





- Наконец пришло время погрузиться в детали машинного обучения!
- Этот раздел поможет Вам настроиться на нужный образ мышления, свойственный для машинного обучения
- Для начала давайте вспомним наши этапы проекта по машинному обучению...







мир

Решить задачу

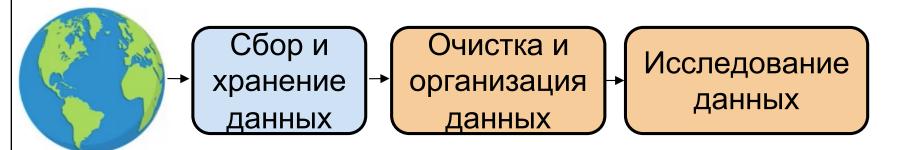
Как починить или поменять X?

Ответить на вопрос

Как изменение в X повлияет на Y?





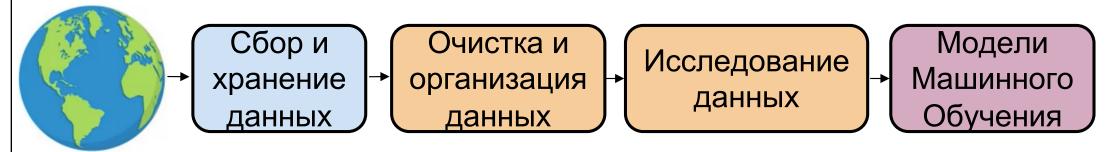


Реальный мир

Статистический анализ, Визуализация данных







Реальный мир

Supervised Learning:

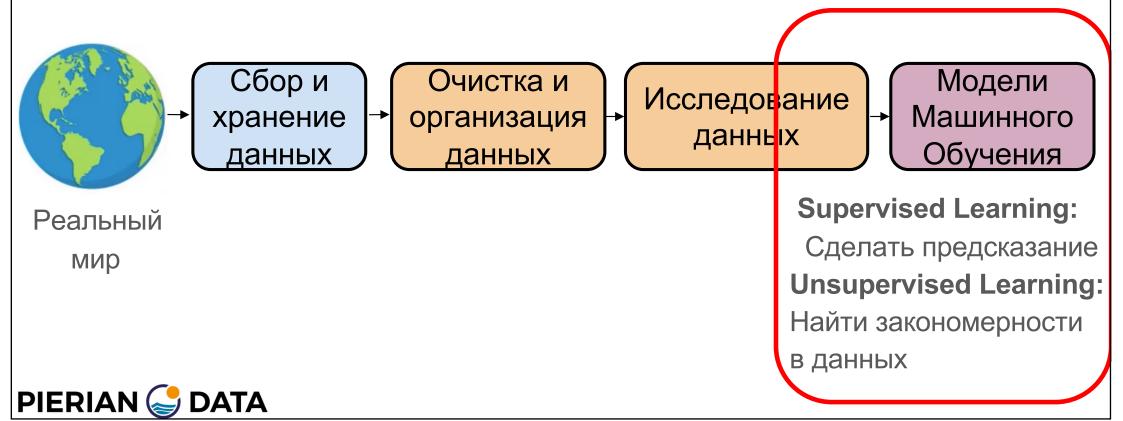
Сделать предсказание

Unsupervised Learning:

Найти закономерности в данных









- В этом разделе мы обсудим:
 - Какие задачи решает машинное обучение
 - Типы машинного обучения:
 - Supervised Learning
 - Unsupervised Learning
 - Процесс машинного обучения для Supervised Learning





- По мере продвижения вперёд мы будем постепенно изучать многие темы:
 - Баланс между смещением и дисперсией
 - Кросс-валидация
 - Построение признаков (Feature Engineering)
 - Scikit-learn
 - Метрики оценки модели и многое другое!





- Разделы по машинному обучению
 - Для каждого вида алгоритма
 - Интуиция и математическая теория
 - Пример написания кода для применения алгоритма
 - Расширения алгоритма
 - Упражнения с решениями





- Разделы по машинному обучению
 - О Для линейной регрессии немного по-другому
 - Интуиция и математическая теория
 - Простая линейная регрессия
 - Scikit-learn и линейная регрессия
 - Регуляризация





- Разделы по машинному обучению
 - Изучение дополнительных тем по машинному обучению
 - Метрики оценки модели
 - Построение признаков (Feature Engineering)
 - Кросс-валидация
 - Применение этих тем для упражнения по линейной регрессии





• Далее мы посмотрим, зачем нужно машинное обучение, и в каких случаях оно применяется!









- Машинное обучение это изучение статистических компьютерных алгоритмов, которые автоматически улучшаются на данных
- В отличие от обычных алгоритмов, где человек должен указывать компьютеру, какое направление выбрать, здесь алгоритмы сами выбирают лучший подход на основе входных данных





- Машинное обучение (machine learning) является частью более общего направления искусственный интеллект (artificial intelligence)
- Алгоритмы машинного обучения не программируются явно на то или иное решение
- Вместо этого, алгоритмы определяют оптимальные решения на основе данных





- Примеры задач, решаемых с помощью машинного обучения:
 - Кредитный скоринг
 - Риски страхования
 - Предсказание цен
 - Фильтрация спама
 - Сегментация клиентов
 - И многое другое!





- Подход к постановке задачи по машинному обучению:
 - На основе признаков (features) в наборе данных, найти требуемую целевую переменную (label/target)
 - Алгоритмы машинного обучения часто ещё называют словом "estimator", поскольку они оценивают (estimate) целевую переменную





- Каким образом алгоритмы машинного обучения могут решать столько разных задач?
- Эти алгоритмы с помощью статистических методов обрабатывают данные, в результате чего выясняют, какие признаки являются важными в данных





- Простой пример:
 - Предсказать цену продажи дома на основе его характеристик (район, количество комнат и т.д.)





- Предсказания цены дома:
 - Обычный алгоритм человек пишет алгоритм, в котором указывает значение важности для каждого признака
 - Алгоритм машинного обучения алгоритм сам определяет важность признаков, на основе самих данных





- Зачем нужно машинное обучение?
 - Многие сложные задачи проще решить с помощью методов машинного обучения
 - Такие задачи, как выявление спама, распознавание написанного вручную текста, эффективно решаются с помощью машинного обучения





- Почему бы не использовать машинное обучение повсеместно?
 - Главное требование алгоритмов ML наличие хороших данных
 - Большая часть времени на построение модели тратится на очистку и подготовку данных, а вовсе не на реализацию алгоритма машинного обучения





- Нужно ли разрабатывать свои алгоритмы машинного обучения?
 - Это нужно очень редко, поскольку существующие алгоритмы хорошо реализованы и задокументированы





• Далее мы обсудим, какие бывают типы алгоритмов машинного обучения!









- В ближайших разделах мы изучим два типа алгоритмов машинного обучения:
 - Supervised Learning обучение с учителем
 - Unsupervised Learning обучение без учителя





- Supervised Learning обучение с учителем
 - Используются размеченные данные или исторические данные, на основе которых делается предсказание целевой переменной
- Unsupervised Learning обучение без учителя
 - Применяется к неразмеченным данным. Модель машинного обучения ищет возможные закономерности в данных





- Supervised Learning обучение с учителем
 - О Требуются исторические размеченные данные
 - Исторические известные результаты из моментов в прошлом
 - Размеченные значения целевой переменной были известны





- Supervised Learning обучение с учителем
 - Целевые переменные бывают двух типов
 - Нужно предсказать категориальную переменную это задача классификации
 - Нужно предсказать непрерывную переменную это задача регрессии





- Supervised Learning обучение с учителем
 - Задачи классификации
 - Предсказать категорию
 - Доброкачественная или раковая опухоль
 - Погашение кредита или неплатежи
 - Категория изображения например, конкретная буква в задаче распознавания текста, написанного от руки





- Supervised Learning обучение с учителем
 - Задачи регрессии
 - Предсказать значение непрерывной переменной
 - Будущие цены
 - Нагрузка на электросети
 - Результаты тестов





- Unsupervised Learning обучение без учителя
 - О Сгруппировать неразмеченные данные
 - Пример:
 - Сгруппировать клиентов на группы, используя их поведенческие характеристики





- Unsupervised Learning обучение без учителя
 - Основная сложность в том, что у нас нет исторических "правильных" данных, и поэтому намного сложнее оценить правильность работы модели машинного обучения





- Разделы по машинному обучению
 - Мы начнём с алгоритмов обучения с учителем (supervised learning), чтобы изучить возможности различных алгоритмов машинного обучения
 - Далее мы перейдём к алгоритмам обучения без учителя (unsupervised learning), для кластеризации данных и уменьшения размерности данных





• Прежде чем заняться написанием кода и линейной регрессией в следующем разделе, давайте детально посмотрим на весь процесс машинного обучения с учителем (supervised learning), чтобы подготовиться к следующим разделам!





Процесс машинного обучения с учителем

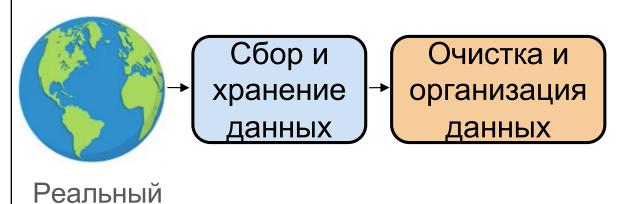








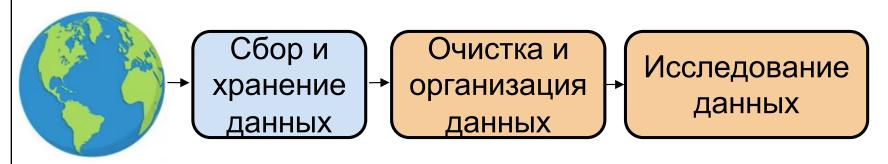




МИР



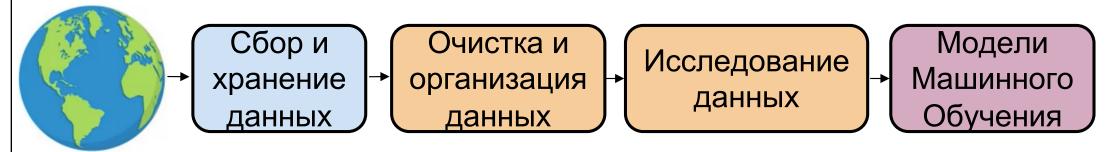




Реальный мир







Реальный мир

Supervised Learning:

Сделать предсказание

Unsupervised Learning:

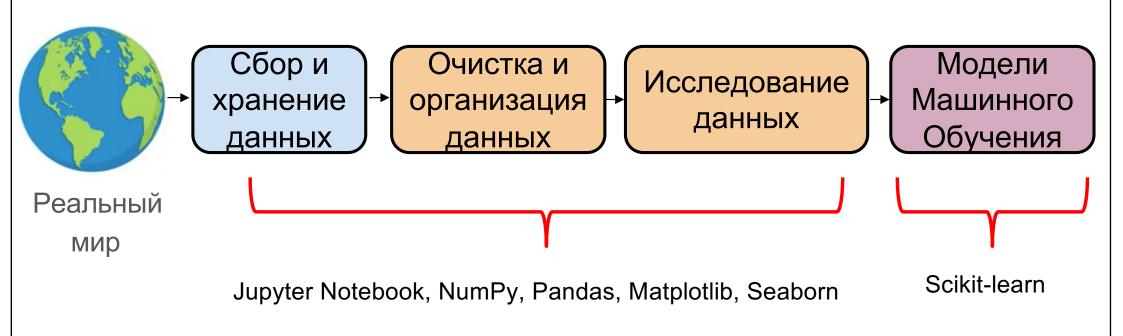
Найти закономерности в данных



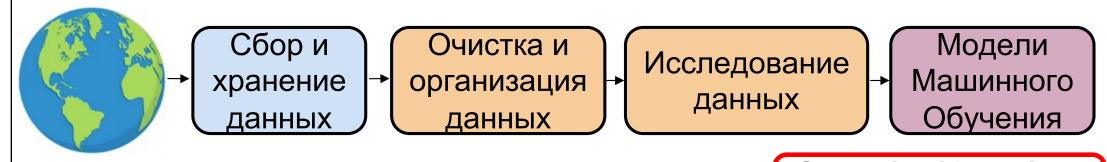


PIERIAN 🈂 DATA

Этапы работ по машинному обучению





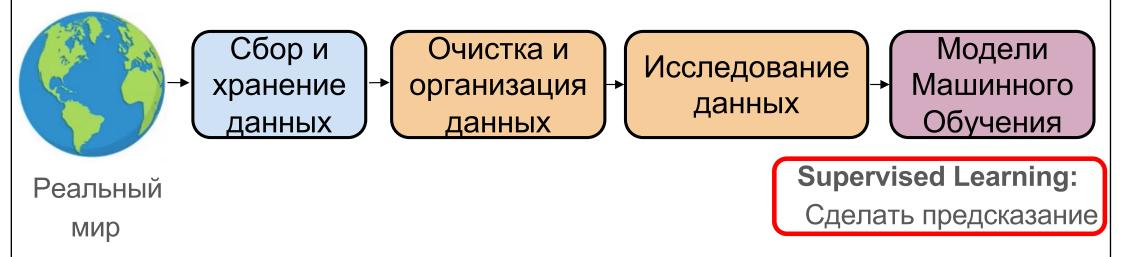


Реальный мир

Supervised Learning: Сделать предсказание







Например: предсказать цену продажи дома





Процесс машинного обучения с учителем

- Начинаем со сбора и структуризации исторических данных
- Предыдущие продажи домов это размеченные данные

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





Процесс машинного обучения с учителем

 Задача – предсказать цену продажи для нового дома, с известными значениями прощади, кол-ва спален и кол-ва

санузлов

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





Процесс машинного обучения с учителем

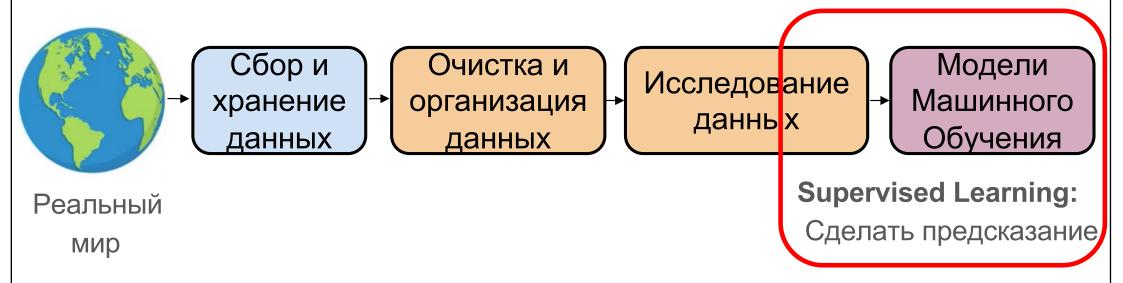
• Продукт на основе данных: на входе параметры дома, на выходе предсказание цены (на основе исторических данных)

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





• Задача предсказания цены дома







• Задача предсказания цены дома

Модели Машинного Обучения

Supervised Learning:

Сделать предсказание





• Задача предсказания цены дома

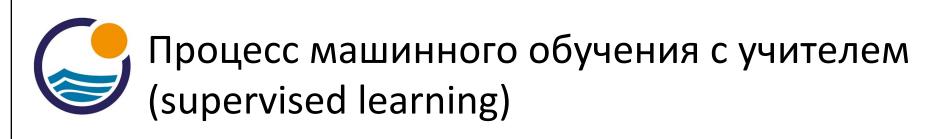
Модели Машинного Обучения

Supervised Learning:

Сделать предсказание

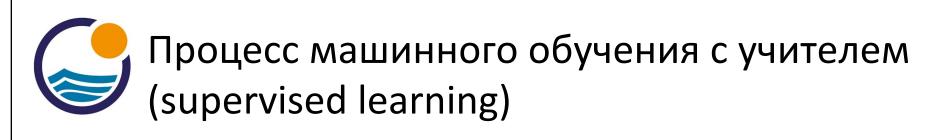
Данные

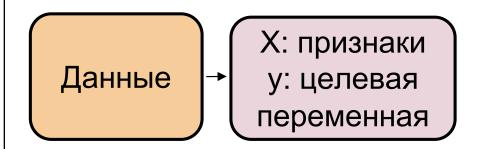
















 Целевая переменная – значение, которое нам нужно предсказать

Данные → Х: признаки у: целевая переменная

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





 Целевая переменная – значение, которое нам нужно предсказать

Данные → X: признаки у: целевая переменная

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





• Признаки (features) – известные характеристики объекта

Данные → Х: признаки у: целевая переменная

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





• Признаки и целевая переменная определяются в соответствии с решаемой задачей

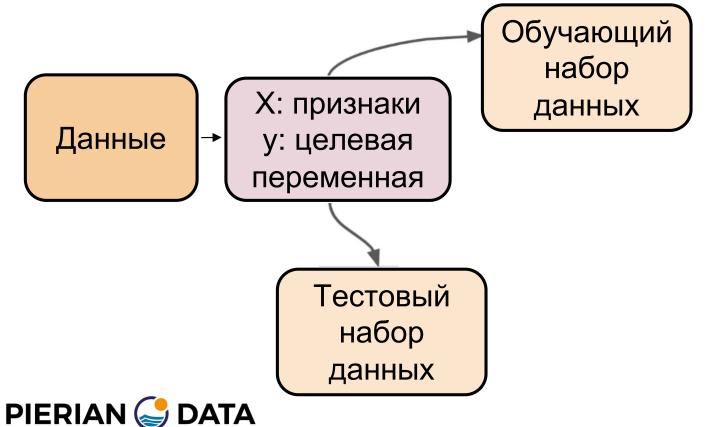
Данные → Х: признаки у: целевая переменная

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000



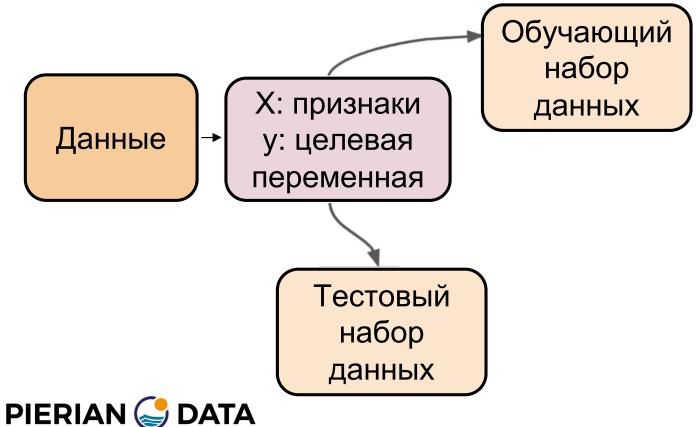


• Разбиваем данные на обучающий и тестовый наборы данных





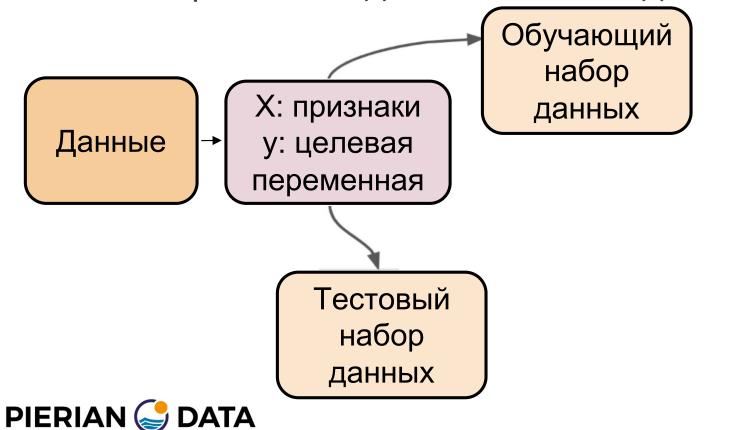
• Разбиваем данные на обучающий и тестовый наборы данных



 Позже мы обсудим кроссвалидацию



• Зачем разбивать данные? Как это делать?





• Зачем разбивать данные? Как это делать?

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





• Как бы Вы оценили работу обычного риэлтора?



Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





Как бы Вы оценили работу обычного риелтора?

• Вы могли бы попросить риелтора посмотреть на

исторические данные...



Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





 Затем дать риелтору характеристики нового дома, и попросить сделать оценку цены



Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





• Но как понять, насколько оценка цены будет корректна?



Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





Но как понять, насколько оценка цены будет корректна?
 Какой дом выбрать для тестирования оценки?



Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





Проверить оценку для нового дома нельзя, потому что его истинная цена продажи ещё неизвестна



Area m²	Bedrooms Bathroom		Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





 Проверить оценку для имеющихся домов тоже нельзя, потому что риелтор мог просто запомнить эти данные



Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





Поэтому следует разбить данные на обучающий и тестовый наборы – давайте посмотрим, зачем...



Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





 Теперь мы разбиваем эти данные на обучающий набор данных (TRAIN) и тестовый набор данных (TEST)

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000



TRAIN



 Теперь мы разбиваем эти данные на обучающий набор данных (TRAIN) и тестовый набор данных (TEST)

	Area m²	Bedrooms	Bathrooms	Price		
	200	3	2	\$500,000		
TRAIN	190	2	1	\$450,000		
	230	3	3	\$650,000		
TEST	180	1	1	\$400,000		
	210	2	2	\$550,000		
PIERIAN 🈂 DATA						



• У нас получилось 4 компонента:

	Area m²	Bedrooms	Bathrooms	Price	
	200	3	2	\$500,000	V TDAIN
X TRAIN	190	2	1	\$450,000	Y TRAIN
	230	3	3	\$650,000	
X TEST	180	1	1	\$400,000	Y TEST
	210	2	2	\$550,000	
PIERIAN 🈂 DATA					



 Итак, вернёмся к вопросу о том, как мы можем проверить работу риелтора...



Area m²	Bedrooms Bathroom		Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





Итак, вернёмся к вопросу о том, как мы можем проверить работу риелтора...



TRAIN

TEST

Area m²	Bedrooms Bathrooms		Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000
180	1	1	\$400,000
210	2	2	\$550,000





 Мы можем дать риелтору только данные из обучающего набора данных (TRAIN) – как X, так и Y



TRAIN

Area m²	Bedrooms	Bathrooms	Price
200	3	2	\$500,000
190	2	1	\$450,000
230	3	3	\$650,000





 После "обучения" на обучающем наборе данных, просим риелтора сделать оценки цены дома для тестовых данных



 Area m²
 Bedrooms
 Bathrooms

 180
 1
 1

 210
 2
 2





- После "обучения" на обучающем наборе данных, просим риелтора сделать оценки цены дома для тестовых данных
- Предоставляем только признаки X, просим оценить у



TEST

	Area m²	Bedrooms	Bathrooms
-	180	1	1
	210	2	2





- После "обучения" на обучающем наборе данных, просим риелтора сделать оценки цены дома для тестовых данных
- Предоставляем только признаки X, просим оценить у
- Риелтор ещё не видел этих данных! И цену этих домов!



TEST

	Area m²	Bedrooms	Bathrooms
-	180	1	1
	210	2	2





Риелтор даёт свои оценки – предсказания цены (predictions)



Predictions	Area m²	Bedrooms	Bathrooms
\$410,000	180	1	1
\$540,000	210	2	2





- Риелтор даёт свои оценки предсказания цены (predictions)
- Далее мы добавляем истинные цены



Predictions	Area m²	Bedrooms	Bathrooms	Price
\$410,000	180	1	1	\$400,000
\$540,000	210	2	2	\$550,000





- Риелтор даёт свои оценки предсказания цены (predictions)
- Далее мы добавляем истинные цены
- И сравниваем предсказания цены и истинные цены



Predictions	Price
\$410,000	\$400,000
\$540,000	\$550,000





ullet Часто обозначения такие — оценка $\hat{oldsymbol{y}}$ и истинное значение $oldsymbol{y}$





Predictions	Price
\$410,000	\$400,000
\$540,000	\$550,000





- ullet Часто обозначения такие оценка $\hat{oldsymbol{y}}$ и истинное значение $oldsymbol{y}$
- Позже мы обсудим различные методы для понимания того, насколько эти оценки хороши!



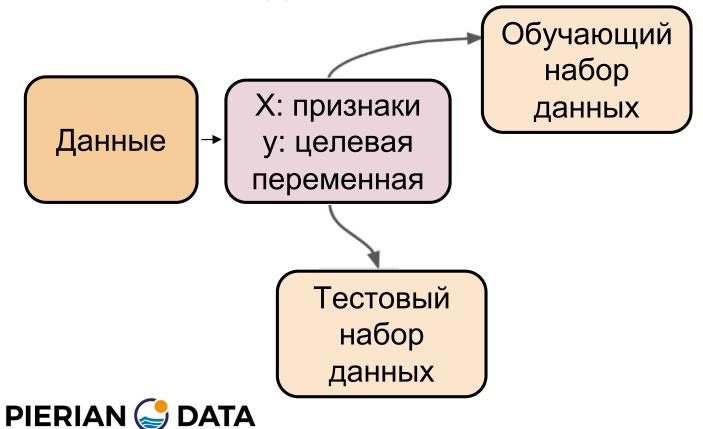


Predictions	Price
\$410,000	\$400,000
\$540,000	\$550,000



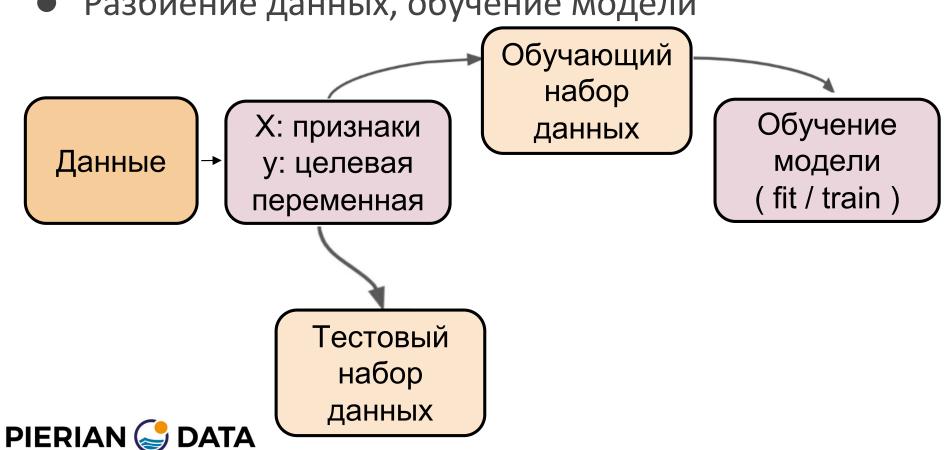


• Разбиение данных



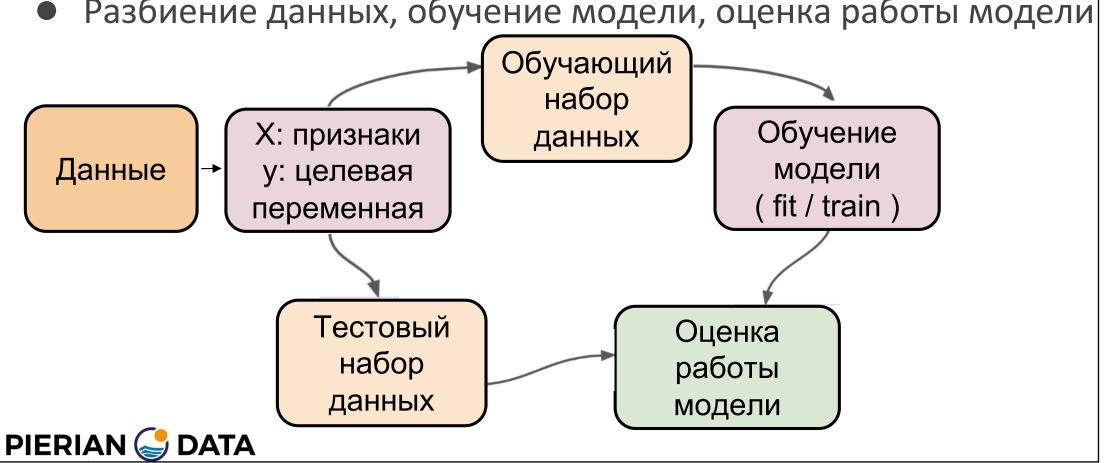


Разбиение данных, обучение модели



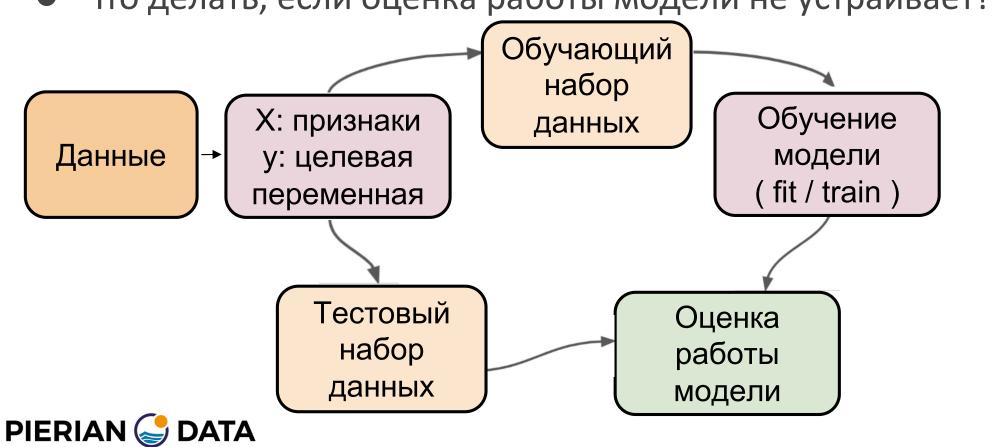


Разбиение данных, обучение модели, оценка работы модели



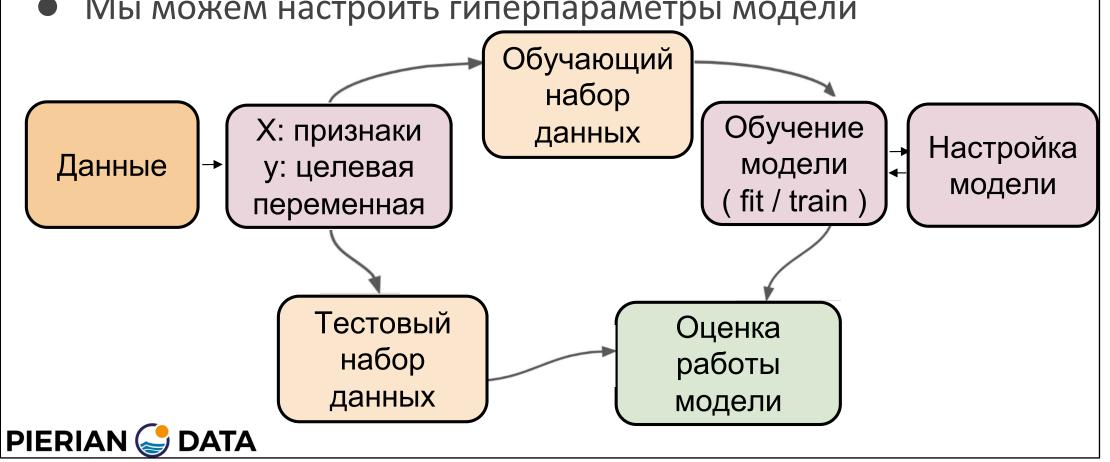


• Что делать, если оценка работы модели не устраивает?





Мы можем настроить гиперпараметры модели





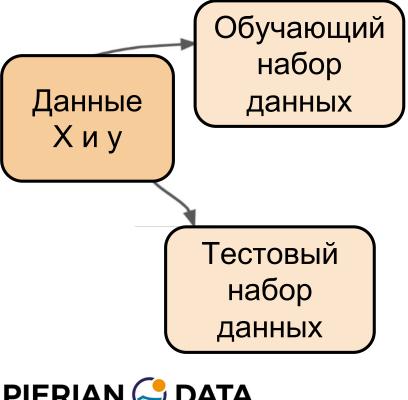
- Давайте быстро повторим весь процесс
- Получаем данные признаки X и целевую переменную у

Данные Хиу





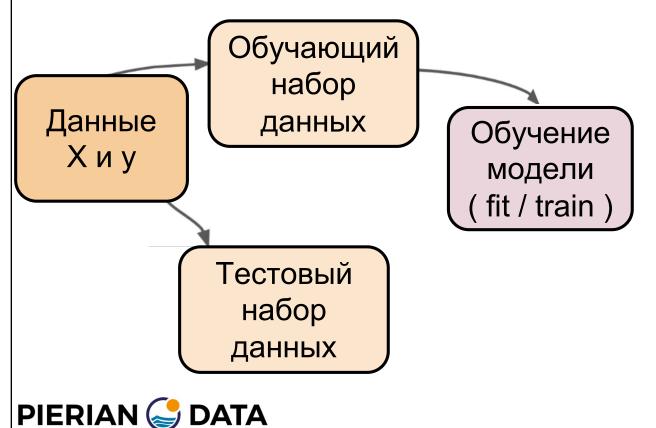
Разбиваем данные на обучающий и тестовый наборы данных





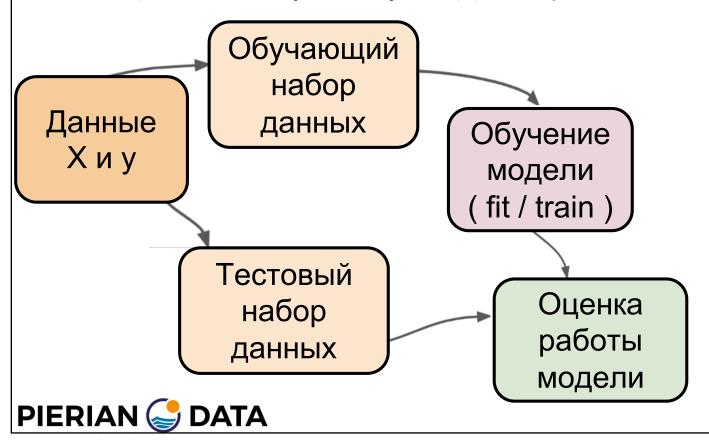


• Обучаем модель на обучающем наборе данных



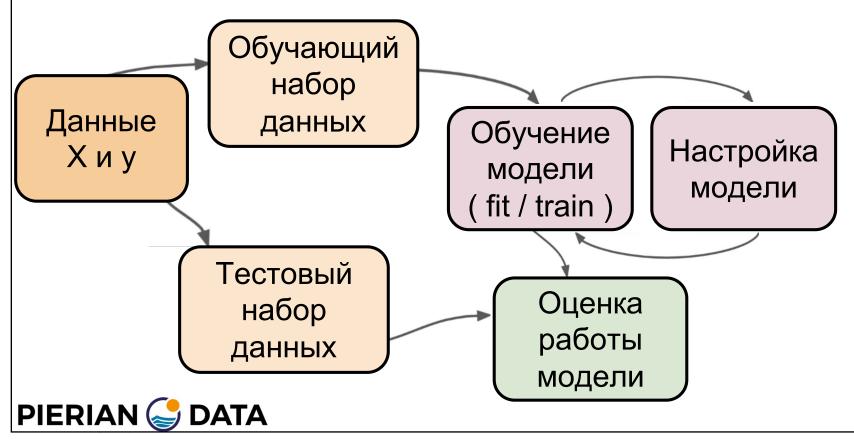


Оцениваем работу модели (evaluate model performance)



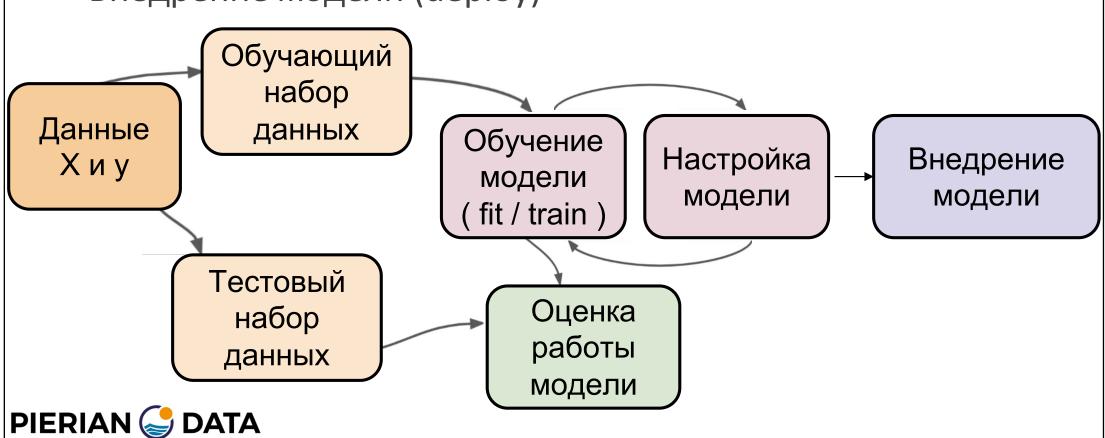


Настраиваем гиперпараметры модели (несколько раз)





Внедрение модели (deploy)





Этапы работ по машинному обучению

Процесс машинного обучения с учителем (supervised learning)

