp. 598–605, 1990.
- [6] A.N. Gorban, “Neuroinformatics and applications”, *Otkrytye Sistemy, SUBD [Open Systems, DBMS]*, 4, 1998. <https://www.osp.ru/os/1998/04/179540/>
- [7] Plato, “Phaedrus”, translated by B. Jowett, The Project Gutenberg EBook of Phaedrus, by Plato, 2013. <http://www.gutenberg.org/files/1636/1636-h/1636-h.htm>
- [8] S.E. Gilev, A.N. Gorban, E.M. Mirkes, “Small experts and internal conflicts in learning neural networks,” *Akademiia Nauk SSSR, Doklady*, vol. 320 (1), pp. 220-223, 1991.
- [9] R.A. Jacobs, M.I. Jordan, S.J. Nowlan, G.E. Hinton, “Adaptive mixtures of local experts,” *Neural computation*, vol. 3 (1), pp. 79-87, 1991.
- [10] S.E. Gilev, A.N. Gorban, D.A. Kochenov, Y.M. Mirkes, S.E. Golovenkin, S.A. Dogadin, E.V. Maslennikova, G.V. Matyushin, K.G. Nozdrachev, D.A. Rossiev, V.A. Shulman, “MultiNeuron neural simulator and its medical applications.” In *Proceedings of International Conference On Neural Information Processing, ICONIP1994, Vol. 3, No. 2, pp. 1261-1264*.
- [11] A.N. Gorban, D.A. Rossiev, S.E. Gilev, M.A. Dorrer, D.A. Kochenov, Y.M. Mirkes, S.E. Golovenkin, S.A. Dogadin, K.G. Nozdrachev, G.V. Matyushin, V.A. Shulman, A.A. Savchenko, “Medical and physiological applications of MultiNeuron neural simulator.” In: DeWitt, J.T. (ed.): *Proceedings of the WCNN’95 (World Congress on Neural Networks’95, Washington DC, July 1995)*, pp. 170–175. Lawrence Erlbaum Associate, 1996.
- [12] A.N. Gorban, D.A. Rossiev, E.V. Butakova, S.E. Gilev, S.E. Golovenkin, S.A. Dogadin, M.G. Dorrer, D.A. Kochenov, A.G. Kopytov, E.V. Maslennikova, G.V. Matyushin, Medical, psychological and physiological applications of MultiNeuron neural simulator. In *Neuroinformatics and Neurocomputers, 1995, Second International Symposium on*, pp. 7–14. IEEE, 1995.
- [13] P.C. Kainen, “Utilizing geometric anomalies of high dimension: when complexity makes computation easier.” In *Computer-Intensive Methods in Control and Signal Processing: The Curse of Dimensionality*. New York, Springer, pp. 283–294. 1997.
- [14] D.L. Donoho, “High-dimensional data analysis: The curses and blessings of dimensionality.” *AMS Math Challenges Lecture*, 1, 32 pp., 2000. <http://statweb.stanford.edu/~donoho/Lectures/AMS2000/Curses.pdf>
- [15] J. Anderson, M. Belkin, N. Goyal, L. Rademacher, J. Voss, “The More, the Merrier: the Blessing of dimensionality for learning large Gaussian mixtures,” *Journal of Machine Learning Research: Workshop and Conference Proceedings* **35**, pp. 1–30, 2014.
- [16] A.N. Gorban, I.Y. Tyukin, I. Romanenko, “The Blessing of Dimensionality: Separation Theorems in the Thermodynamic Limit,” *IFAC-PapersOnLine* **49-24**, 064–069, 2016.
- [17] A. N. Gorban, I.Y. Tyukin, “Blessing of dimensionality: mathematical foundations of the statistical physics of data.” *Phil. Trans. R. Soc. A*, 2018. DOI 10.1098/rsta.2017.0237, <https://arxiv.org/abs/1801.03421>
- [18] Ledoux M. 2001 *The Concentration of Measure Phenomenon*. (Mathematical Surveys & Monographs No. 89). Providence: AMS.
- [19] D. Donoho, J. Tanner, “Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing.” *Phil. Trans. R. Soc. A* **367**(1906), 4273–4293, 2009.
- [20] A. N. Gorban, I.Y. Tyukin, “Stochastic separation theorems.” *Neural Netw.* **94**, pp. 255–259, 2017.
- [21] A. N. Gorban, I. Romanenko, R. Burton, I.Y. Tyukin. “One-trial correction of legacy AI systems and stochastic separation theorems.” 2016 <https://arxiv.org/abs/1610.00494>
- [22] I.Y. Tyukin, A. N. Gorban, K. Sofeikov, I. Romanenko. “Knowledge transfer between artificial intelligence systems.” 2017 <https://arxiv.org/abs/1709.01547>
- [23] I.Y. Tyukin, A. N. Gorban, C. Calvo, J. Makarova, V.A. Makarov. “High-dimensional brain. A tool for encoding and rapid learning of memories by single neurons.” 2017 <https://arxiv.org/abs/1710.11227>
- [24] F. Ricci, L. Rokach, B. Shapira (Eds.), *Recommender Systems Handbook*, Springer, 2015.
- [25] I. Bárány, F. Zoltán. “On the shape of the convex hull of random points.” *Probability theory and related fields*, **77**(2), (1988): 231-240.
- [26] I. Bárány, F. Zoltán. “Approximation of the sphere by polytopes having few vertices.” *Proceedings of the American Mathematical Society*, **102**(3), (1988): 651-659.
- [27] P. Kainen, V. Kůrková. “Quasiorthogonal dimension of Euclidian spaces”. *Appl. Math. Lett.*, **6**, (1993): 7-10.
- [28] A.N. Gorban, I.Y. Tyukin, D.V. Prokhorov, K.I. Sofeikov. “Approximation with random bases: Pro et contra.” *Information Sciences*, **364**, (2016): 129-145.
- [29] V. Kůrková, M. Sanguinetti. “Probabilistic lower bounds for approximation by shallow perceptron networks.” *Neural Networks* **91**, (2017): 34-41.
- [30] M. Talagrand, “Concentration of measure and isoperimetric inequalities in product spaces.” *Publications Mathematiques de l’IHES* **81**, 73–205, 1995 (doi:10.1007/BF02699376)
- [31] O. Guédon, E. Milman. “Interpolating thin-shell and sharp large-deviation estimates for isotropic log-concave measures.” *Geometric and Functional Analysis*, **21**(5), (2011): 1043-1068.
- [32] L. Lovász, S. Vempala. “The geometry of logconcave functions and sampling algorithms.” *Random Structures & Algorithms* **30**(3), (2007): 307-358.
- [33] S. Bobkov, M. Ledoux. “From Brunn-Minkowski to Brascamp-Lieb and to logarithmic sobolev inequalities.” *Geometric & Functional Analysis* **10**(5) (2000): 1028-1052.
- [34] S. Bobkov, “Isoperimetric and analytic inequalities for log-concave probability measures.” *The Annals of Probability* **27**(4) (1999): 1903-1921.
- [35] S. Brazitikos, Giannopoulos, A., Valettas, P., Vritsiou, B. “Geometry of isotropic convex bodies.” *Mathematical Surveys and Monographs*, vol. 196. American Mathematical Soc., 2014.
- [36] G. Paouris. “Small ball probability estimates for log-concave measures.” *Transactions of the American Mathematical Society* **364**(1), (2012): 287-308.
- [37] P. Stavrakakis and P. Valettas. “On the geometry of log-concave probability measures with bounded log-Sobolev constant.” In: Ludwig M., Milman V., Pestov V., Tomczak-Jaegermann N. (eds) *Asymptotic Geometric Analysis*. Fields Institute Communications, vol 68, 2013. 359–380.
- [38] Berwald, L. “Verallgemeinerung eines Mittelwertsatzes von J. Favard für positive konkave Funktionen.” *Acta Mathematica* **79**(1) (1947): 17–37.
- [39] B. Klartag, E. Milman. “Inner regularization of log-concave measures and small-ball estimates.” In: Klartag B., Mendelson S., Milman V. (eds) *Geometric Aspects of Functional Analysis*. Lecture Notes in Mathematics, vol 2050. Springer, Berlin, Heidelberg, 2012. 267–278.
- [40] D. Wang, M. Li. “Stochastic configuration networks: Fundamentals and algorithms.” *IEEE Trans. On Cybernetics*, **47**(10) (2017): 3466–3479.
- [41] S. Scardapane, D. Wang. “Randomness in neural networks: an overview.” *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, (2017) **7**(2).