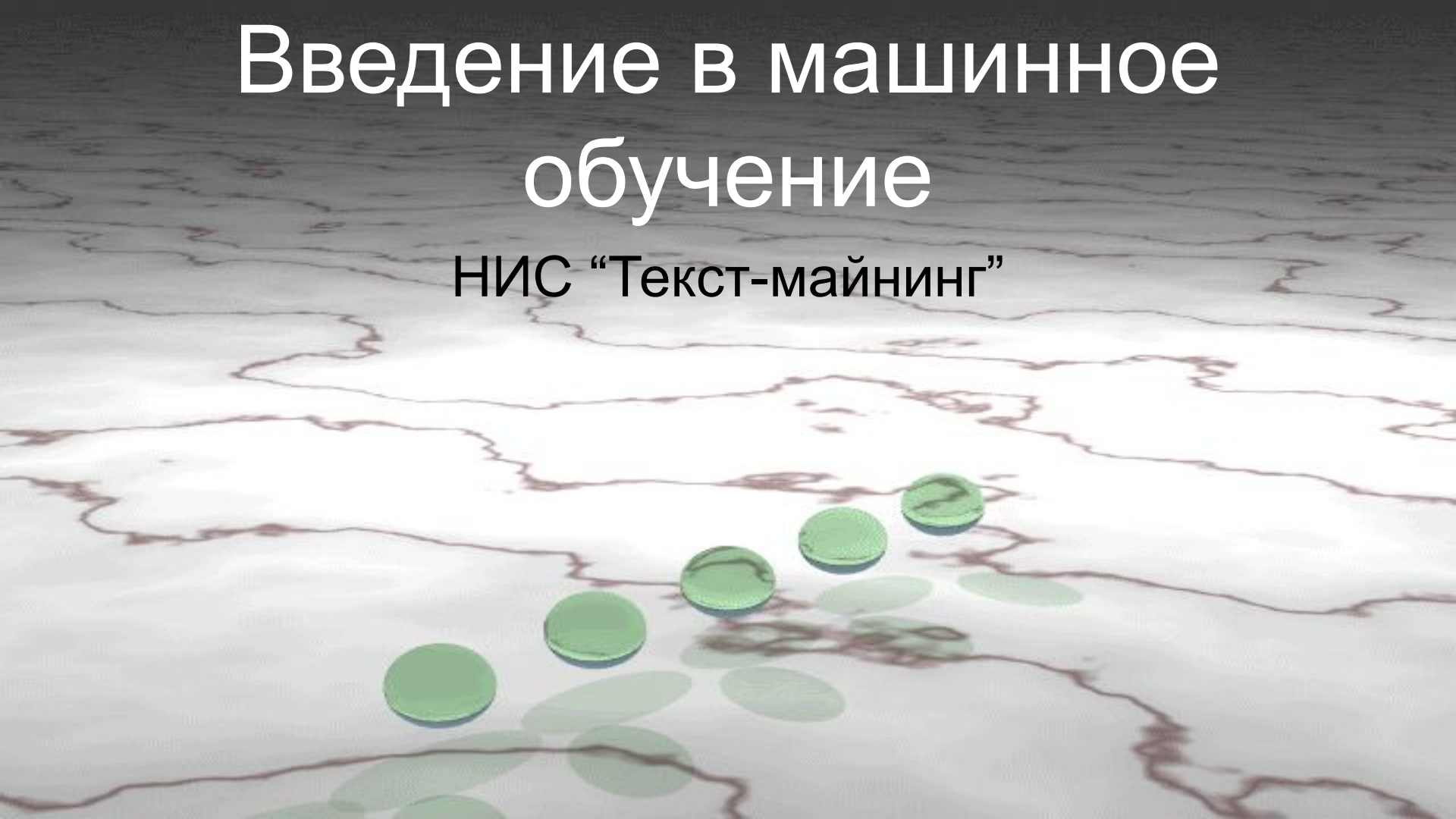
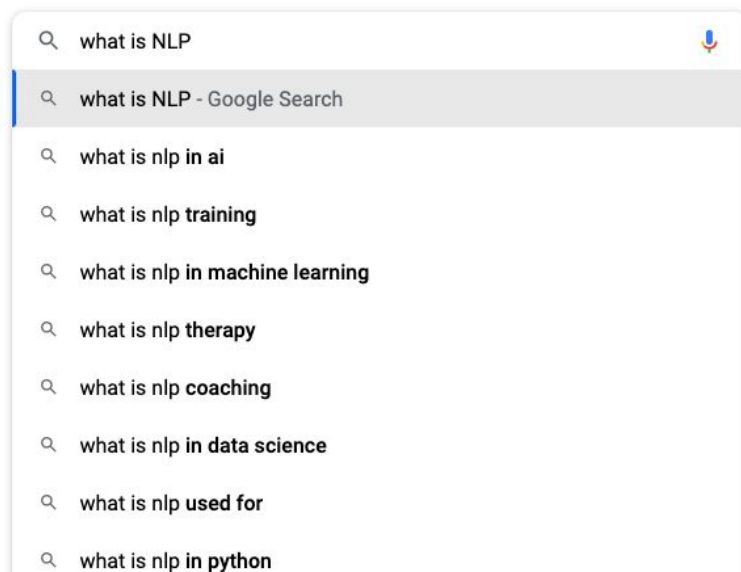


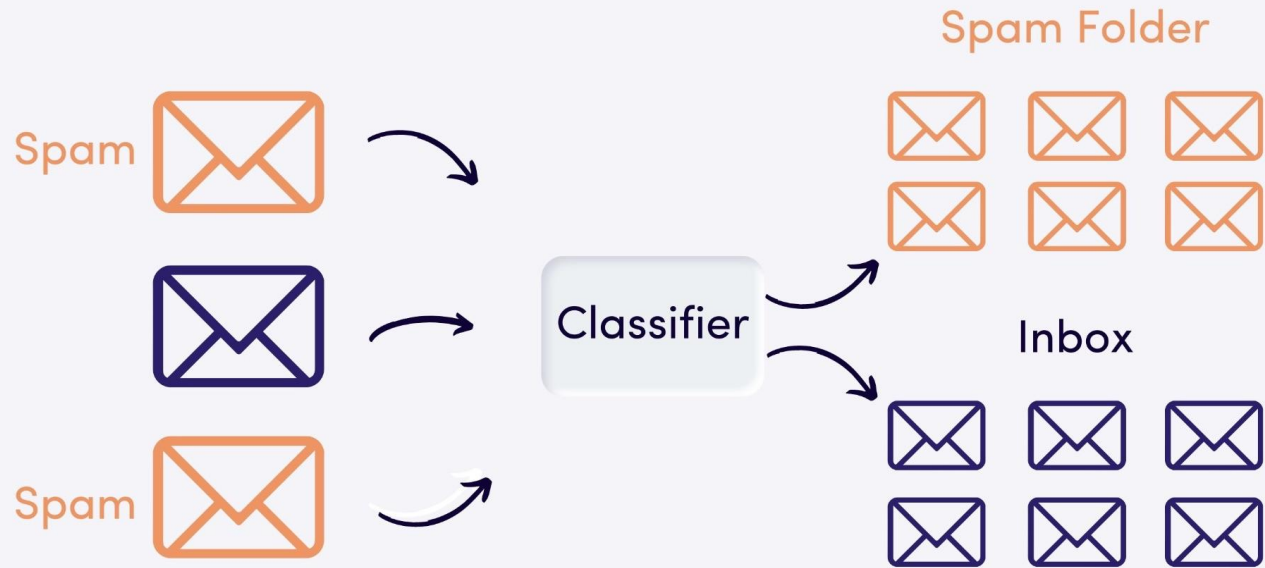
Введение в машинное обучение

НИС “Текст-майнинг”



Машинное обучение и тексты







COCA-COLA

Coca-Cola Soda

Overall

KPIs

Pros & cons

Relevance

Predictive

Wonder

MORE

Notifications

< Back

Wonder

Wonder

Insights powered by Wonderflow AI. Ask anything you would like to know, we have the answers.



Performance

How is this product performing?



What are the main areas of improvements?



What are the characteristics that customers love?



Advanced

What macro-topics are more relevant for consumers?



What problems, if fixed, generate the largest increase in star rating?



How is this product performing against its competitors?



Prescriptive

What problems should we address first?



What should we do to make this product the best in its category?



What topics should we advertise to generate maximum return from consumers?



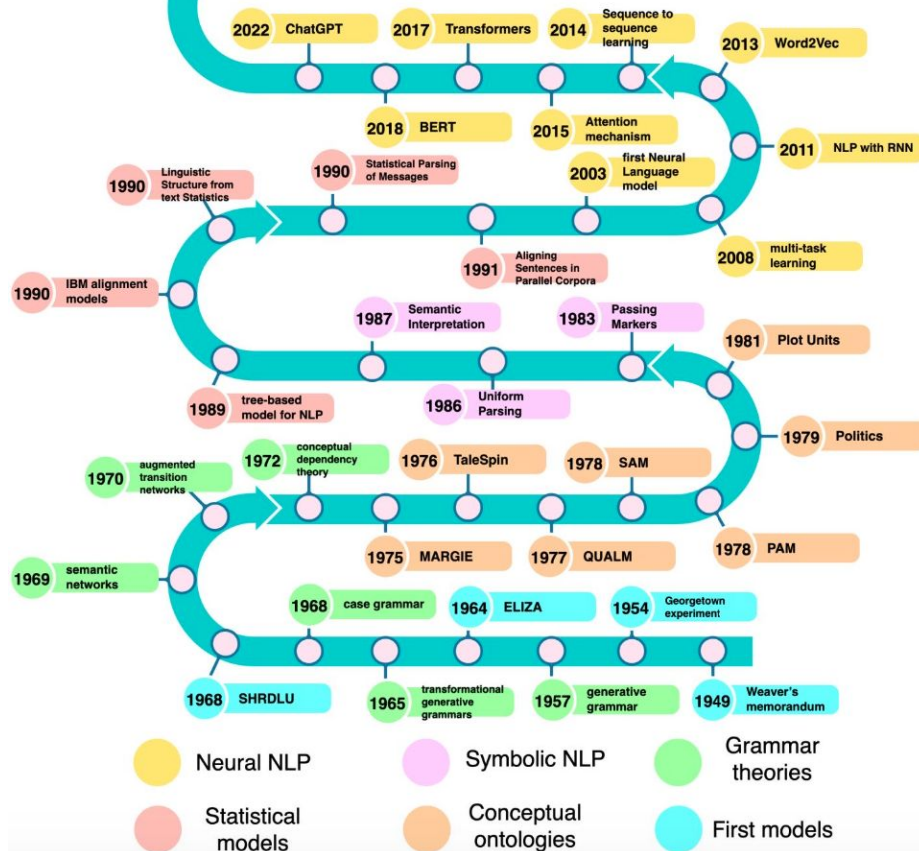
Ask me anything



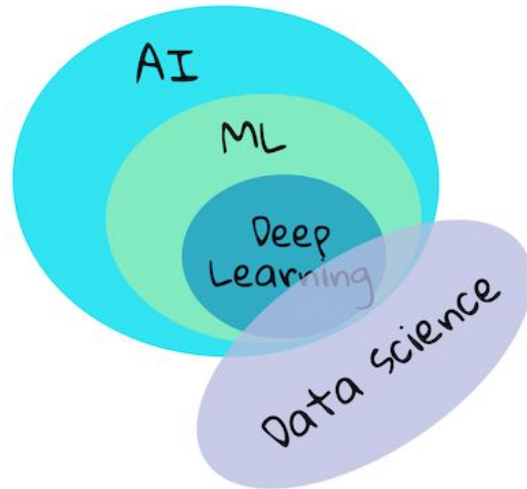
History of NLP

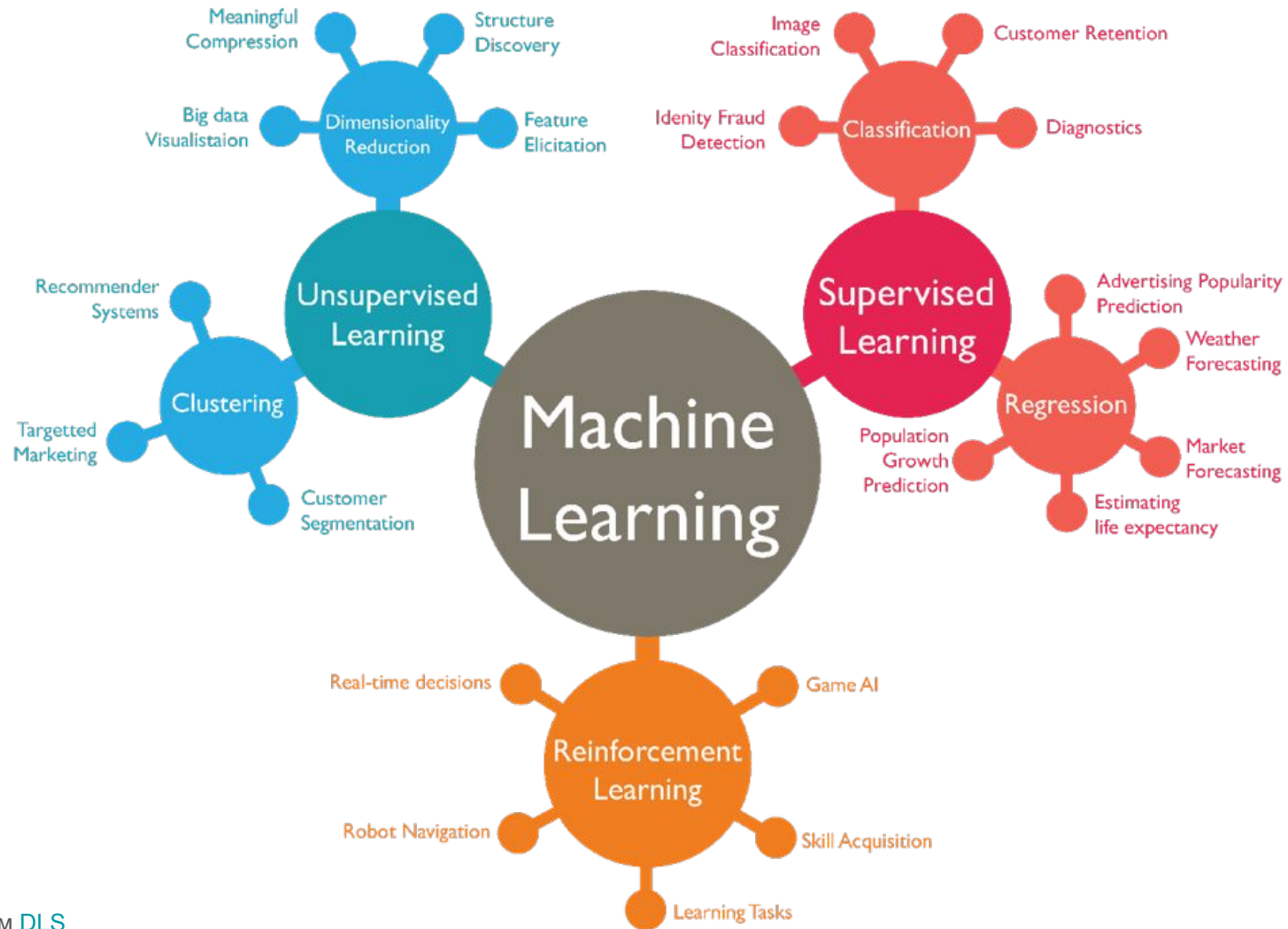
TheAiEdge.io

Natural Language Processing

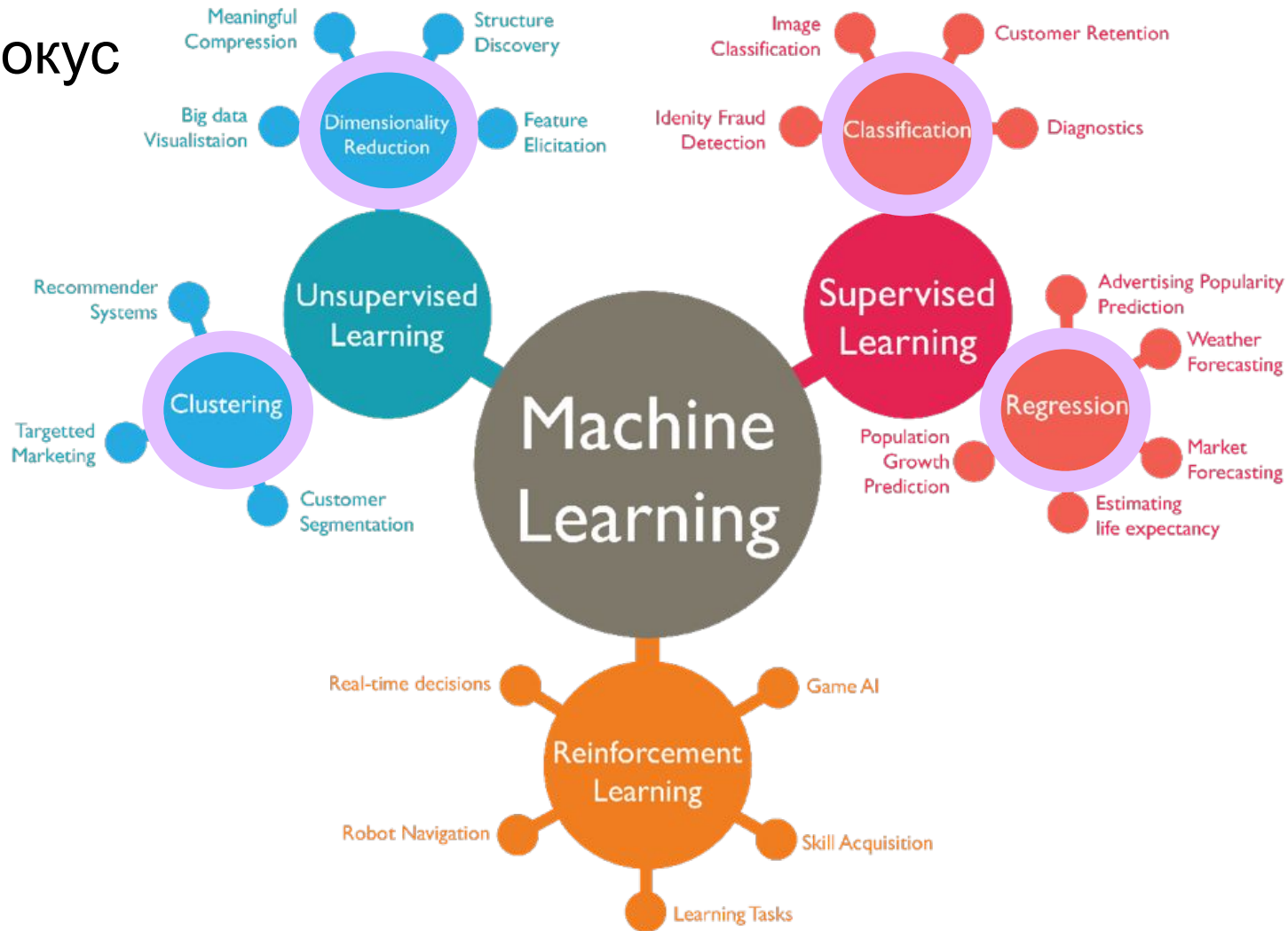


Области машинного обучения





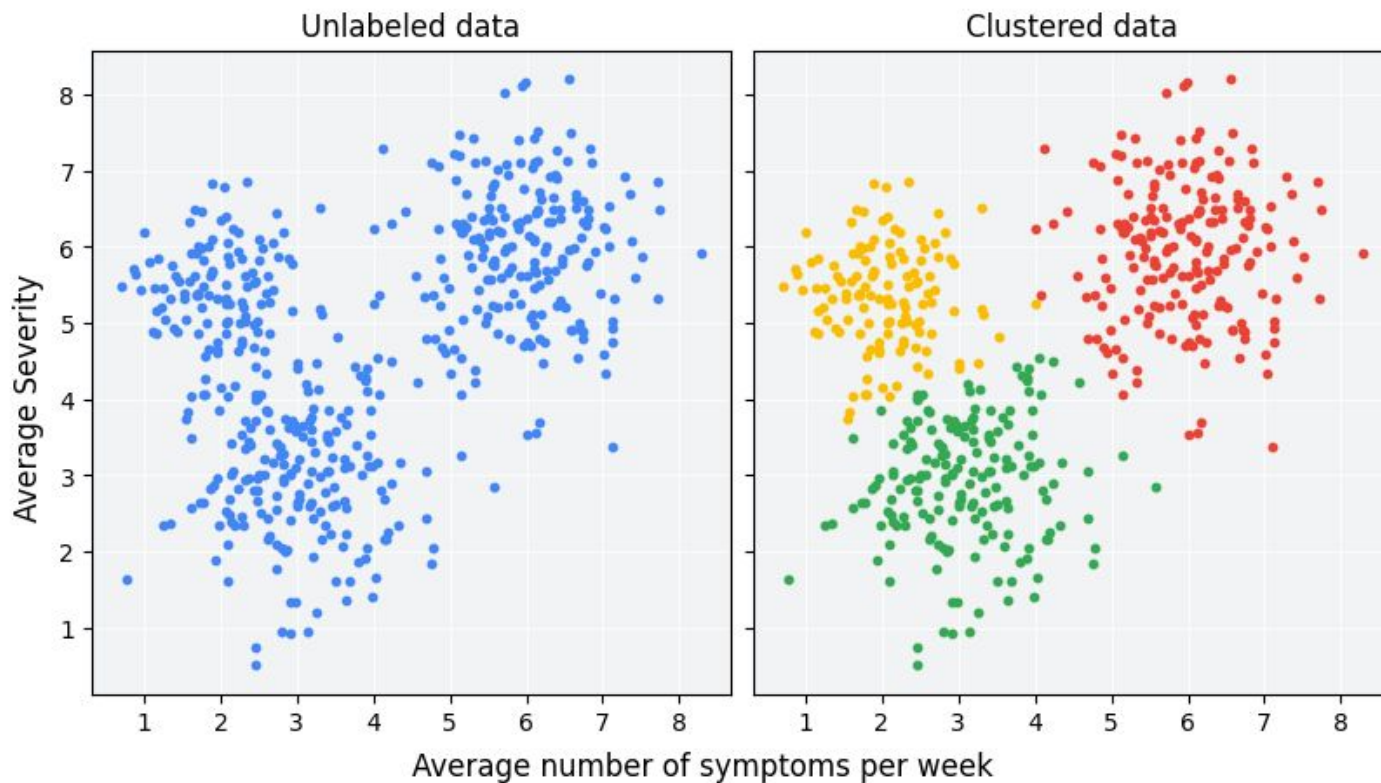
Наш фокус



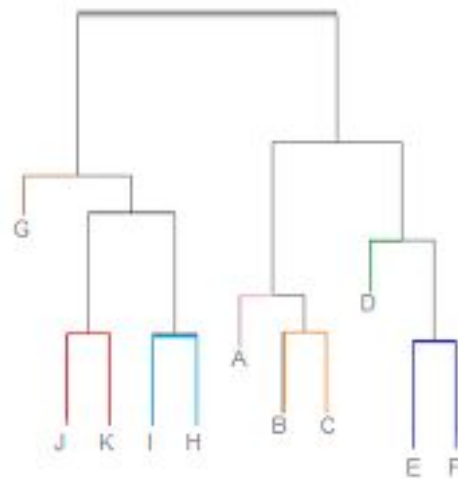
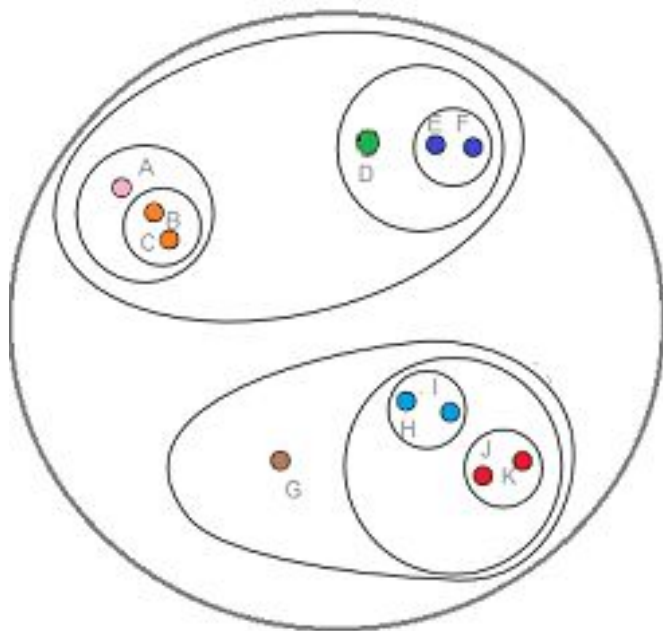
Классическое Обучение



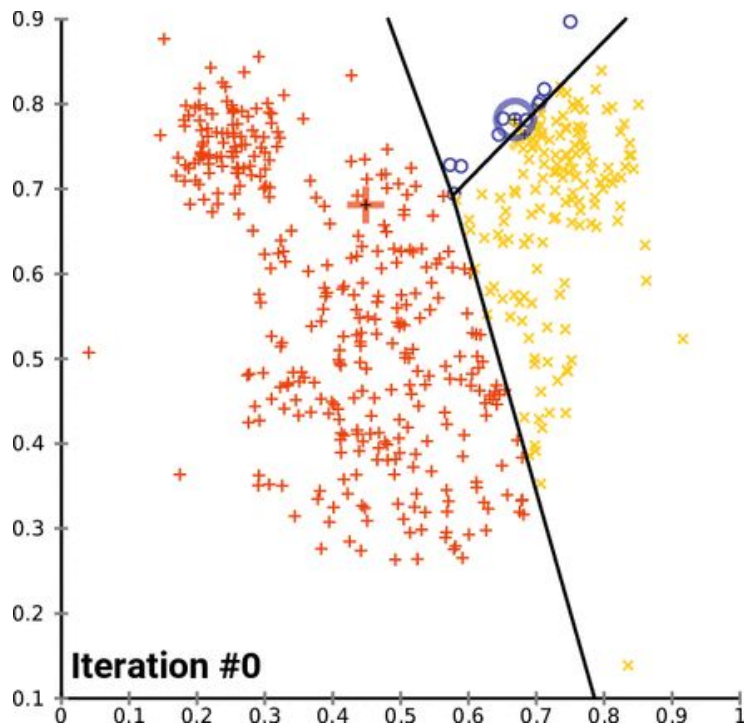
Кластеризация



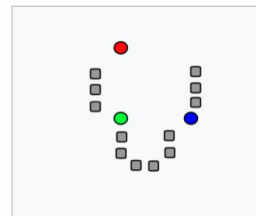
Алгоритмы кластеризации



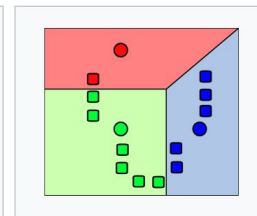
K-means



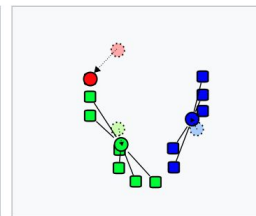
Demonstration of the standard algorithm



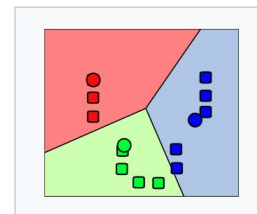
1. k initial "means" (in this case $k=3$) are randomly generated within the data domain (shown in color).



2. k clusters are created by associating every observation with the nearest mean. The partitions here represent the [Voronoi diagram](#) generated by the means.

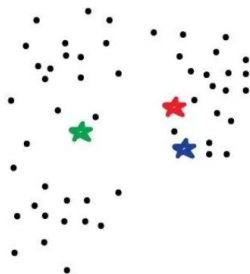


3. The [centroid](#) of each of the k clusters becomes the new mean.

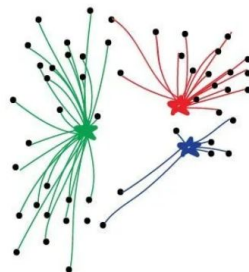


4. Steps 2 and 3 are repeated until convergence has been reached.

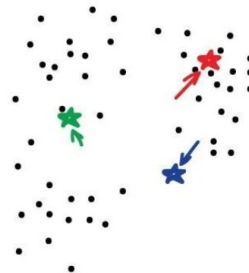
Ставим три ларька с шаурмой оптимальным образом (иллюстрируя метод К-средних)



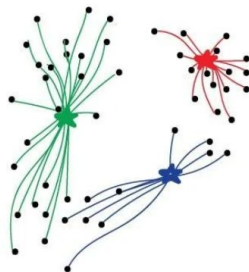
1. Ставим ларьки с шаурмой
в случайных местах



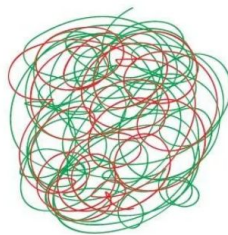
2. Смотрим в какой
кому ближе идти



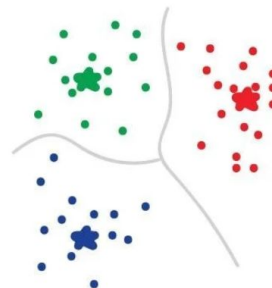
3. Двигаем ларьки ближе
к центрам их популярности



4. Снова смотрим и двигаем



5. Повторяем много раз



6. Готово, вы великолепны!

Примеры задач кластеризации

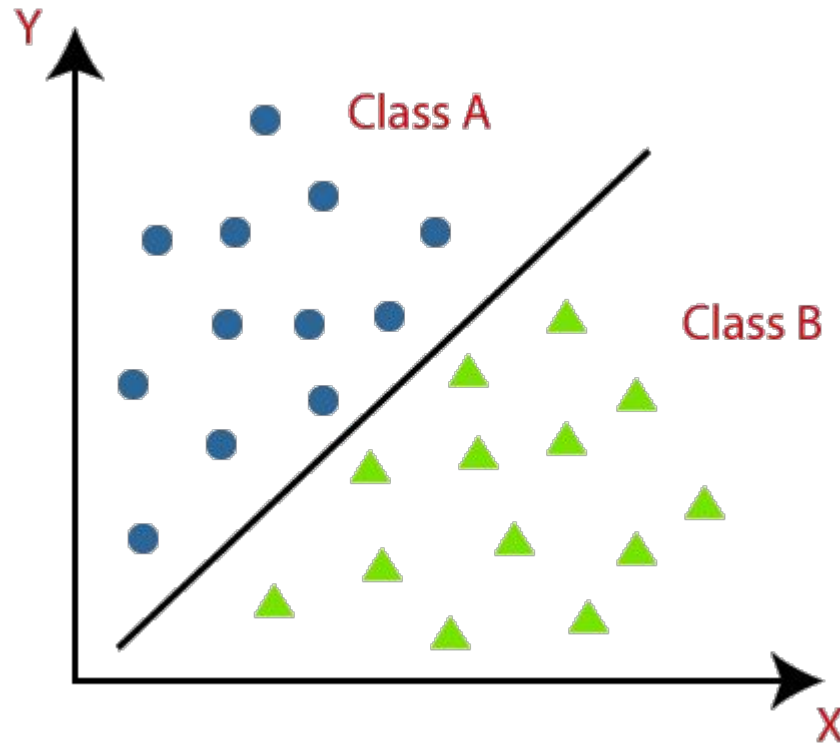
- NLP:

- Анализ отзывов
- Тематическое моделирование
- Определение тональности и настроений

- Другие области:

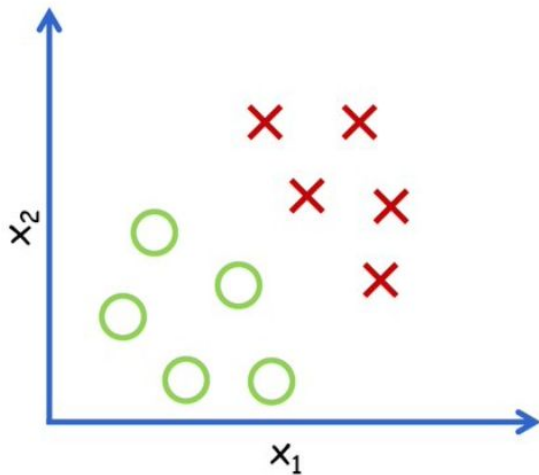
- Кластеризация генов или белков на основе их функций или экспрессии в геномных данных для выявления групп сходных биологических функций
- Группировка клиентов по потребительским привычкам
- Кластеризация пользователей по общим интересам

Классификация - заранее знаем классы!

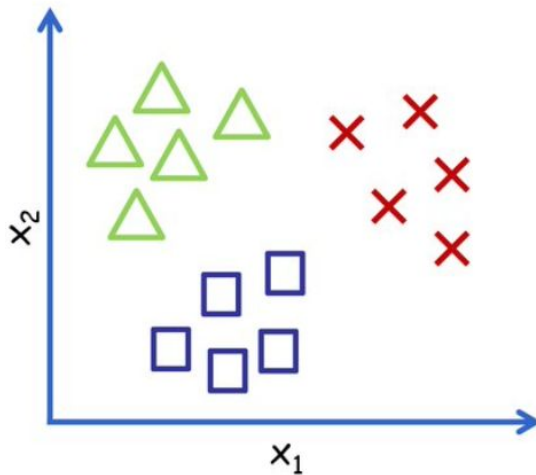


Бинарная и многоклассовая

Бинарная классификация



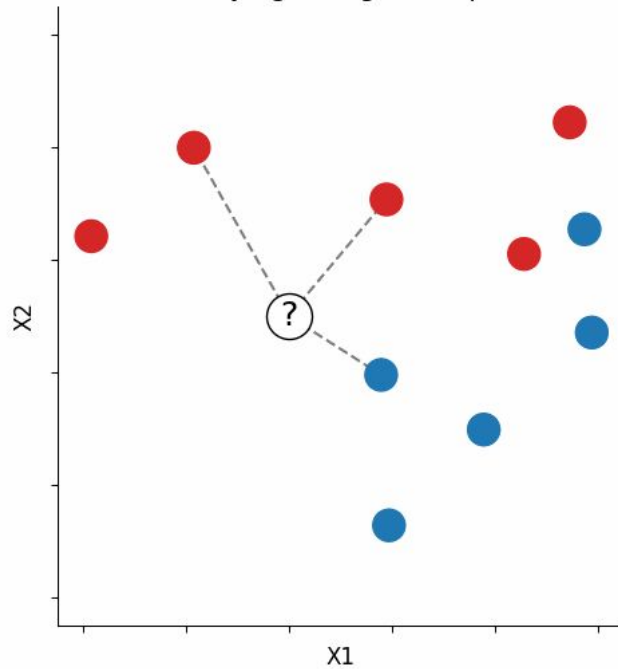
Многоклассовая классификация



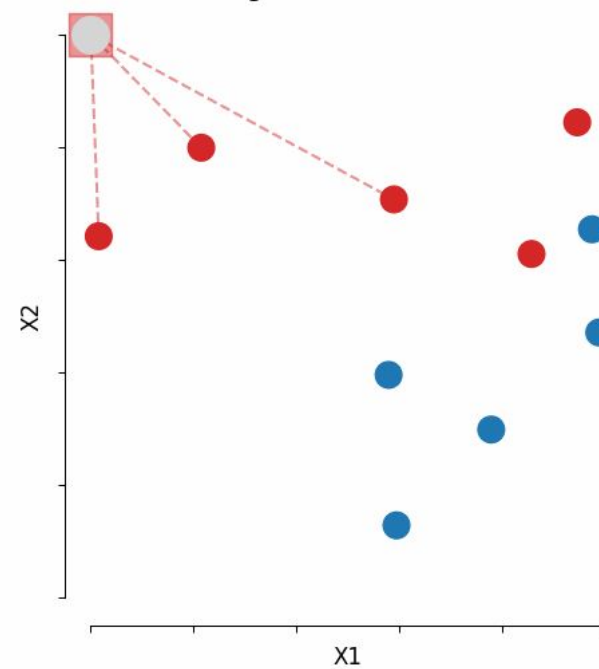
KNN

kNN classifier | $k = 3$

Classifying a single test point



Drawing the decision surface



Примеры задач классификации

- NLP:

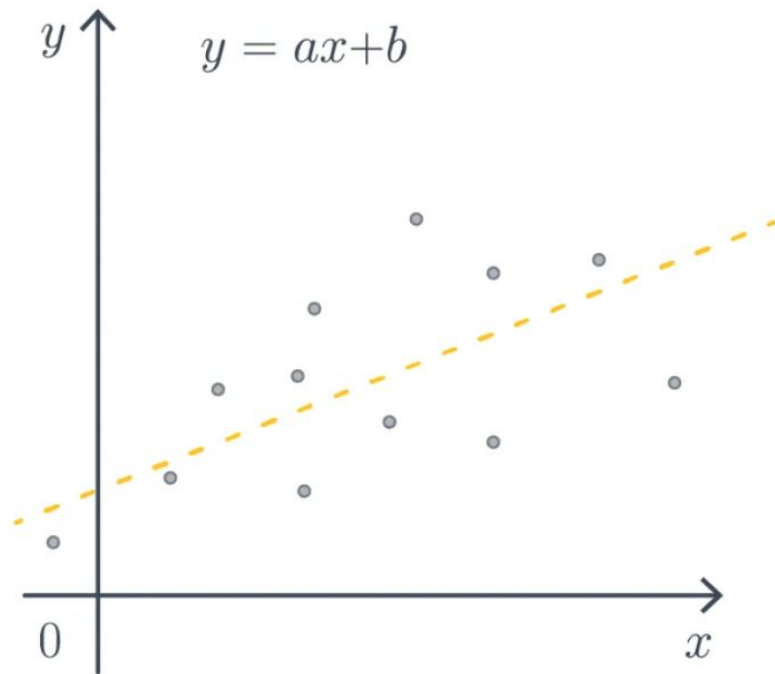
- Выявление фейковых и реальных новостей (бинарная классификация)
- Токсичность
- Спам

- Другие области:

- Кредиты
- Наличие заболевания

Немного про регрессию

Регрессия



x — (единственный)
признак

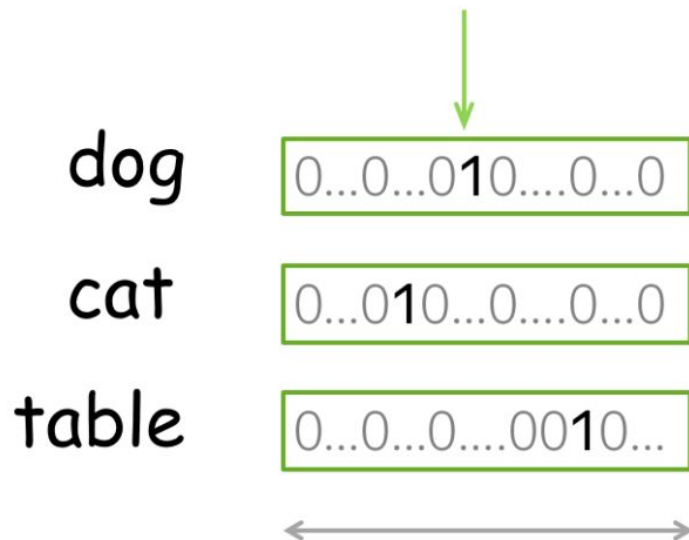
y — таргет

Как работать с данными?

- Выполняем предобработку, которой мы научились + **векторизация**
- EDA (exploratory data analysis)
 - На этом этапе **важно**:
 - Обработать NaN (заменить средним / медианой / dropna)
 - Обработать выбросы (слишком высокие или слишком низкие значения)
 - Проверить дубликаты (как правило их удаляют)
 - Постройте визуализации, которые помогут лучше понять ваши данные
- Разбиваем данные на train и test
- Обучение
- Валидация
- Наслаждаемся жизнью модель делает все за нас))) (если бы...)

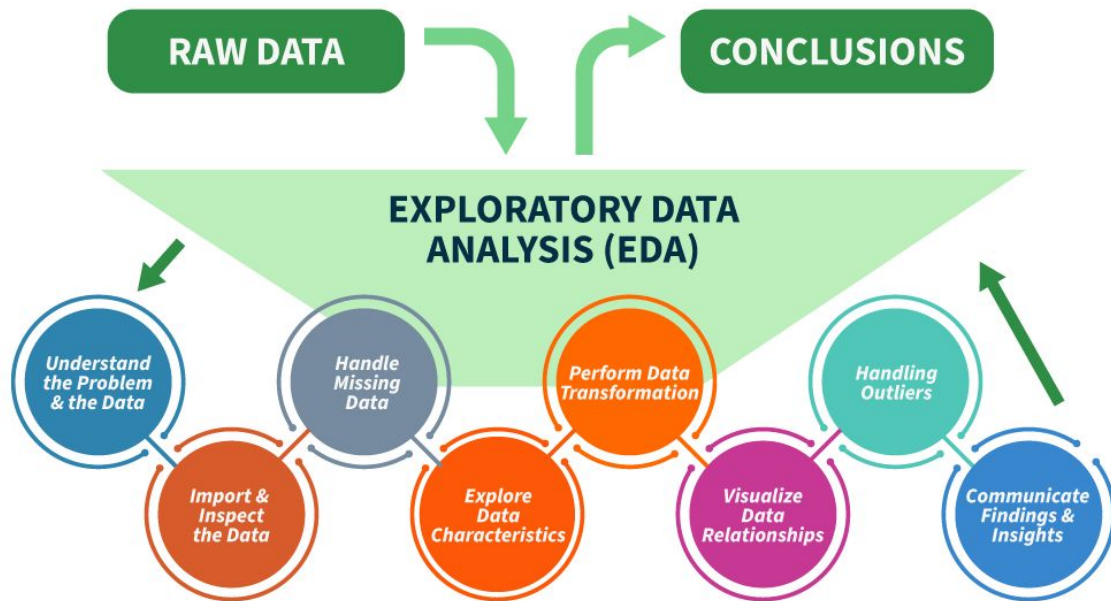
Векторизация

https://lena-voita.github.io/nlp_course/word_embeddings.html



Подробнее о EDA

Steps for Performing Exploratory Data Analysis



Начинаем обучение!



Representing Data

one sample

$$X = \begin{pmatrix} 1.1 & 2.2 & 3.4 & 5.6 & 1.0 \\ 6.7 & 0.5 & 0.4 & 2.6 & 1.6 \\ 2.4 & 9.3 & 7.3 & 6.4 & 2.8 \\ 1.5 & 0.0 & 4.3 & 8.3 & 3.4 \\ 0.5 & 3.5 & 8.1 & 3.6 & 4.6 \\ 5.1 & 9.7 & 3.5 & 7.9 & 5.1 \\ 3.7 & 7.8 & 2.6 & 3.2 & 6.3 \end{pmatrix}$$

one feature

$$y = \begin{pmatrix} 1.6 \\ 2.7 \\ 4.4 \\ 0.5 \\ 0.2 \\ 5.6 \\ 6.7 \end{pmatrix}$$

outputs / labels

Train и Test - для обучения с учителем



A diagram showing a dataset split into two parts. A large rectangle is divided into two equal halves by a vertical line. The left half is light gray and contains the text 'TRAIN - на этой части датасета мы обучаем нашу модель'. The right half is light purple and contains the text 'TEST - на этой части датасета мы тестируем нашу модель'.

TRAIN - на этой части датасета мы обучаем нашу модель

TEST - на этой части
датасета мы тестируем
нашу модель

Training and Test Data

training set

$$X = \begin{pmatrix} 1.1 & 2.2 & 3.4 & 5.6 & 1.0 \\ 6.7 & 0.5 & 0.4 & 2.6 & 1.6 \\ 2.4 & 9.3 & 7.3 & 6.4 & 2.8 \\ 1.5 & 0.0 & 4.3 & 8.3 & 3.4 \\ 0.5 & 3.5 & 8.1 & 3.6 & 4.6 \\ 5.1 & 9.7 & 3.5 & 7.9 & 5.1 \\ 3.7 & 7.8 & 2.6 & 3.2 & 6.3 \end{pmatrix}$$

test set

$$y = \begin{pmatrix} 1.6 \\ 2.7 \\ 4.4 \\ 0.5 \\ 0.2 \\ 5.6 \\ 6.7 \end{pmatrix}$$

Обучаем модель



Fit - обучаем

Predict - делаем предсказания
(первым делом как раз на тестовой
выборке и смотрим метрики)

Тестируем нашу модель

**Model Performance
on Training Data**



**Model Performance
on Test Data**



классификация

Confusion Matrix

<u><i>Actual</i></u> \ <u><i>Predict</i></u>	0	1
0	TN	FN
1	FP	TP

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1\ Score = \frac{2 * Precision * Recall}{Precision + Recall}$$

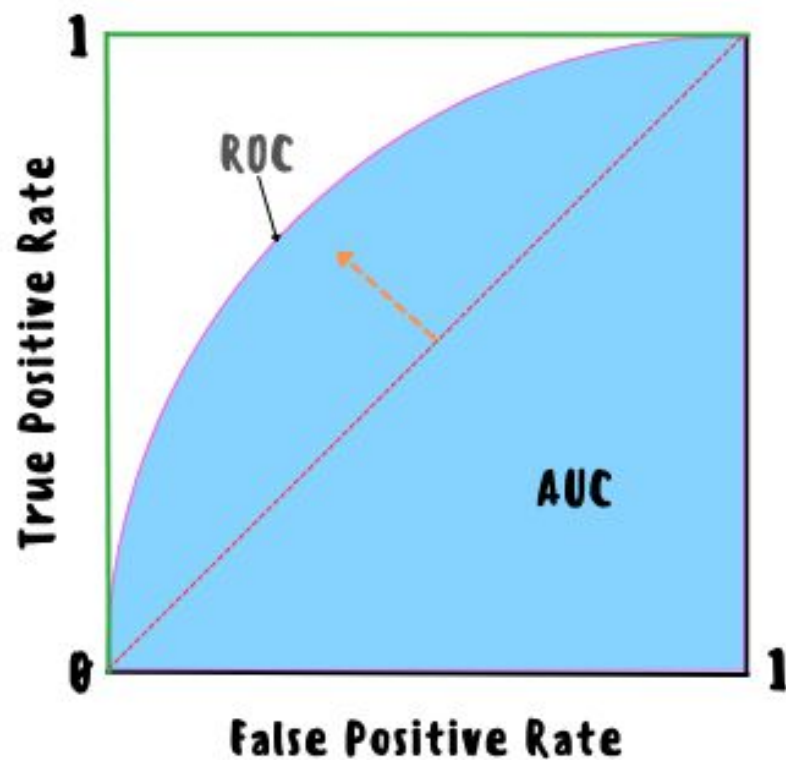
$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

Accuracy дает общую картину того, насколько можно полагаться на прогноз модели. Этот показатель не учитывает разницу между классами и типами ошибок. Поэтому он недостаточно хорош для несбалансированных наборов данных.

Precision показывает сколько реальных объектов класса среди всех тех, что классификатор отнес к этому классу.

Recall измеряет способность модели обнаруживать выборки, относящиеся к классу Positive.

F1-score - это среднее гармоническое значение между precision и recall.



- Perfect prediction
- Normal prediction
- - Random prediction
- - Better quality direction

Немного про LLM



**NOOO YOU CAN'T JUST MIX
UP ALL THE STEPS OF YOUR TASK
AND ASK AN LLM TO DO IT ALL.
HOW WILL YOU EVER MAKE A RELIABLE
AND EXTENSIBLE SYSTEM THAT WAY?**

imgflip.com



HAHA LLM GO BRRR

GPT-4o : The flagship model across audio, vision, and text by OpenAI	Grok-2 : Grok-2 by xAI	Claude 3.5 : Claude by Anthropic
Llama 3.1 : Open foundation and chat models by Meta	Gemini : Gemini by Google	Mixtral of experts : A Mixture-of-Experts model by Mistral AI
GPT-4-Turbo : GPT-4-Turbo by OpenAI	Jamba 1.5 : Jamba by AI21 Labs	Gemma 2 : Gemma 2 by Google
Claude : Claude by Anthropic	DeepSeek Coder v2 : An advanced code model by DeepSeek	Nemotron-4 340B : Cutting-edge Open model by Nvidia
Llama 3 : Open foundation and chat models by Meta	Athene-70B : A large language model by NexusFlow	Qwen Max : The Frontier Qwen Model by Alibaba
GPT-3.5 : GPT-3.5-Turbo by OpenAI	Yi-Large : State-of-the-art model by 01 AI	Yi-Chat : A large language model by 01 AI
Phi-3 : A capable and cost-effective small language models (SLMs) by Microsoft	Reka Core : Frontier Multimodal Language Model by Reka	Reka Flash : Multimodal model by Reka
Command-R-Plus : Command R+ by Cohere	Command R : Command R by Cohere	Qwen 1.5 : The First 100B+ Model of the Qwen1.5 Series
GLM-4 : Next-Gen Foundation Model by Zhipu AI	DBRX Instruct : DBRX by Databricks Mosaic AI	InternLM : A multi-language large-scale language model (LLM), developed by SHLAB.

HuggingFace

<https://huggingface.co/models>

