



# Прикладная статистика и анализ данных

Съезд V

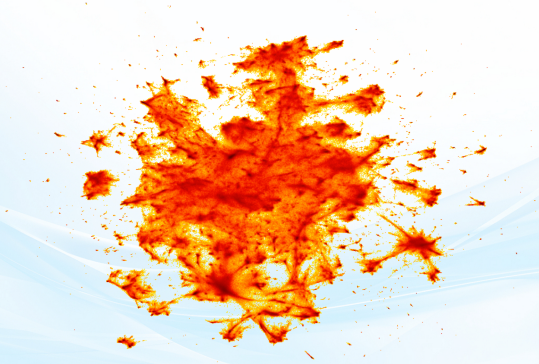
# Метод UMAP (2018)

Uniform Approximation and Projection —

метод, выполняющий нелинейное снижение размерности.

Создан Лилендом Макиннесом + его коллегами из Таттского инст.

**Цель:** было получить метод, похожий на t-SNE,  
но с более сильным математическим обоснованием.



Визуализация 3 млн. слов из GoogleNews dataset



## Метод UMAP (2018)

Пусть  $X = (X_1, \dots, X_n)$  — выборка в пространстве  $\mathcal{X}$ .

$\rho(x_1, x_2)$  — метрика в  $\mathcal{X}$

Определим **случайный ориентир. граф** на  $m$ -ве вершин  $X$ .

Считаем, что каждое ребро появляются независимо от других.

Рассмотрим вершину  $X_i$ .

Пусть  $X_{i_1}, \dots, X_{i_k}$  —  $k$  ближайших соседей объекта  $X_i$  в выборке  $X$ .

$r_i = \min_s \rho(X_i, X_{i_s})$  — расстояние до ближайшего соседа.

Вероятность ребра из  $X_i$  в  $X_{i_s}$  [для остальных вер-ть = 0]:

$$P(X_i \rightarrow X_{i_s}) = \exp \left( -\frac{\rho(X_i, X_{i_s}) - r_i}{\sigma_i} \right) \in [0, 1],$$

где  $\sigma_i$  подбирается как решение уравнения  $\sum_{s=1}^k P(X_i \rightarrow X_{i_s}) = \log_2 k$ ,  
т.е.  $\sigma_i$  играет роль нормировки вероятностей ребер.



## Метод UMAP (2018)

На основе ориентированного графа построим **неориентированный**.

$X$  — множество вершин

Вероятность ребра между  $X_i$  и  $X_j$ :

$$\begin{aligned} P(X_i - X_j) &= P(\{X_i \rightarrow X_j\} \cup \{X_j \rightarrow X_i\}) = \\ &= P(X_i \rightarrow X_j) + P(X_j \rightarrow X_i) - P(X_i \rightarrow X_j) P(X_j \rightarrow X_i). \end{aligned}$$

Пусть  $Y = (Y_1, \dots, Y_n)$  — **вложение**  $X$  в маломерное пр-во  $\mathbb{R}^d$ .

Величины  $Y_1, \dots, Y_n$  неизвестны, их нужно подобрать.

На  $Y$  зададим случайный неор. граф с вероятностями

$P(Y_i - Y_j) = (1 + a\|Y_i - Y_j\|_2^{2b})^{-1}$ , где  $a$  и  $b$  — гиперпараметры.

**Минимизируем** дивергенцию Кульбака-Лейблера:  $KL(P_X, P_Y) =$

$$= \sum_{i,j} \left[ P(X_i - X_j) \log \frac{P(X_i - X_j)}{P(Y_i - Y_j)} + (1 - P(X_i - X_j)) \log \frac{1 - P(X_i - X_j)}{1 - P(Y_i - Y_j)} \right]$$



# Компоновка графа в пространстве низкой размерности

Для расположения точек в пространстве низкой размерности используется компоновка графа **в теории**.

- ▶ Сила притяжения вдоль ребра из  $Y_i$  в  $Y_j$

$$\frac{-2ab\|Y_i - Y_j\|_2^{2(b-1)}}{1 + \|Y_i - Y_j\|_2^2} P(X_i \rightarrow X_j)(Y_i - Y_j),$$

где  $a$  и  $b$  — гиперпараметры.

- ▶ Сила отталкивания между вершинами  $Y_i$  и  $Y_j$

$$\frac{b}{(\varepsilon + \|Y_i - Y_j\|_2^2)(1 + \|Y_i - Y_j\|_2^2)} (1 - P(X_i \rightarrow X_j))(Y_i - Y_j),$$

где  $\varepsilon = 0.001$ .

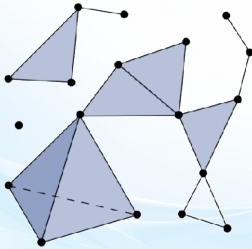
- ▶ Сходимость обеспечивается за счет медленного уменьшения сил притяжения и отталкивания.
- ▶  $P(Y_i - Y_j)$  — аппроксимация этих сил.



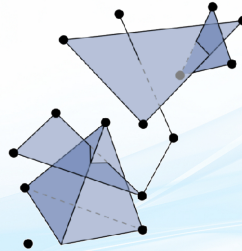
# Рассмотрим теоретическое обоснование UMAP

**Симплициальный комплекс** — топологическое пространство с заданной на нем триангуляцией, т.е., склеенное из топологических симплексов по определенным правилам.

Simplicial complex



Collection of Simplices





# Немного топологии

**Покрытие** топологического пространства — семейство множеств, таких, что их объединение содержит заданное множество.

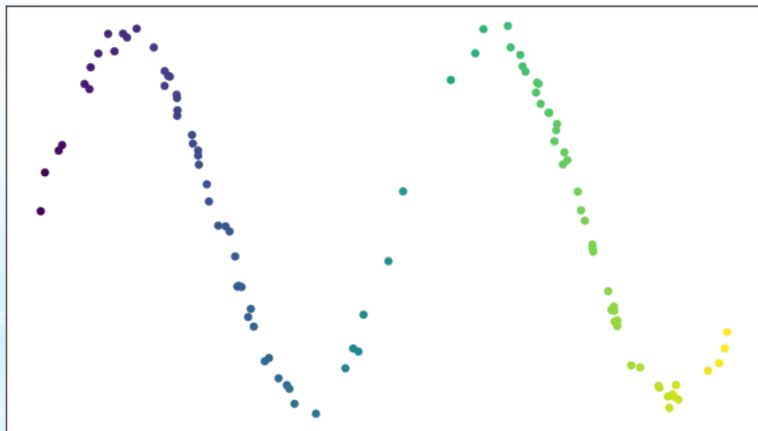
Пусть  $\{W_\alpha\}$  — конечное покрытие топологического пространства  $X$ .

**Нерв покрытия**  $\{W_\alpha\}$  — это абстрактный симплициальный комплекс  $N$ , множество вершин которого отождествлено

с множеством индексов покрытия, при этом  $N$  содержит симплекс

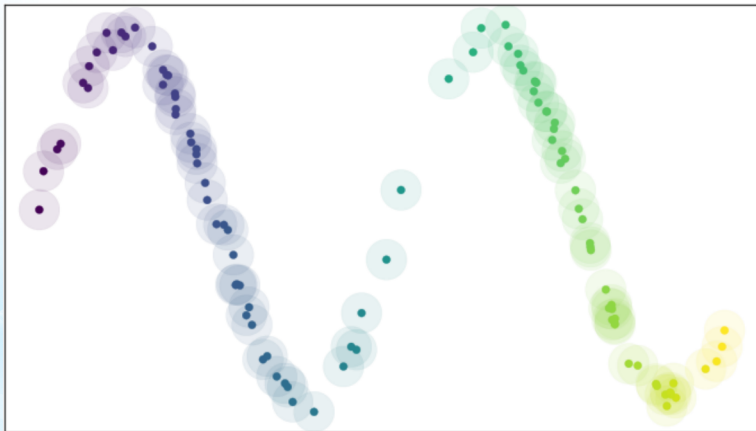
с вершинами  $\alpha_1, \alpha_2, \dots, \alpha_n$  тогда и только тогда, когда  $\bigcap_{i=1}^n W_{\alpha_i} \neq \emptyset$ .

# Топологическое множество

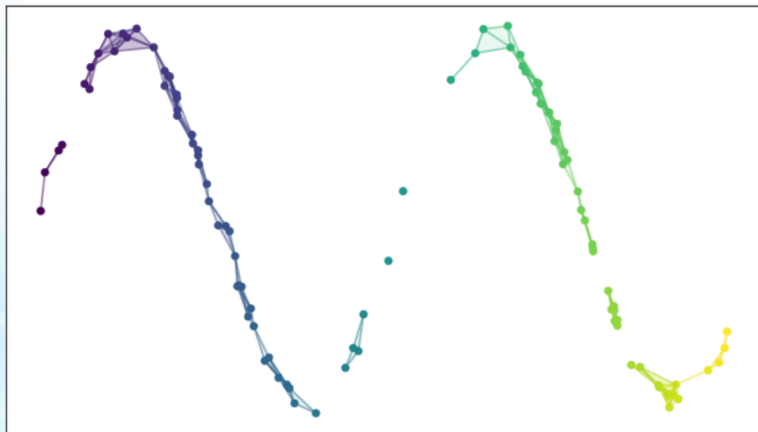




# Покрывтие топологического множества



# Нерв покрытия





## Теорема о нерве

Если пространство  $X$  триангулируемо и  $\{W_\alpha\}$  — конечное покрытие замкнутыми множествами, причём все непустые пересечения стягиваемы, то нерв покрытия гомотопически эквивалентен  $X$ .

Гомотопия из  $X$  в  $Y$  — непрерывное отображение  $f : [0, 1] \times X \rightarrow Y$ .

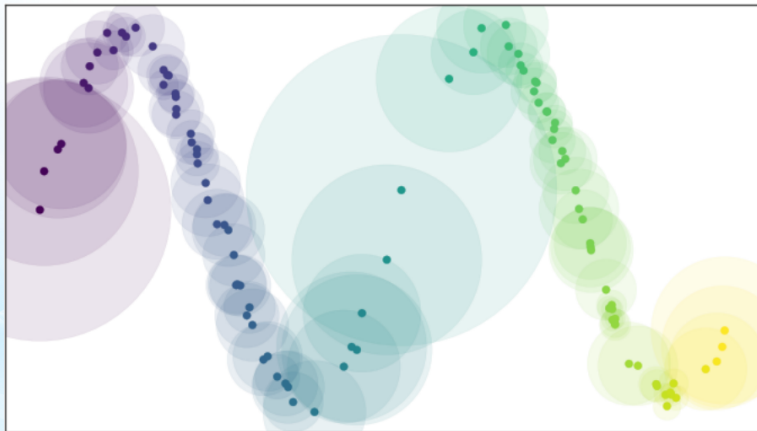
Если данные равномерно распределены по многообразию, то покрытие будет “хорошим”.

Когда данные так хорошо себя ведут?

*Предположение:* данные равномерно распределены на многообразии!

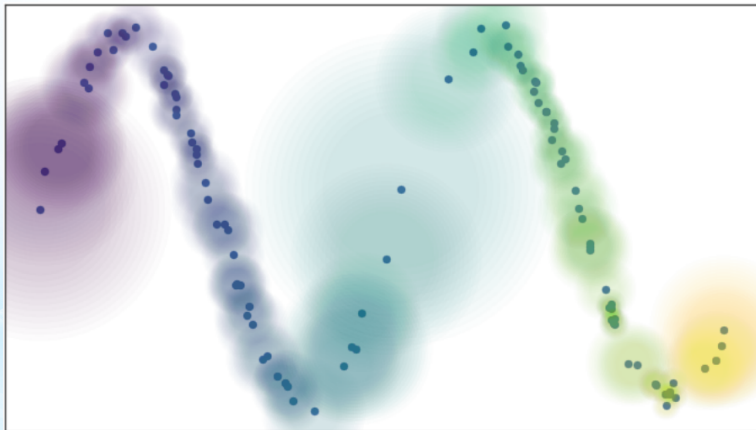
Определим Риманову метрику на многообразии, чтобы сделать это предположение истинным.

## Покрытие с разными радиусами

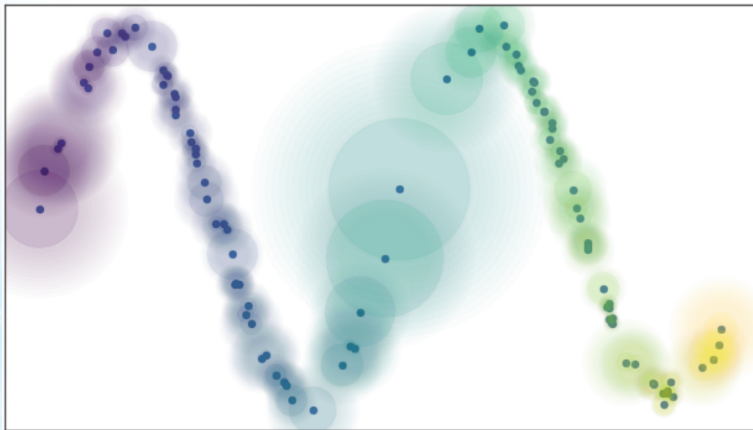




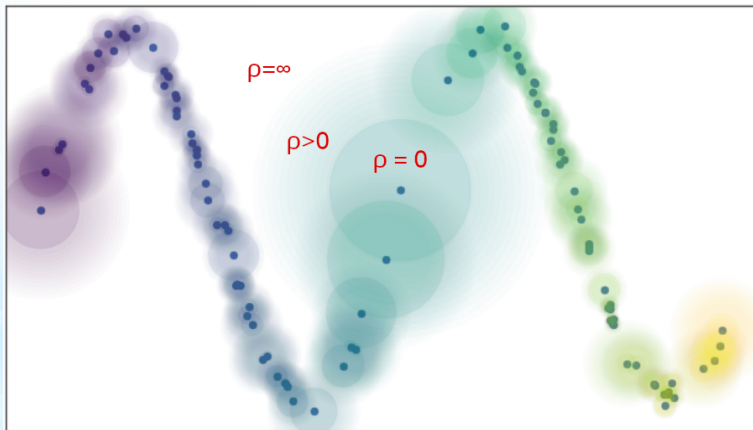
Но почему радиус фиксирован? Берем случайный:



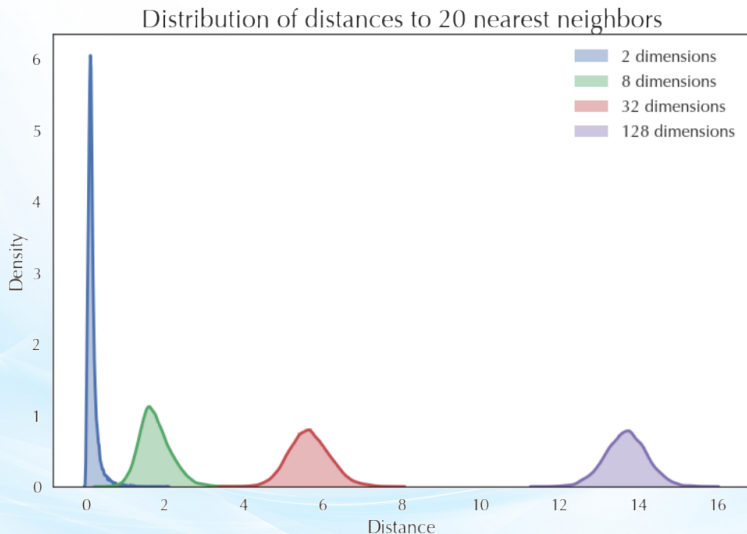
Предположение: многообразие локально связно



Предположение: многообразие локально связно



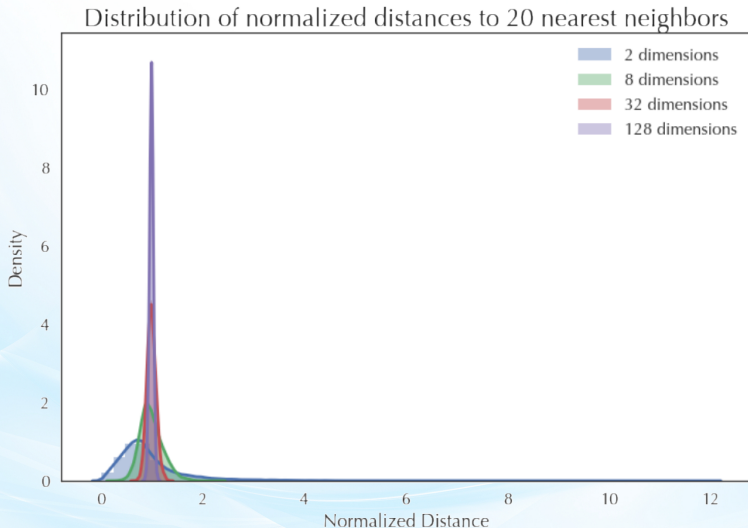
# Распределение расстояний







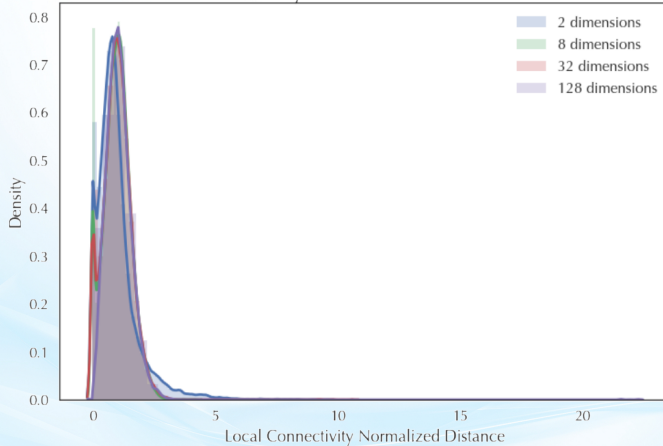
# Распределение нормированных расстояний





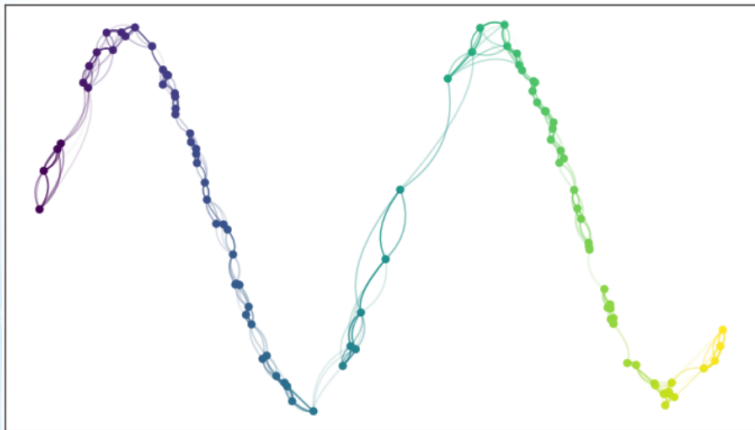
# Распр. норм. расстояний на основе лок. связности

Distribution of local connectivity normalized distances to 20 nearest neighbors

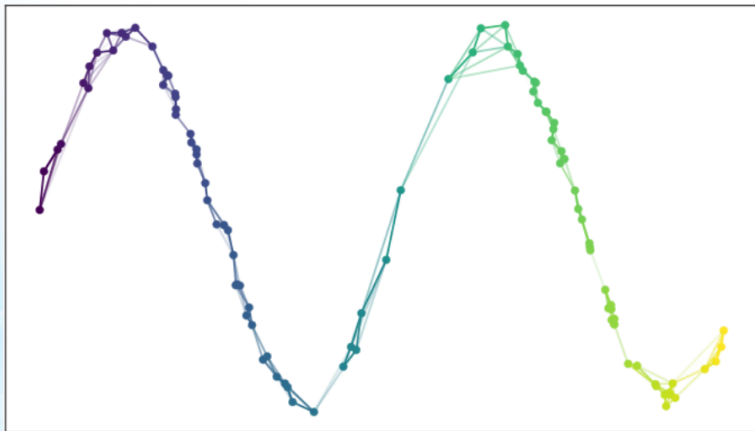




Но наша локальная метрика не симметрична!



Сделаем ее симметричной



# Далее

## Топология, теория представлений, нечеткие множества...

**Definition 7.** Define the functor  $\text{FinReal} : \mathbf{Fin-sFuzz} \rightarrow \mathbf{FinEPMet}$  by setting

$$\text{FinReal}(\Delta_{<a}^n) \triangleq (\{x_1, x_2, \dots, x_n\}, d_a),$$

where

$$d_a(x_i, x_j) = \begin{cases} -\log(a) & \text{if } i \neq j, \\ 0 & \text{otherwise} \end{cases}.$$

and then defining

$$\text{FinReal}(X) \triangleq \text{colim}_{\Delta_{<a}^n \rightarrow X} \text{FinReal}(\Delta_{<a}^n).$$

Similar to Spivak's construction, the action of  $\text{FinReal}$  on a map  $\Delta_{<a}^n \rightarrow \Delta_{<b}^m$  where  $a \leq b$  defined by  $\sigma : \Delta^n \rightarrow \Delta^m$ , is given by

$$(\{x_1, x_2, \dots, x_n\}, d_a) \mapsto (\{x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(n)}\}, d_b),$$

which is a non-expansive map since  $a \leq b$  implies  $d_a \geq d_b$ .

Since  $\text{FinReal}$  preserves colimits it admits a right adjoint, the fuzzy singular set functor  $\text{FinSing}$ . We can then define the (finite) fuzzy singular set functor in terms of the action of its image on  $\Delta \times I$ , analogously to  $\text{Sing}$ .

**Definition 8.** Define the functor  $\text{FinSing} : \mathbf{FinEPMet} \rightarrow \mathbf{Fin-sFuzz}$  by

$$\text{FinSing}(Y) : ([n], [0, a]) \mapsto \text{hom}_{\mathbf{FinEPMet}}(\text{FinReal}(\Delta_{<a}^n), Y).$$

We then have the following theorem.

**Theorem 1.** The functors  $\text{FinReal} : \mathbf{Fin-sFuzz} \rightarrow \mathbf{FinEPMet}$  and  $\text{FinSing} : \mathbf{FinEPMet} \rightarrow \mathbf{Fin-sFuzz}$  form an adjunction with  $\text{FinReal}$  the left adjoint and  $\text{FinSing}$  the right adjoint.

*Proof.* The adjunction is evident by construction, but can be made more explicit as follows. Define a functor  $F : \Delta \times I \rightarrow \mathbf{FinEPMet}$  by

$$F([n], [0, a]) = (\{x_1, x_2, \dots, x_n\}, d_a),$$

where

$$d_a(x_i, x_j) = \begin{cases} -\log(a) & \text{if } i \neq j, \\ 0 & \text{otherwise} \end{cases}.$$

Now  $\text{FinSing}$  can be defined in terms of  $F$  as

$$\text{FinSing}(Y) : ([n], [0, a]) \mapsto \text{hom}_{\mathbf{FinEPMet}}(F([n], [0, a]), Y).$$

where the face maps  $d_i$  are given by pre-composition with  $Fd^i$ , and similarly for degeneracy maps, at any given value of  $a$ . Furthermore post-composition with  $F$  level-wise for each  $a$  defines maps of fuzzy simplicial sets making  $\text{FinSing}$  a functor.

We now construct  $\text{FinReal}$  as the left Kan extension of  $F$  along the Yoneda embedding:

$$\begin{array}{ccc} & \mathbf{Fin-sFuzz} & \\ \uparrow \eta & \nearrow \text{FinReal} & \\ \Delta \times I & \xrightarrow{F} & \mathbf{FinEPMet} \end{array}$$

Explicitly this results in a definition of  $\text{FinReal}$  at a fuzzy simplicial set  $X$  as a colimit:

$$\text{FinReal}(X) = \text{colim}_{y([m], [0, a]) \rightarrow X} F([m]).$$

Further, it follows from the Yoneda lemma that  $\text{FinReal}(\Delta_{<a}^n) \cong F([n], [0, a])$ , and hence this definition as a left Kan extension agrees with Definition 7, and the definition of  $\text{FinSing}$  above agrees with that of Definition 8. To see that  $\text{FinReal}$  and  $\text{FinSing}$  are adjoint we note that



**ВСЁ!**