

## ГЕОМЕТРИЧЕСКИЕ АСПЕКТЫ АЛГОРИТМОВ РАБОТЫ С ВЫСОКОРАЗМЕРНЫМИ ДАННЫМИ В МАТЕМАТИЧЕСКОМ МОДЕЛИРОВАНИИ

Александр Бернштейн

[a.bernstein@skoltech.ru](mailto:a.bernstein@skoltech.ru)

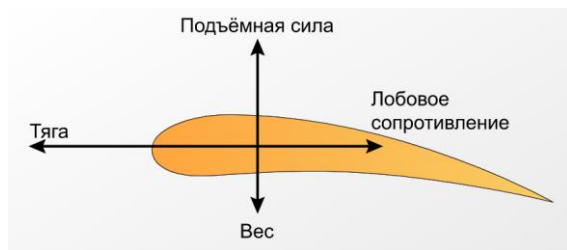
- решение задач математического моделирования – один из важнейших стимулов развития вычислительной математики и программирования  
**Трифонов Н.П. Решение на ЭВМ задач структурного анализа кристаллов (кандидатская диссертация, руководитель С.Л. Соболев)**
- моделирование на основе данных– один из современных трендов в математическом моделировании и предиктивной аналитике как частей искусственного интеллекта
- современные алгоритмы работы с данными являются «математикоемкими» и используют широкий спектр глубоких математических методов, включая дифференциально-геометрические и топологические методы

# Математические модели на первых принципах



## Аэродинамические характеристики

- подъемная сила ( $L$ ),
- лобовое сопротивление ( $C$ ), ...



## Летательный аппарат      Режим полета

- поверхность ( $S$ ),
- скорость ( $V$ )
- тяга двигателя,
- угол атаки ( $\alpha$ ), ...
- вес, ...

## Аэродинамическое проектирование: построить $S$ :

- $L \rightarrow \max$
- 3 •  $C \leq C_{\text{пред}}$

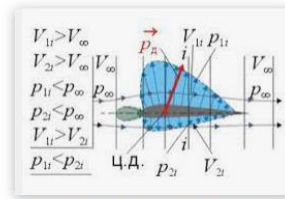
## Аэродинамические расчетные модели

$$L = F_L(S, V, \alpha)$$

$$C = F_C(S, V, \alpha)$$

$$\begin{aligned} X - \frac{1}{\rho} \frac{\partial p}{\partial x} - \frac{\partial (\bar{V}^2)}{\partial x} &= 2(\eta w - \zeta v), \\ Y - \frac{1}{\rho} \frac{\partial p}{\partial y} - \frac{\partial (\bar{V}^2)}{\partial y} &= 2(\zeta u - \xi w), \\ Z - \frac{1}{\rho} \frac{\partial p}{\partial z} - \frac{\partial (\bar{V}^2)}{\partial z} &= 2(\xi v - \eta u), \\ \frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} + \frac{\partial w}{\partial z} &= 0, \end{aligned}$$

Аэродинамика (I)



5.4. Основные законы аэродинамики

## Оптимизационная задача: построить $S$ :

$$S = \arg \max_S F_L(S, V, \alpha)$$

$$F_C(S, V, \alpha) \leq C_{\text{пред}}$$

## Математические модели на первых принципах:

### аэродинамические расчетные модели (CFD коды)

- пример: уравнения Эйлера с нерегулярной сеткой
- время расчёта порядка 10 ч.

## EU FP7 Project FFAST (2010 – 2013): Future Fast Aeroelastic Simulation Technologies

- about **100 000** design layout releases are required for a new passenger aircraft
- each layout should be analyzed for 15 to 20 control and environment scenarios

**Number of conditions (cases)** that are required in the Critical Loads Analysis of a large civil aircraft:

50 - flight points (altitude and speed)

100 - mass cases (loaded weight and weight distribution)

4 - control laws

10 - control surface configurations

5 - manoeuvres and gusts (gradient lengths)

$$\text{Total analysis cost/time} \approx N_{\text{case}} \times C_{\text{sim}} \times \frac{N_{\text{sim}}}{N_{\text{span}}}$$

$N_{\text{case}}$  = Number of test cases required to define the envelope of critical loads

$C_{\text{sim}}$  = Cost/time of each full order simulation

$N_{\text{sim}}$  = Number of full order simulations required to construct each reduced order model

$N_{\text{span}}$  = Number of loads cases spanned by each reduced order model



**TOTAL NUMBER OF REQUIRED CASES:  $N_{\text{case}} = 10\,000\,000$**

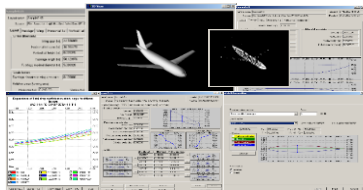
- **Boeing:** a total of **800,000 cray-hours** were spent for simulation of all the systems in new B-787
- Упрощение моделей, основанных на первых принципах («не учет» некоторых физических феноменов, поиск приближенных решений, и т.п.), не решает проблему кардинально

# Построение аэродинамической модели по данным (1)

- **Family of Fast Aerodynamic Surrogate Models** developed in **Russian Academy of Sciences**
- Implemented in **AIRBUS Engineering Tools for Aerodynamic design**

2004 – 2007

- **РАН:** Центр программных технологий РАН → ФИЦ ИУ РАН (ИСА), ИППИ РАН
- **ЦАГИ**



**Airbus experts:** application of such surrogate models “**provides the reduction of up to 10% of lead time and cost in several areas of the aircraft design process**”

## 1) Генерация данных $\{(S_i, V_i, \alpha_i, L_i \approx F_L(S_i, V_i, \alpha_i), C_i \approx F_C(S_i, V_i, \alpha_i)), i = 1, 2, \dots, N\}$

- построена математическая параметрическая модель 3D-поверхности самолета
- сгенерирована множество  $\{S\}$  поверхностей (~ 6000, включая ~ 27 000 аэродинамических профилей крыла)

**технологии генеративного искусственного интеллекта**

- с помощью CFD-кода (основан на решении уравнений полного потенциала)

**- ЦАГИ**

- для каждой поверхности  $S$

- для разных значений скорости ( $V$ ) и угла атаки ( $\alpha$ )

были вычислены подъемная сила  $L \approx F_L(S, V, \alpha)$ , сопротивление  $C \approx F_C(S, V, \alpha)$ , и другие характеристики (~ 20)

# Построение аэродинамической модели по данным (2)

## 2) Построение по данным неизвестных зависимостей (регрессионные методы/машинное обучение):

данные  $\{(S, V, \alpha); (L \approx F_L(S, V, \alpha), C \approx F_C(S, V, \alpha), \dots)\} \rightarrow$  зависимости  $L = F_L(S, V, \alpha), C = F_C(S, V, \alpha), \dots$

**Построенная по данным аэродинамическая модель – «заменитель» математической модели «на первых принципах» (суррогатная модель, метамодель)**

### Построенная суррогатная модель:

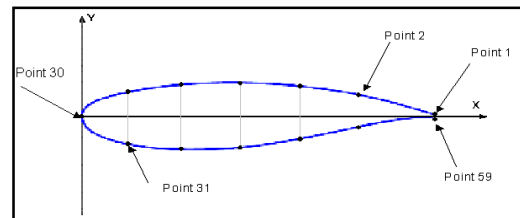
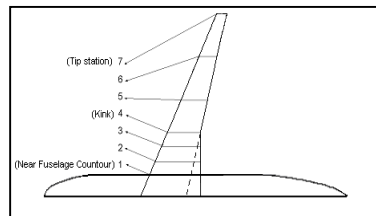
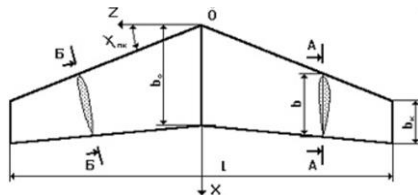
- ускорение вычислений в 360 000 раз:** время расчета **20 аэродинамических характеристик** одной компоновки для **~ 65 различных режимов полета**:
  - CFD-код (уравнение полного потенциала) ~ **2 CPU hours**,
  - построенная суррогатная модель ~ **20 CPU milliseconds**
- относительная ошибка ~ 1%**
- A.V. Bernstein, A.P. Kuleshov, Yu.N. Sviridenko, V.V. Vishinsky. Fast Aerodynamic Model for Design Technology. Workbook «West-East High Speed Flow Field Conference» (November 19-22, **2007**, Moscow, Russia). 2007, 125-126
- А.В. Бернштейн, В.В. Вышинский, А.П. Кулешов, Ю.Н. Свириденко. Быстрый метод аэродинамического расчета для задач проектирования. Труды ЦАГИ “Применение искусственных нейронных сетей в задачах прикладной аэродинамики” 2678, **2008**, с. 35—45.

Reynolds number	Re = $3 \times 10^6$	Re = $3 \times 10^7$
Lift coefficient (CL)	0.848	0.796
Lateral Force coefficient (CY)	0.843	0.721
Profile drag coefficient (CDFP)	0.666	0.776
Friction drag coefficient (CDF)	0.319	0.234
Total drag coefficient (CD)	1.039	1.547
Derivatives of CL with respect to Angle of Attack	1.076	1.094

**ИИ- интерпретация построенной модели**

# Работа с многомерными данными (1)

Поверхность самолета - вектор высокой размерности



## Профили крыла:

- профиль – многомерный вектор  $p \sim 50 \div 200$  (ЦАГИ:  $p = 59$ )
- 7 профилей крыла:  $59 \times 7 = 413$

## Проблемы построения математических моделей по данным

- мир многомерен – реальные объекты описываются векторами высокой размерности
- многомерные данные «трудны» для анализа и обработки (проклятие размерности, феномен «пустого пространства», ...)

## Работа с многомерными данными (2)

**Проклятие размерности** (Ричард Беллман, 1961): экспоненциальный рост

- числа необходимых экспериментальных данных в зависимости от размерности пространства при решении вычислительных, оптимизационных, вероятностно-статистических задач регрессии, распознавания образов, классификации, машинного обучения, дискриминантного анализа, и др.;
- числа вариантов в комбинаторных задачах в зависимости от размерности исходных данных, что приводит к соответствующему росту сложности переборных алгоритмов

### Иллюстративный пример

- $\psi(X): [0, 1]^p \rightarrow \mathbb{R}^1$  – неизвестная Липшицева функция ( $\psi \in \mathbf{Lip}$ ) на единичном  $p$ -мерном кубе
- $\hat{\psi}(X)$  – произвольная оценка построенная по зашумленной выборке  $\{(X_i, \psi(X_i)), i = 1, 2, \dots, n\}$

**Неасимптотическая нижняя граница точности** (Ибрагимов, Хасьминский (1979), Stone (1982)):

$$\sup_{\psi \in \mathbf{Lip}} E(\psi(X) - \hat{\psi}(X))^2 \geq \text{Const} \times n^{-2/(2+p)} \quad (\text{константа не зависит от } n \text{ и } p)$$

$p = 10$ :  $n = 10\,000$  измерений для достижения **заданной точности**

$p = 20$ :  $n \sim 10\,000\,000$  измерений для достижения **той же точности**



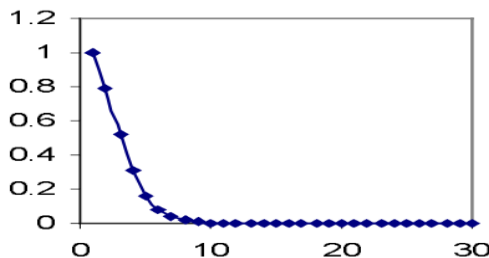
## Феномен «пустого пространства»

### 1. Объемы гиперкуба и вписанного в него шара

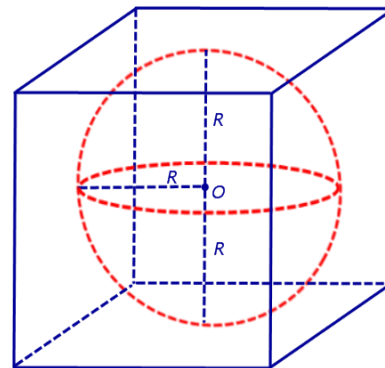
$C(R, p) = [-R, R]^p$  -  $p$ -мерный куб  $V(C(R, p)) = (2R)^p$

$B(R, p) = \{x \in \mathbb{R}^p: |x| \leq R\}$  -  $p$ -мерный шар  $V(B(R, p)) = \frac{\pi^{p/2} \times R^p}{\Gamma(\frac{p}{2} + 1)}$

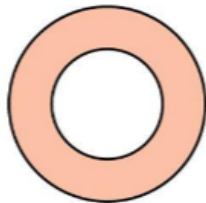
$$\lim_{p \rightarrow \infty} \frac{V(B(R, p))}{V(C(R, p))} = 0$$



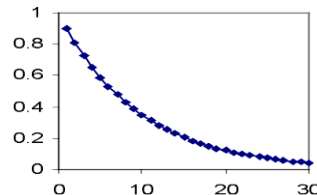
0,1      для  $p = 6$   
0,0025      для  $p = 10$



### 2. Объем тонкой сферической оболочки



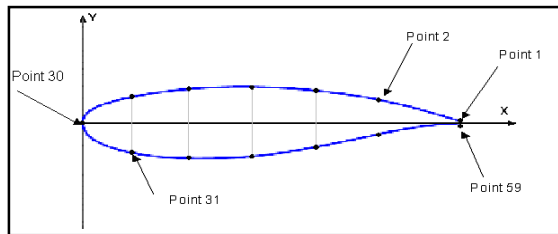
$$\lim_{p \rightarrow \infty} \frac{V(B(1, p)) - V(B(1-\varepsilon, p))}{V(B(1, p))} = 1$$



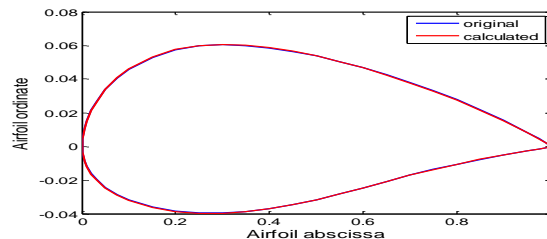
## Работа с многомерными данными (4)

Концентрация мер, разреженная структура, **низкоразмерная внутренняя структура «реальных» данных**:

- носители **реальных** данных лежат на **малой** части «объемлющего» многомерного пространства,
- могут быть описаны/параметризованы **небольшим числом параметров**



59 → 6



A. Bernstein, E. Burnaev, S. Chernova, F. Zhu, N. Qin. **Comparison of Three Geometric Parameterization methods and Their Effect on Aerodynamic Optimization**. Eurogen-2011, Sira, Italy, September 14-16, 2011, pp. 758–772

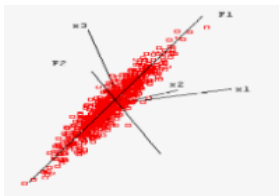
- вычислительно эффективные алгоритмы** работы с многомерными данными (предиктивная аналитика, машинное обучение, Искусственный интеллект, ...) в явном или неявном виде **используют низкоразмерную внутреннюю структуру носителей данных**
- нахождение** по данным **внутренней структуры носителей данных** и их использование в алгоритмах – отдельное научное **направление исследований**

# Методы работы с многомерными данными (20 век)

## Многомерный статистический анализ

### Снижение размерности

PCA - метод главных компонент (1903)

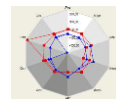
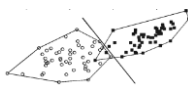


Регрессионный анализ (1795/1805/...)

Дискриминантный анализ (1936)

Дисперсионный анализ (1918)

Факторный анализ (1904), ...



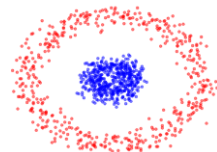
Линейные модели данных  
Гауссовские распределения

PCA → Спектральное разложение матриц/тензоров → «тензорные поезда» →

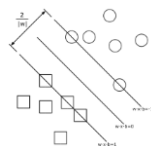
Эвристические алгоритмы работы с «нелинейными» многомерными данными

Метод потенциальных функций (1964) → Kernel PCA (1998)

(методы гильбертова пространства с воспроизводящим ядром)



Метод опорных векторов (1963)



Искусственные нейронные сети/автоэнкодеры (1974 → 1991 → ...)

Отсутствие математических  
моделей многомерных данных

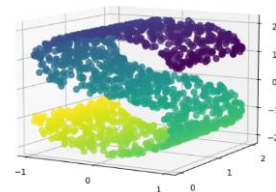


# Новые «математикоемкие» методы работы с многомерными данными

Моделирование многообразий (Manifold Learning) (2000)

Топологический анализ данных (Topological Data Analysis) (2009), ...

**Нелинейная модель многомерных данных:** многомерные данные лежат на (или вблизи) нелинейного многообразия невысокой размерности (**Многообразия данных**), вложенного в высокоразмерное пространство наблюдений Seung, Lee: The Manifold Ways of Perception. Science, 2000



**Гипотеза многообразия:** выборка  $p$ -мерных данных  $\mathbf{X}_n = \{X_1, X_2, \dots, X_n\} \subset \mathbb{R}^p$ :

- лежит на **неизвестном** Многообразии данных  $\mathbf{M}$  вложенном в  $p$ -мерное пространство ( $\mathbf{M} \subset \mathbb{R}^p$ )
- многообразие имеет **неизвестную** внутреннюю размерность  $q = \text{Dim } \mathbf{M}$
- получена в соответствии с **неизвестным** вероятностным распределением  $\mu$  на  $\mathbf{M}$

# Моделирование многообразий: задачи

По заданной выборке решаются различные задачи:

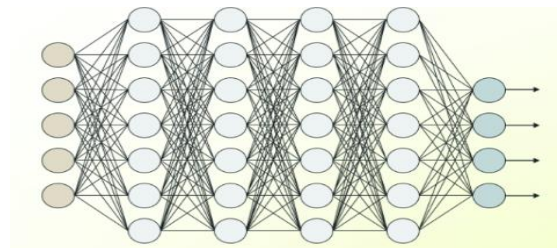
- **снижение размерности:** при выбранном  $q$ , построить низкоразмерные описания данных

$$\mathbf{X}_n = \{X_1, X_2, \dots, X_n\} \subset \mathbb{R}^p \rightarrow \mathbf{Y}_n = \{y_1, y_2, \dots, y_n\} \subset \mathbb{R}^q:$$

(построить отображение вложения  $h: X \in \mathbf{M} \subset \mathbb{R}^p \rightarrow y = h(X) \in \mathbb{R}^q$  с заданными свойствами)

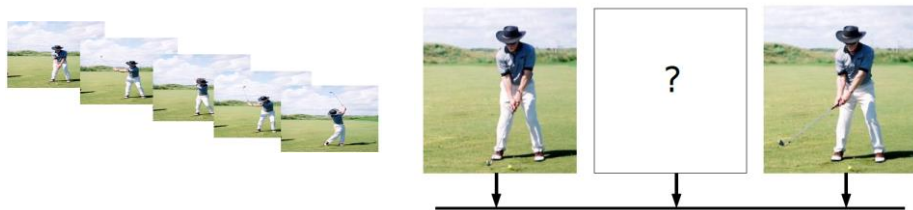
- оценить внутреннюю размерность  $q = \text{Dim } \mathbf{M}$
- построить оценку  $\hat{\mathbf{M}}$  многообразия данных  $\mathbf{M}$
- построить оценки различных элементов многообразия (касательных пространств  $T_{\mathbf{M}}(X)$ , риманова тензора  $Q_{\mathbf{M}}(X)$ , ..., в различных точках многообразия  $X \in \mathbf{M}$ )
- оценить вероятностное распределение  $\mu$  на  $\mathbf{M}$
- решение различных статистических задач (регрессии, классификации, ...), носители данных которых лежат на неизвестном Многообразии данных

**Алгоритмы снижения размерности применяются не только к исходным данным, но и к данным, возникающим на промежуточных слоях глубоких нейронных сетей**



# Моделирование многообразий: пример 1 (геодезические линии)

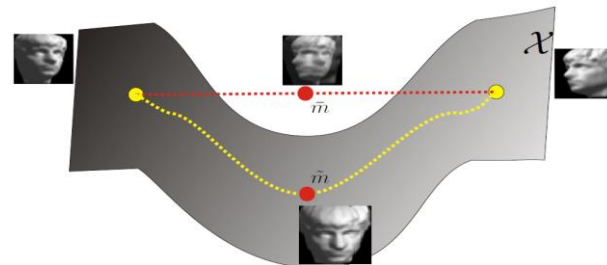
## Задачи повышения четкости изображений



## Линейная интерполяция



Линейные методы не учитывают, что реальные изображения «живут» на **нелинейных** низкоразмерных структурах



нелинейные методы

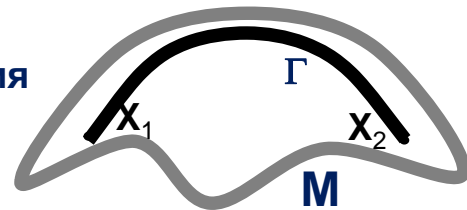


# Моделирование многообразий: примеры (1)

Геодезические линии на многообразии – кратчайшие пути → расстояния

ISOMetric MAPing (ISOMAP)

Tehenbaum, de Silva, Langford: A global geometric framework for nonlinear dimensionality reduction, 2000



Многомерное шкалирование: сохранение расстояний в низкоразмерных данных

$$\Delta_{\text{MetricMDS}} = \sum_{i,j=1}^n \left( \|X_i - X_j\|^2 - \|y_i - y_j\|^2 \right)^2$$

**ISOMAP:** евклидовы расстояния заменяются длинами геодезических  $D(X_i, X_j)$  (оцениваются)

$$\Delta_{\text{ISOMAP}} = \sum_{i,j=1}^n \left( \left( D(X_i, X_j) \right)^2 - \|y_i - y_j\|^2 \right)^2$$

ISOMAP-решение: гильбертово пространство со специфическим воспроизводящим ядром

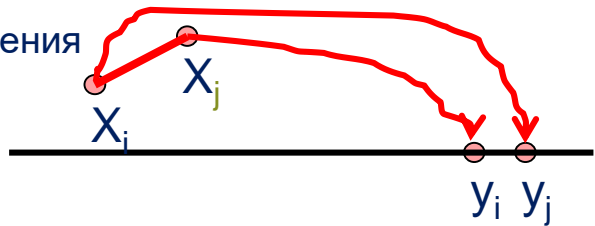
$$K_{\text{ISOMAP}}(X_i, X_j) = - D^2(X_i, X_j) - \sum_{k=1}^n D^2(X_k, X_i) - \sum_{k=1}^n D^2(X_j, X_k) + \sum_{k,s=1}^n D^2(X_s, X_k)$$

LogMap (2005), Riemannian manifold learning (2006, 2008), ...

# Моделирование многообразий: пример 2 (уравнения на многообразиях)

$\mathbf{X}_n$  – обучающая выборка,  $h: \mathbf{M} \subset \mathbb{R}^p \rightarrow \mathbb{R}^q$  – желаемое отображение вложения

$$|h(X') - h(X)| \leq |\nabla_{\mathbf{M}} h(X)| \times |X' - X| + o(|X' - X|) \quad \nabla h(X) = \begin{pmatrix} \frac{\partial h}{\partial X_1} \\ \dots \\ \frac{\partial h}{\partial X_p} \end{pmatrix}$$



$$F(h) = \int_{\mathbf{M}} |\nabla_{\mathbf{M}} h(X)|^2 \text{mes}(dX) \rightarrow \min$$

Теорема Стокса:  $F(h) = \int_{\mathbf{M}} (h \times \Delta_{\mathbf{M}} h)(X) \text{mes}(dX)$

$$\Delta_{\mathbf{M}}(h): h(X) \rightarrow -\text{div}(\nabla_{\mathbf{M}} h(X)) = -\sum_{k=1}^p \frac{\partial^2 h(X)}{\partial X_k^2}$$

- оператор Лапласа-Бельтрами

Компоненты оптимального вложения  $h(X) = \begin{pmatrix} h_1(X) \\ \dots \\ h_q(X) \end{pmatrix}$  состоят из  $q$  **собственных функций оператора**

**Лапласа-Бельтрами**, отвечающих  $q$  минимальным собственным числам

**Laplacian Eigenmaps** (Belkin, Niyogi: Laplacian Eigenmaps for Dimensionality Reduction and Data Representation, 2003)

- строится  $n \times n$  **матрица Лапласиана**  $L_{LE}(\mathbf{X}_n)$  – выборочный аналог оператора Лапласа-Бельтрами
- $n$ -мерные векторы  $(h_k(X_1), h_k(X_2), \dots, h_k(X_n))^T$ , состоящие из  $k$ -ых компонент векторов  $\{h(X_i)\}$ ,  $k = 1, 2, \dots, q$ , являются собственными векторами матрицы Лапласиана, отвечающих  $q$  минимальным собственным числам



# Моделирование многообразий: пример 3 (касательные расслоения)

Оценка обобщающей способности методов снижения размерности в произвольной точке  $X \in M$

$$\gamma_- \times d_{P,2}(T_M(X), T_{\hat{M}}(\hat{X})) \leq \delta^*(X) \leq \gamma_+ \times d_{P,2}(T_M(X), T_{\hat{M}}(\hat{X}))$$

- $\delta^*(X)$  – некоторая характеристика обобщающей способности
- $T_M(X)$  - касательное пространство к многообразию  $M$  данных в точке  $X$
- $T_{\hat{M}}(\hat{X})$  - касательное пространство к построенной по данным оценке  $\hat{M}$  многообразия данных в «восстановленной» точке  $\hat{X}$
- $d_{P,2}(L, \hat{L})$  – расстояние между  $q$ -мерными линейными пространствами в  $R^p$  (проекционная 2-норма на многообразия Грассмана)
- $\gamma_-$  и  $\gamma_+$ ,  $0 \leq \gamma_- \leq \gamma_+ \leq 1$ , - константы

Bernstein, Kuleshov. **Manifold Learning: generalizing ability and tangent proximity**, 2013

**Задача оценивания касательного расслоения многообразия данных** - построение алгоритма снижения размерности (**Grassmann&Stiefel Eigenmaps**), обеспечивающий близости между

- $M \approx \hat{M}$  - между точками  $X \in M$  и  $\hat{X} \in \hat{M}$  исходного и восстановленного по данным многообразий  $M$  и  $\hat{M}$
- $T_M(X) \approx T_{\hat{M}}(\hat{X})$  - между касательными пространствами  $T_M(X)$  и  $T_{\hat{M}}(\hat{X})$  к многообразиям  $M$  и  $\hat{M}$

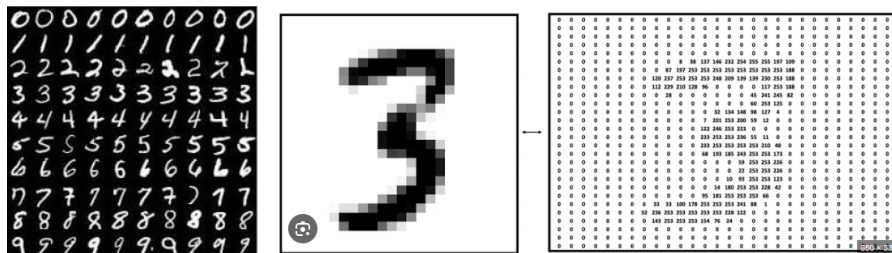
$TB(M) = \{(X, T_M(X), X \in M)\}$  - касательное расслоение многообразия данных  $M$

Bernstein, Kuleshov. **Tangent Bundle Manifold Learning via Grassmann & Stiefel Eigenmaps**, 2012

$$\mathbf{X}_n = \{X_1, X_2, \dots, X_n\} \subset \mathbb{R}^p \rightarrow \mathbf{Y}_n = \{y_1, y_2, \dots, y_n\} \subset \mathbb{R}^q$$

- $\mathbf{X}_n \rightarrow \Gamma(\mathbf{X}_n)$  – граф с вершинами в точках выборки, ребра связывают каждую вершину с  $k$  ближайшими соседями, веса ребер определяются нормированными ядрами «теплопроводности»:  
 $(i, j) \rightarrow P_{ij} = \text{const} \times \exp \left\{ -\frac{1}{t} |X_i - X_j|^2 \right\} \rightarrow$  вероятности появления ребер ( $\sum_{i,j} P_{ij} = 1, t = 2\sigma^2$ )
- $\mathbf{Y}_n \rightarrow \Gamma(\mathbf{Y}_n)$  – граф строится подобным образом по «искомым» низкоразмерным представлениям с вероятностями появления ребер  $\{Q_{ij}\}$
- **Stochastic Neighbor Embedding (SNE):** G.E. Hinton and S.T. Roweis. Stochastic Neighbor Embedding. In: Advances in Neural Information Processing Systems, vol. 15, pp. 833–840, Cambridge, MA, USA, 2002
- **Symmetric SNE (SSNE):** J.A. Cook, I. Sutskever, A. Mnih, and G.E. Hinton. Visualizing similarity data with a mixture of maps. In: Proc. of the 11th International Conference on Artificial Intelligence and Statistics, vol. 2, pp. 67–74, 2007  
Минимизация расстояния Кульбака-Лейблера:  $L_{\text{SSNE}}(\mathbf{Y}_n | \mathbf{X}_n) = \text{KL}(P \parallel Q) = \sum_i \sum_j P_{ij} \log_2 \frac{P_{ij}}{Q_{ij}} \rightarrow \min$
- **t-distributed SNE (t-SNE):** L.J.P. van der Maaten, G.E. Hinton. Visualizing High-Dimensional Data Using t-SNE. Journal of Machine Learning Research, 9 (Nov), pp. 2579-2605, 2008  
 $L_{\text{t-SNE}}(\mathbf{Y}_n | \mathbf{X}_n) = \text{KL}(P \parallel Q)$ , распределение  $Q$  имеет t-распределение Стьюдента
- **Uniform Manifold Approximation and Projection (UMAP):** L. McInnes, et. al. UMAP: Uniform Manifold Approximation and Projection. The Journal of Open Source Software, 3(29), p. 861, 2018

## MNIST

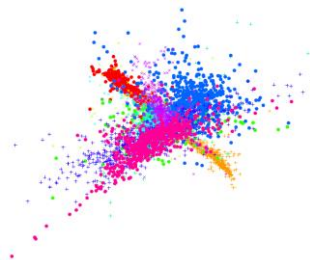


Построение двумерных представлений  
756-мерных векторов по 6 000 цифрам

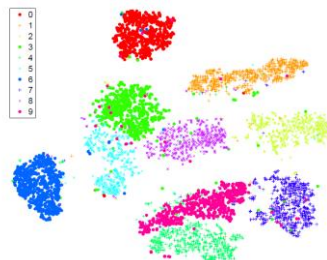
ISOMAP



Laplacian  
Eigenmaps



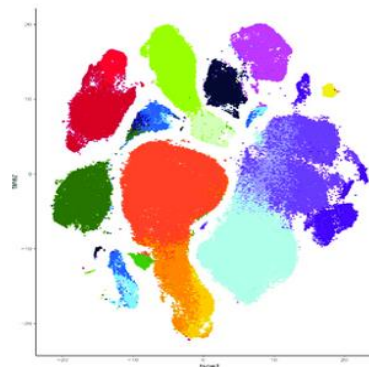
t-SNE



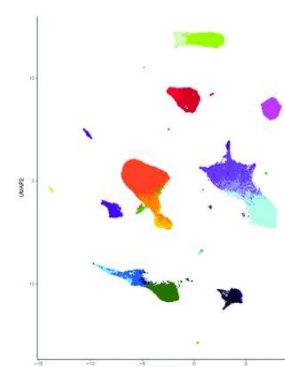
$28 \times 28 = 756$  пикселей

Рукописная цифра – 756-мерный вектор

t-SNE



UMAP



# Моделирование многообразий: исследования

- Bernstein A.V., Kuleshov A.P. **Low-Dimensional Data Representation in Data Analysis. Lecture Notes in Artificial Intelligence**, vol. 8774 “**Artificial Neural Networks in Pattern Recognition**”, Springer International Publishing, Switzerland, pp. 47-58, 2014
- A.P. Kuleshov, A.V. Bernstein. **Statistical Learning on Manifold-valued Data. Lecture Notes in Artificial Intelligence Series**, vol. 9729 ‘**Machine Learning and Data Mining in Pattern Recognition**,’ Switzerland, Springer International Publishing, pp. 311–325, 2016.
- Kuleshov, A., Bernstein, A. **Nonlinear multi-output regression on unknown input manifold. Annals of Mathematics and Artificial Intelligence**, vol. 81, №1-2, pp. 209-240, 2017
- Kuleshov, A., Bernstein, A., Burnaev, E., Yanovich, Yu. **Machine Learning in Appearance-based Robot Self-localization. Proceedings of the 16<sup>th</sup> International IEEE Conference on Machine Learning and Applications (ICMLA-2017)**, IEEE Computer Society, Mexico, December 18-21, pp. 106-112, 2017



NEW ORLEANS IEEE International Conference on Data Mining

## High-dimensional Density Estimation for Data Mining Tasks

Alexander Kuleshov, Alexander Bernstein, Yury Yanovich  
Skolkovo Institute of Science and Technology  
Moscow, Russia

November 18, 2017



ANNPR 2018

8th IAPR TC3 Workshop on Artificial Neural Networks in Pattern Recognition  
September 19-21, 2018, Siena, Italy

## Manifold learning regression with non-stationary kernels

Alexander Bernstein, Alexander Kuleshov, Evgeny Burnaev  
Skolkovo Institute of Science and Technology  
Moscow, Russia

September 21, 2018



## Manifold Learning in Machine Vision and Robotics

Alexander Bernstein  
Skolkovo Institute of Science and Technology  
Moscow, Russia

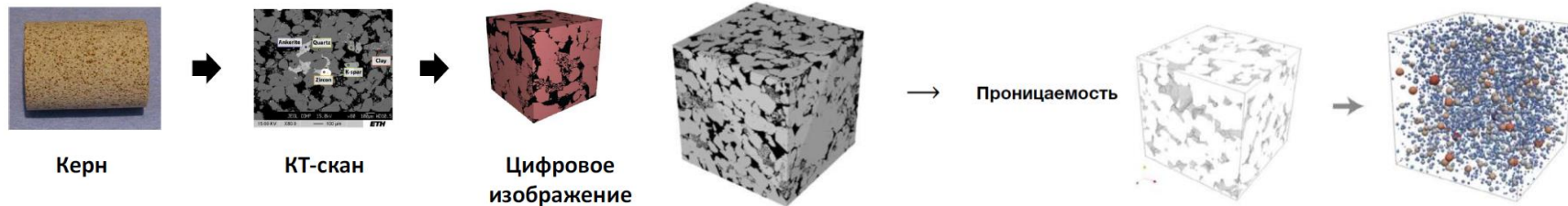
**Skoltech**

Skolkovo Institute of Science and Technology

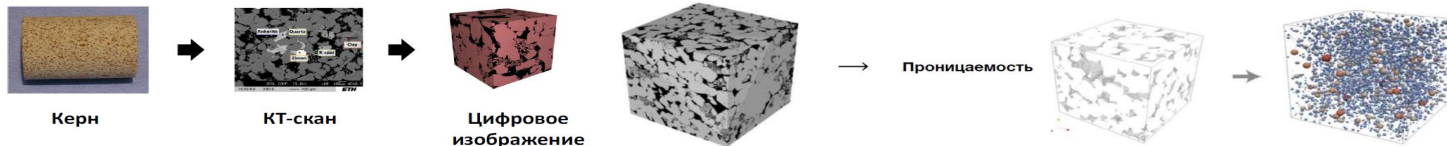
# Топологический Анализ Данных: мотивация

- в задачах ИИ (предиктивной аналитики, машинного обучения, суррогатного моделирования, ...)
  - обучающая информация представлена в виде наборов **дискретных данных/облаков точек**
  - информация в компьютере представлена **в дискретной форме**
- **топология** - часть математики/геометрии, изучающая в самом общем виде **явление непрерывности**, а также свойства обобщённых геометрических объектов, не меняющиеся при малых деформациях и не зависящие от способа их задания

## Предсказания проницаемости пористых и гранулярных сред

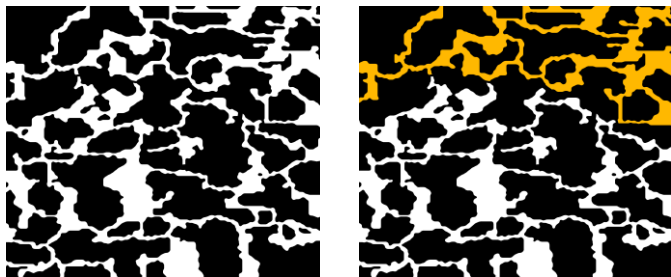


## Предсказания проницаемости пористых и гранулярных сред



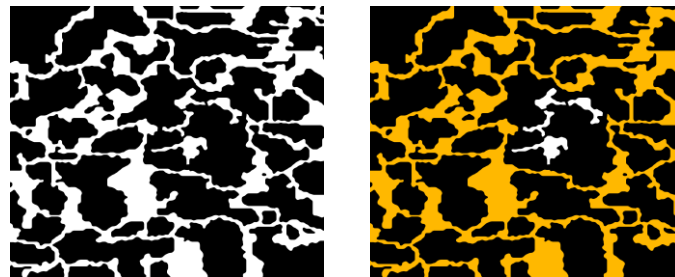
- Геометрические характеристики (пористость, извилистость, ...)
- Функционалы Минковского (площади, периметры, эйлеровы характеристики, ...)

### Образец А



Нулевая проницаемость (в вертикальном направлении)

### Образец Б



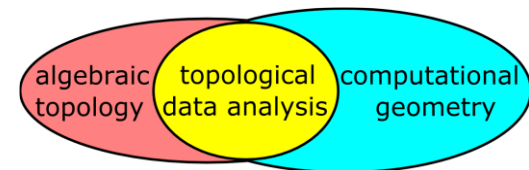
Не нулевая проницаемость

Значения  
функционалов Минковского

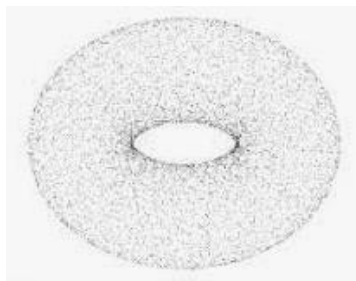
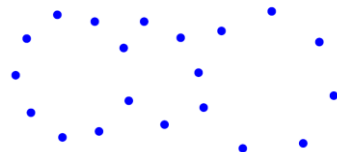
	Образец А	Образец Б	Разница, %
Площадь А	1593603	1593394	0.01
Периметр S	58087	58112	0.04
Эйлерова характеристика $\chi$	-32	-32	—

**Топологический анализ данных: объединение дискретных фрагментов в непрерывные образы (глобальные структуры) и вычисление их различных характеристик**

- замена набора элементов данных некоторым семейством **симплициальных комплексов** в соответствии с параметром близости.
- анализ этих топологических комплексов с помощью алгебраической топологии (новой теорией **персистентных гомологий**)
- перекодировка устойчивой гомологии набора данных в параметризованную версию **чисел Бетти** (баркоды, персистентные диаграммы, ...)



**Облако точек** (пиксели/воксели изображения)  
критерий близости



- выявление **скрытых структур**
- построение новых **признаков** (дескрипторов),  
называемых **топологическими инвариантами**, ...



# Топологический Анализ Данных: пример (фильтрация)

## Числа Бетти:

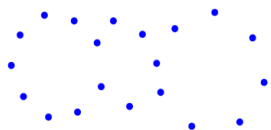
- $\beta_0$  – число связных компонент (0-мерные гомологии)
- $\beta_1$  – число «щелей» (1-мерные гомологии)



$$\beta_0 = 1, \beta_1 = 2$$

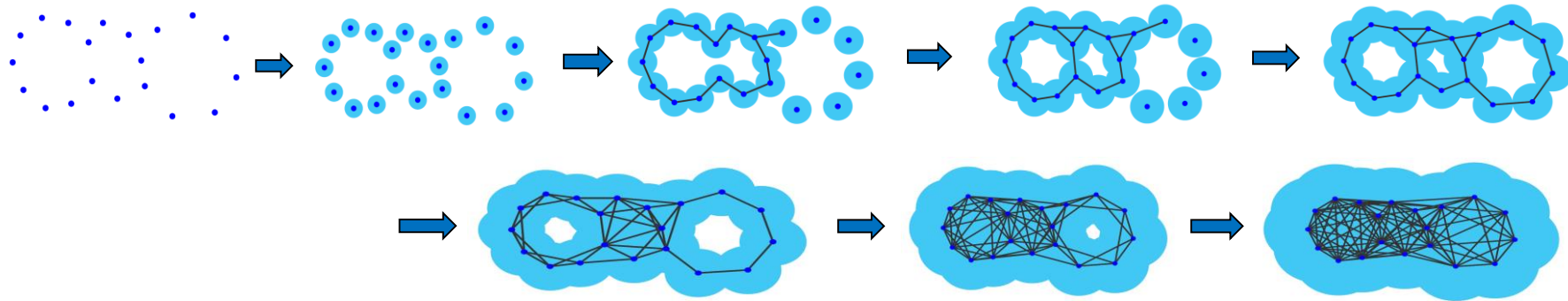


$$\beta_0 = 2, \beta_1 = 1$$



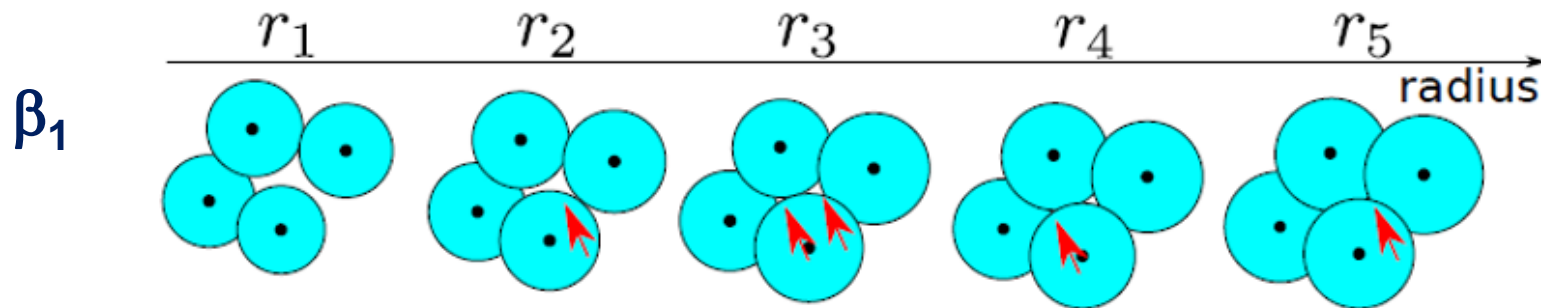
**фильтрация:** точка  $\rightarrow$  шар радиуса  $r \rightarrow$  увеличивая  $r$ , «следим» за:

- количеством связных компонент ( $\beta_0$ )
- количеством щелей ( $\beta_1$ )
- «временами жизни» связных компонент и щелей





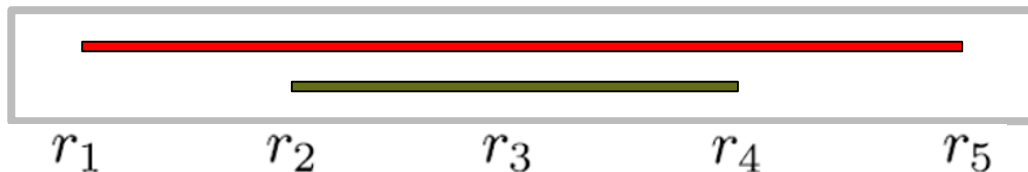
# Топологический Анализ Данных: пример (процессы «рождения и смерти»)



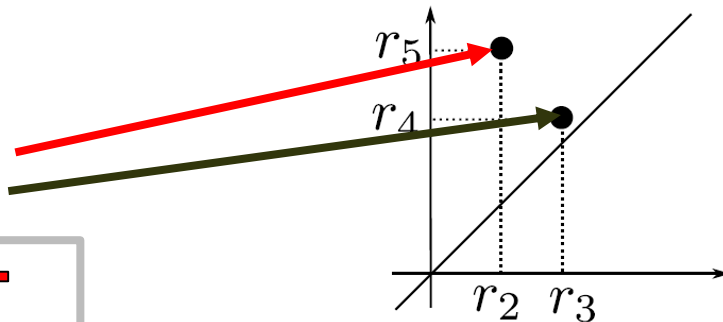
- В момент  $r_2$  «родилась» Щель 1
- В момент  $r_3$  «родилась» Щель 2
- В момент  $r_4$  «умерла» Щель 2
- В момент  $r_5$  «умерла» Щель 1



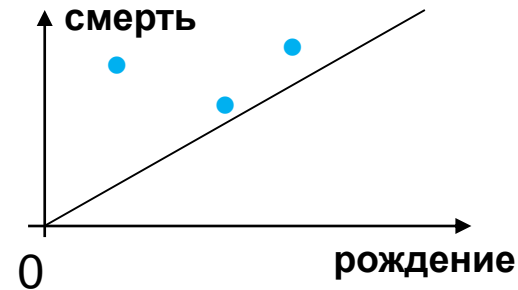
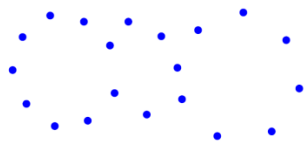
- **Щель 1:** родилась в момент  $r_2$  и умерла в момент  $r_5$
- **Щель 2:** родилась в момент  $r_3$  и умерла в момент  $r_4$



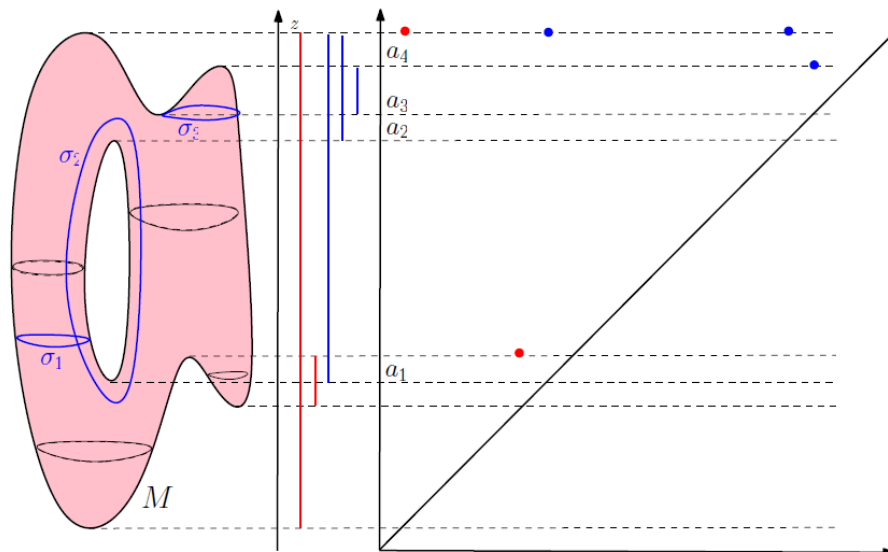
Персистентная диаграмма  
(persistence diagram)



Баркод (barcode)



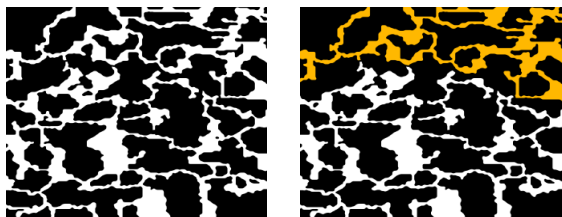
## Фильтрация по «линиям уровня»



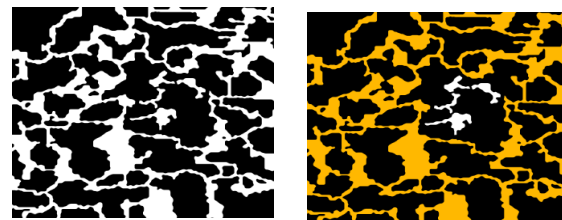
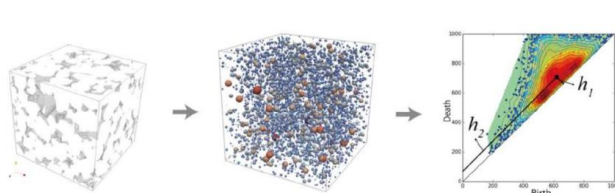
## Баркод и персистентная диаграмма:

- **Связные компоненты**
- **цели**

# ТАД: предсказания проницаемости пористых и гранулярных сред



Образец А (нулевая проницаемость)

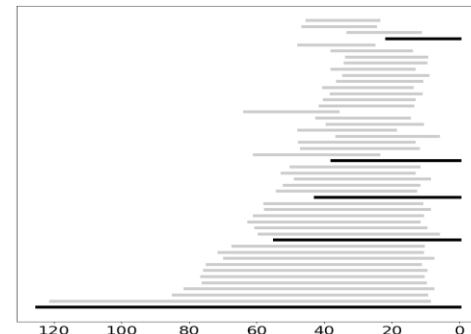
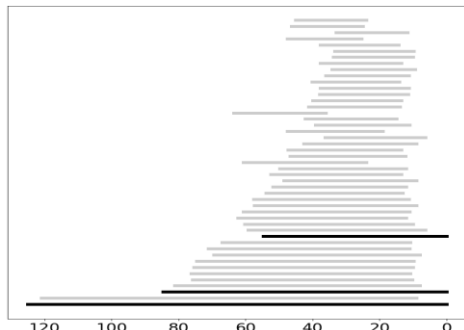


Образец Б (не нулевая проницаемость)

## Баркоды и персистентные диаграммы (размерности 0)

в зависимости от характерного размера

- Черные полосы: топологические характеристики, соответствующие компонентам связности множества пор
- Для непроницаемой/проницаемой пород персистентные диаграммы существенно различаются



Топологические признаки для предсказания проницаемости пористых и гранулярных сред

Topological Characteristics for prediction of permeability of porous and granular media

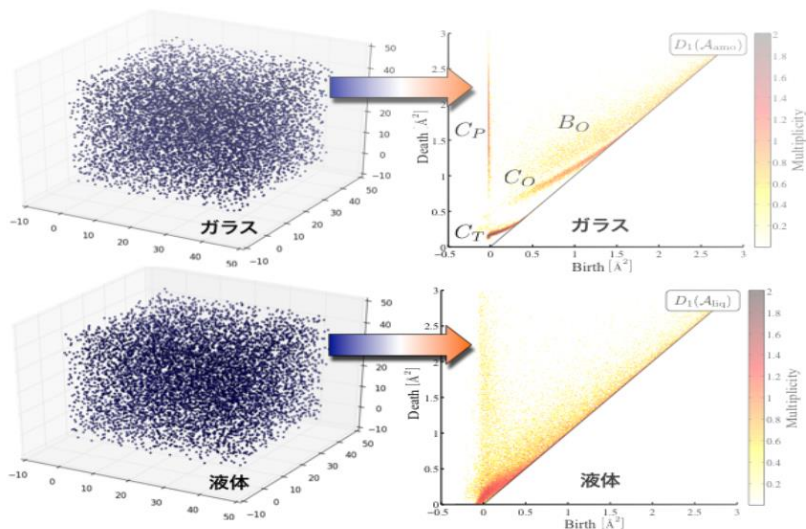
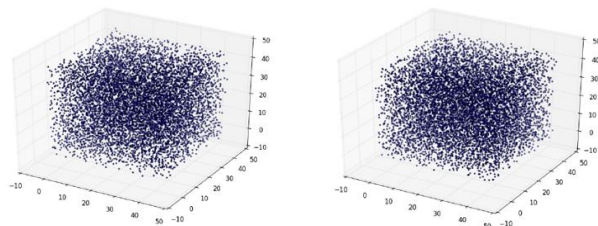
Бернштейн А., Бурнаев Е., Осипцов А., Баранников С., Качан О.  
Сколковский институт науки и технологий

Сравнение стандартных характеристик (функционалов Минковского) с современными подходами на основе топологического анализа данных

Бернштейн А., Бурнаев Е., Осипцов А., Баранников С., Качан О.  
Сколтех

# Топологический Анализ Данных: машинное зрение

## Атомные конфигурации различных химических соединений (аморфный и жидкий кремнеземы)



### Topological Data Analysis in Machine Vision

Alexander Bernstein

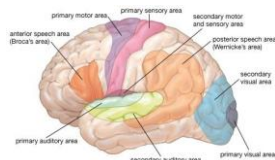
Joint work with E. Burnaev, M. Sharaev, E. Kondratyeva, O. Kachan

Skolkovo Institute of Science and Technology, Moscow, Russia

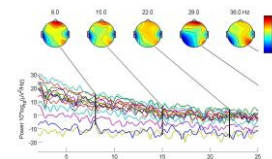
**SPIE.** DIGITAL LIBRARY

Amsterdam, The Netherlands, November 16, 2019

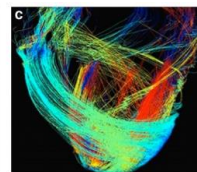
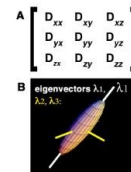
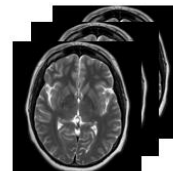
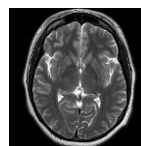
## Нейрональные зоны мозга



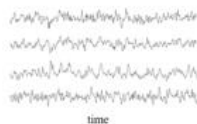
ЭЭГ



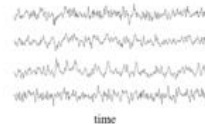
МРТ/фМРТ/диффузионно-тензорная МРТ



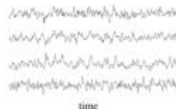
Регион 1



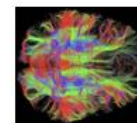
Регион 2



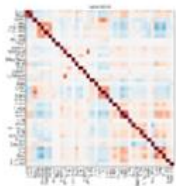
... Регион N



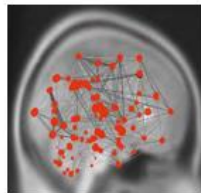
матрицы связей  
(коннективности)



Матрица  
коннективности



Граф  
коннективности



Здоровый



Больной

Первая международная конференция  
**Теоретическая физика  
и математика мозга:**  
междисциплинарные контакты

**Топологический анализ матриц коннективности  
в задачах медицинской диагностики**

А. Бернштейн<sup>1</sup>, В. Бухштабер<sup>1,2</sup>, Е. Бурнаев<sup>1</sup>, М. Шараев<sup>1</sup>,  
О. Качан<sup>1</sup>, Е. Стрельцова<sup>2</sup>, М. Поминова<sup>1</sup>



<sup>1</sup>Сколковский институт науки и технологий



<sup>2</sup>Московский Государственный Университет им. М.В. Ломоносова



4 декабря 2019 г.

Data Science

Building new models based on  
biomedical training data and preliminary  
constructed predictive models.

Student: *Nikolay Skuratov*  
Research Advisor: *Alexander Bernstein*

Skoltech

June, 2020



**Диагностика депрессии по ЭЭГ в состоянии покоя  
(19 каналов) (Resting-state)**

<b>Точность:</b>	
по 8 стандартным характеристикам	<b>0,79</b>
<b>+ 2 топологических характеристики</b> (Entropy Complex on Envelopes, 0,64)	<b>0.85</b>
<b>+ еще 2 топологических характеристики</b> (Persistent Diagram Bins, 0.63)	<b>0.88</b>

---

# СПАСИБО ЗА ВНИМАНИЕ