



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ  
УНИВЕРСИТЕТ

# Компьютерная лингвистика и информационные технологии

Неделя 1: Предобработка данных



# Описание курса

- Блоки “Информационный поиск” и “Компьютерная лингвистика”;
- Преп.: Олег Сериков (@oserikov), Мария Пономарева (@MashPo), Влад Михайлов (@vmkhlv);
- Ассист.: Дарья Самсонова (@SamsonovaDaria), Кирилл Коновалов (@kirillknv);
- [https://github.com/vmkhlv/hse\\_compling\\_and\\_it](https://github.com/vmkhlv/hse_compling_and_it);
- 1 модуль:  $0.8 * \text{ДЗ} + 0.2 * \text{Тесты}$ ;
- 3 модуль:  $0.5 * \text{ДЗ} + 0.1 * \text{Тесты} + 0.2 * \text{Контроль (Модуль 2)} + 0.2 * \text{Экзамен (Модуль 3)}$ ;



# Блок “Компьютерная лингвистика”

- Предобработка данных;
- Базовые понятия машинного обучения, алгоритмы машинного обучения для задач классификации, кластеризации и регрессии;
- Архитектуры нейронных сетей: FFN, CNN, RNN, LSTM;
- Sequence-to-sequence, Sequence Tagging, Language Modeling;
- Статические эмбединги: Word2Vec, FastText;
- Контекстуальные эмбединги: BERT, ELMo.



# Фреймворки



# Keras

A deep learning library

# PYTORCH

Deep Learning with PyTorch





hows life baby gorilla

babygirl\*

About 634,000,000 results (0.62 seconds)

Russian – detected ▼



English ▼

Можно поспать?



Can I sleep?

Mozhno pospat'?

Говорят ты  
сегодня нагрубил



Feedback

✓ 21:56

Олег

Роза моего искусственного  
интеллекта цвела не ради  
этого вопроса.



21:56



# Как учится машина? Основные понятия

- Объект (картинка, текст, аудио и д.р.);
- Целевая функция / таргет (класс объекта / чиселка);
- Признаки объекта (возраст, пол, уровень дохода и т.д.);
- Матрица объектов-признаков;
- Обучающая / валидационная / тестовая выборки (train / val / test);
- Модель / алгоритм – это решающая функция;
- Функция потерь / ошибок;
- Метрика.



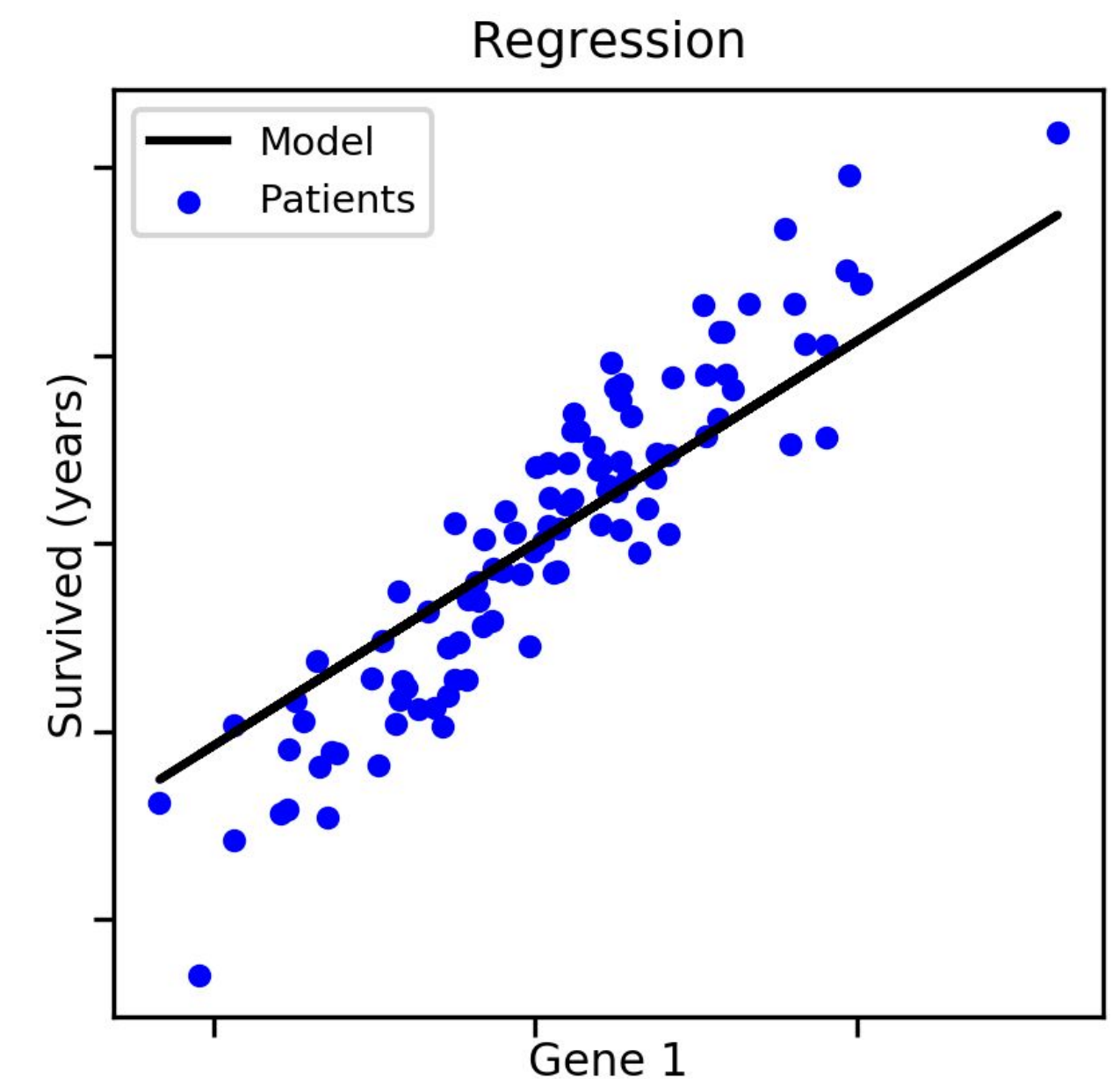
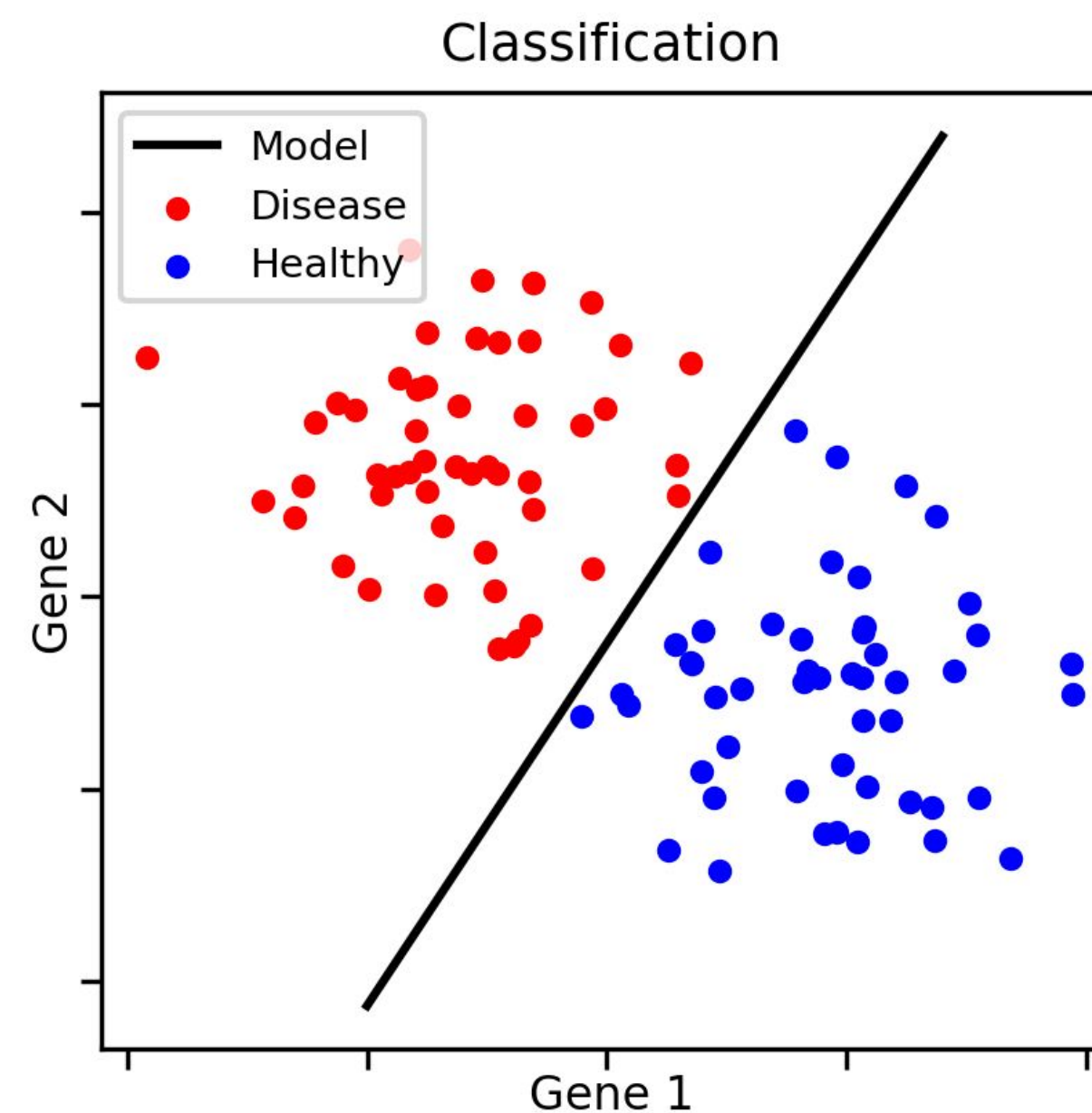
# Целевая функция определяет задачу

## Классификация

- Фильтрация спама;
- Определение темы сообщения;
- Определение языка.

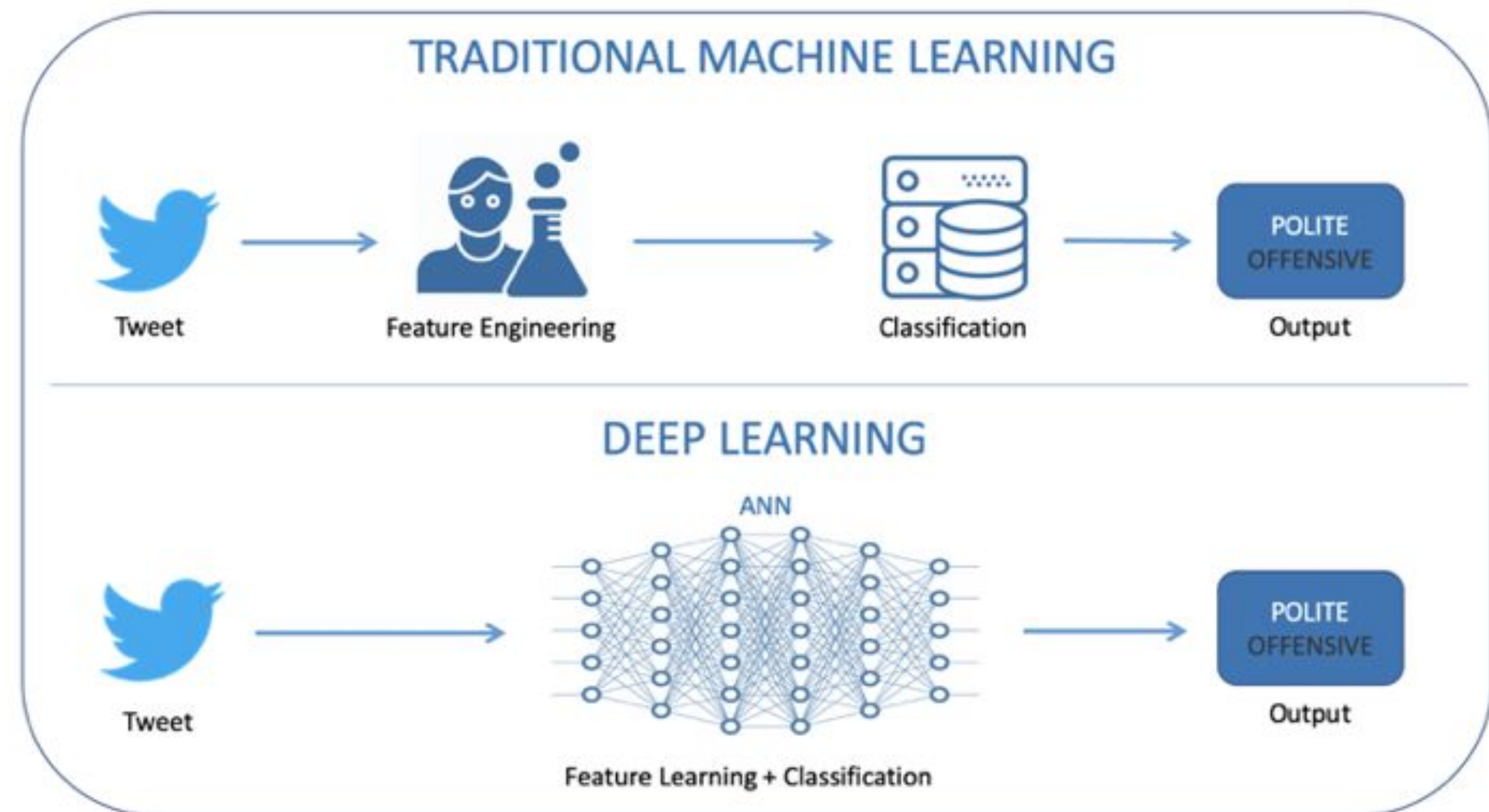
## Регрессия

- Количество лайков;
- Количество зараженных COVID-19.



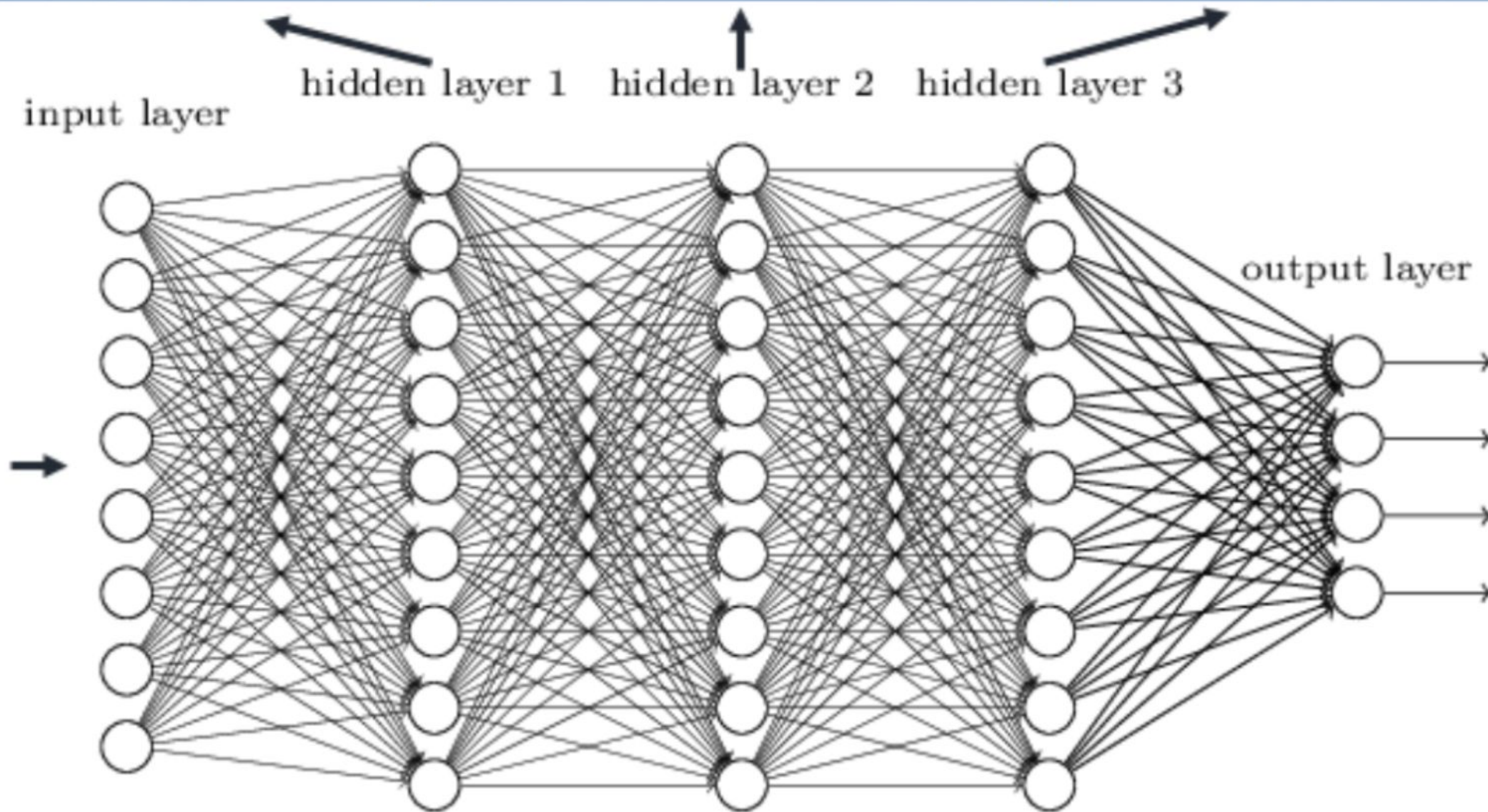
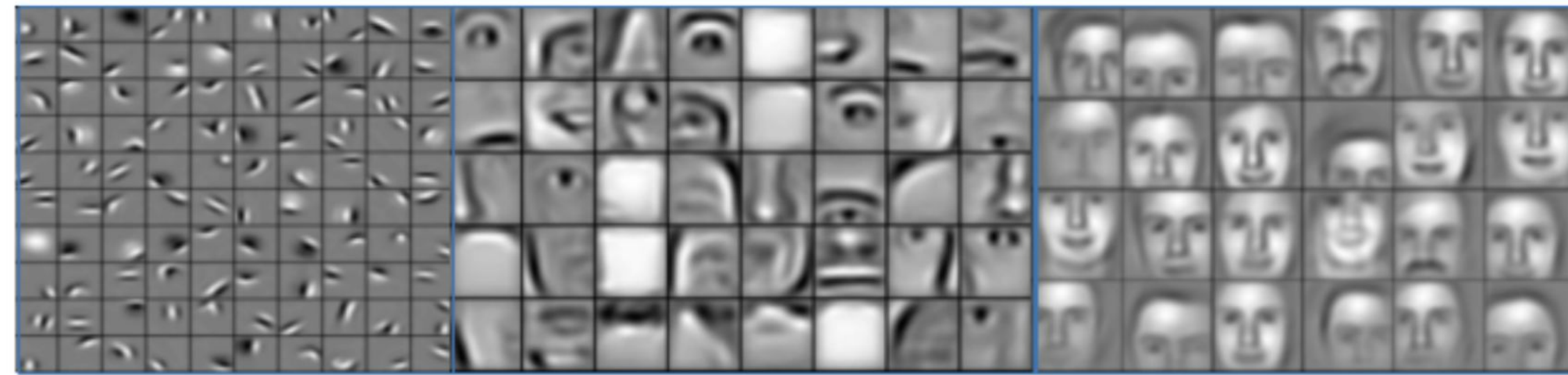
# Важность признаков

- Конструирование;
- Извлечение;
- Отбор;
- Определение важности;
- Интерпретация поведения.





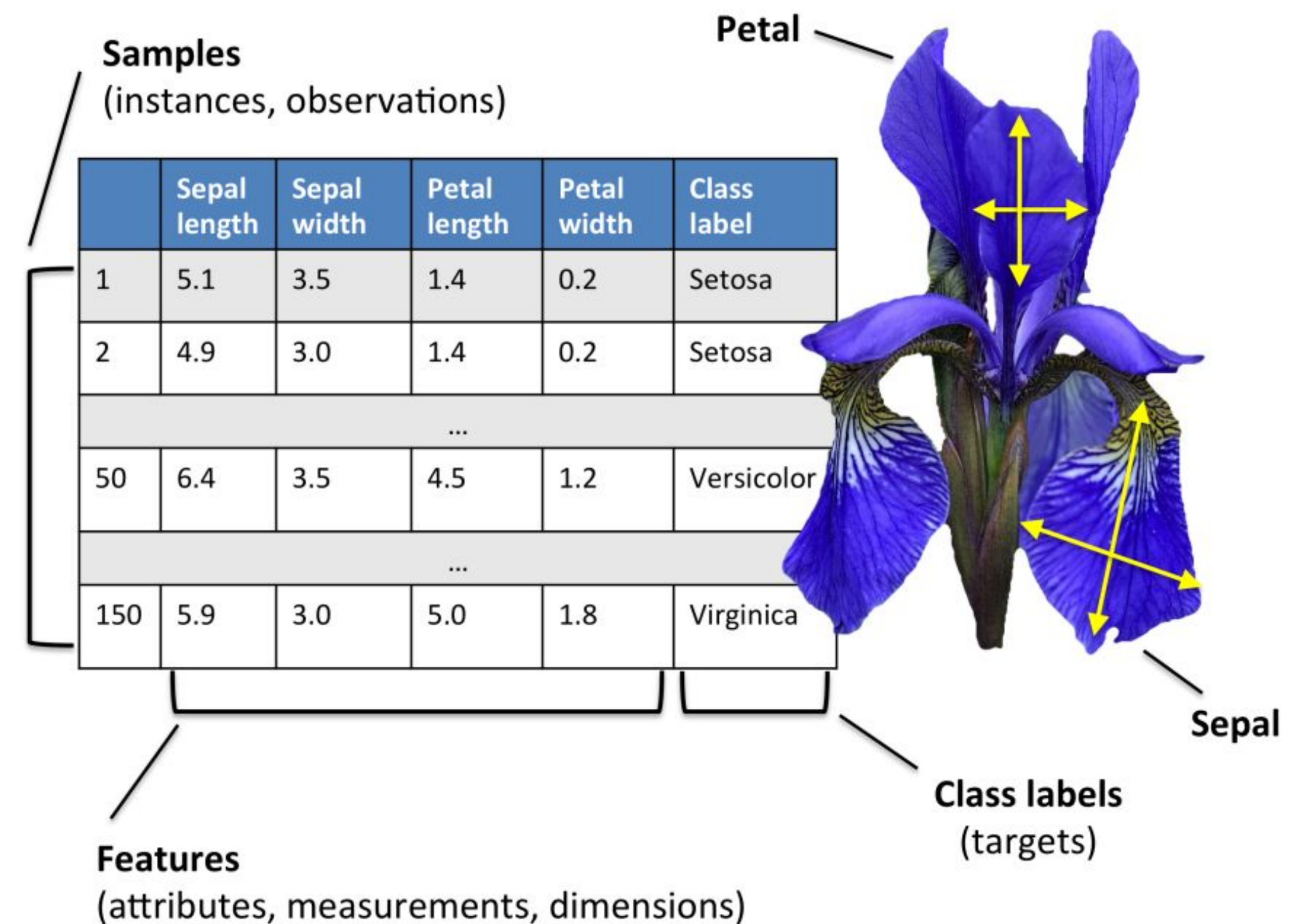
Deep neural networks learn hierarchical feature representations





# Типы признаков

- Чиселки (numerical) – доход, объем, масса;
- Бинарные (binary) – пришли / не пришли на пару;
- Категориальные (categorical) – пол, уровень английского;
- Порядковые – курс.



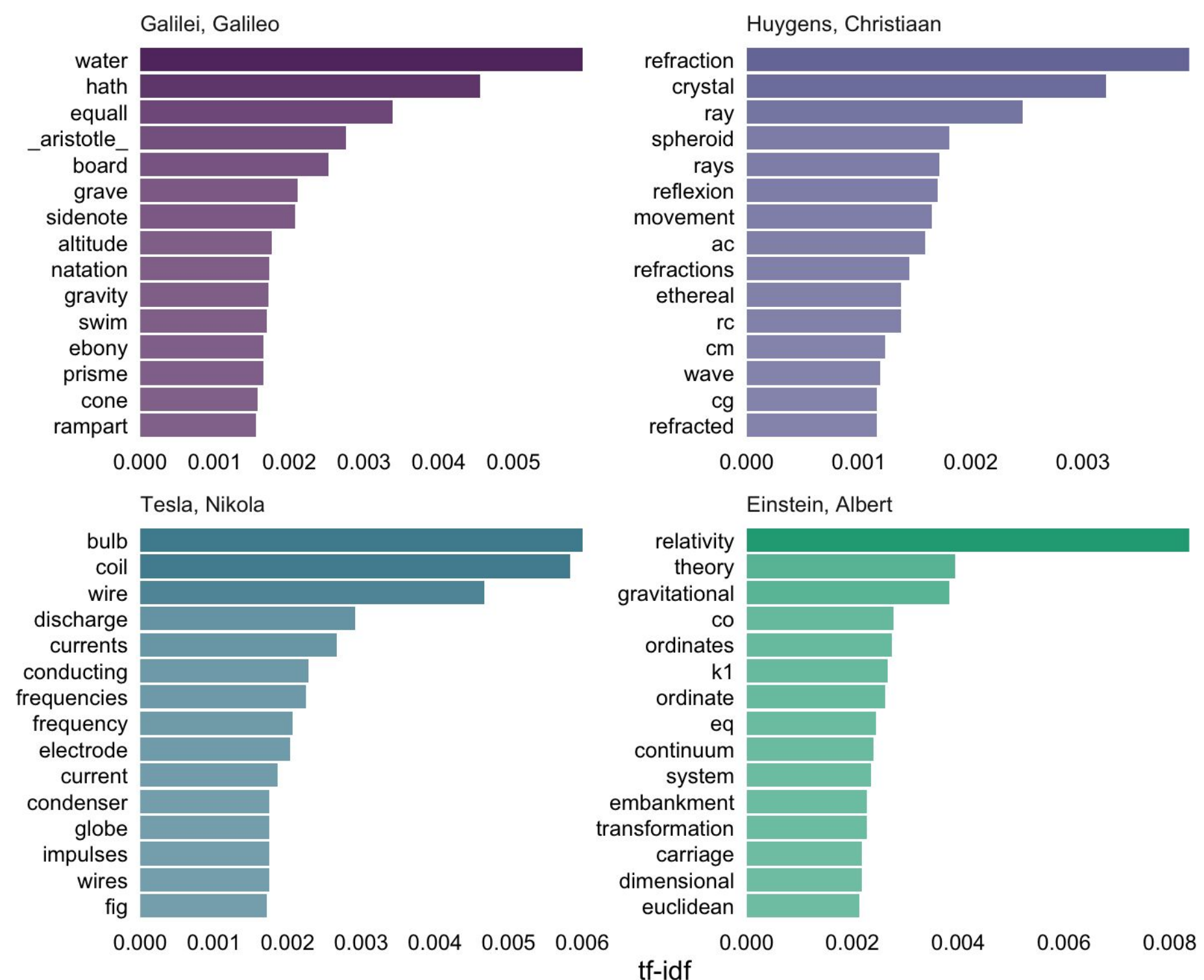
# Текстовые признаки

Текст – это упорядоченная *последовательность*.

*Например:*

- Средняя длина;
- Средняя частотность по коллекции;
- IPM (instance per million);
- Среднее количество сущностей / документ;
- Частотность пос-тегов по коллекции;
- Частотность N-грамм пос-тегов / лемм;
- Automated Readability Index (ARI).

Highest tf-idf words in Classic Physics Texts





# Решаем ЕГЭ

- CatBoostClassifier;
- POS N-grams;
- Bag of words;
- Ensemble + Voting.

## 1. Задание #T29871



Расставьте знаки препинания.

Он стоял перед дворцом во время обеда государя (1) глядя в окна дворца (2) ожидая чего-то ещё и (3) завидуя одинаково и сановникам (4) подъезжавшим к крыльцу (5) и (6) мелькавшим в окнах (7) камер-лакеям.

Укажите цифру(-ы), на месте которой(-ых) в предложении должна(-ы) стоять запятая(-ые).

Ответ



Проверить ответ

[Показать разбор и ответ](#)





# Определение автора

- **Topic or Style? Exploring the Most Useful Features for Authorship Attribution (Sari et al., 2018)**

Type	Group	Category	#	Description
Style	Lexical	Word-level	2	Average word length, number of short words
		Char-level	2	Percentage of digits, percentage of upper-case letters
		Letters	26	Letter frequency
		Digits	10	Digit frequency
	Syntactic	Vocabulary richness	2	Richness (hapax-legomena and dis-legomena)
		Function words	174	Frequency of function words
		Punctuation	12	Occurrence of punctuation
Content	Word <i>n</i> -gram	Words unigrams	100	Frequency of 100 most common word unigrams
		Words bigrams	100	Frequency of 100 most common word bigrams
		Word trigrams	100	Frequency of 100 most common word trigrams
Hybrid	Char <i>n</i> -gram	Char bigrams	100	Frequency of 100 most common character bigrams
		Char trigrams	100	Frequency of 100 most common character trigrams

Table 3: Authorship attribution feature sets.



# Предобработка

- Сегментация: *rusenttokenize, razdel*
- Токенизация: *spacy\_russian\_tokenizer, razdel*
- Стоп-слова: *nltk stopwords*
- Морфология: *rumorphy2, mystem, slovnet, RNNMorph*
- Синтаксис: *UDPipe, slovnet, stanza, GramEval2020*

	raw_word	stemmed_word
0	..trouble..	..trouble..
1	trouble<	trouble<
2	trouble!	trouble!
3	<a>trouble</a>	<a>trouble</a>
4	1.trouble	1.troubl

# Зачем все это?

- Предварительный анализ данных;
- Фильтрация данных;
- Составление датасетов;
- Аугментация данных;
- Конструирование признаков;
- Составление выборок;
- Анализ поведения модели и ошибок.

