

TECHNOLOGIES IN EDUCATION
UNIVERSITY NSU

MICROELECTRONICS
INNOVATIONS
CATALYTIC
MATERIALS
ASSEMBLY
POINT **DRUG**
DESIGN

SCIENTIFIC
LABORATORY
HYBRID
MATERIALS
GEOPHYSICS
ENGINEERING
ENERGY CONSERVATION
BIOTECHNOLOGY
GEOCHEMISTRY
NANOTECHNOLOGY

HIGH
ENERGIES
SEMIOTICS
SCIENCE
MATHEMATICAL MODELING

IT
DEEP
LEARNING
BRAIN
STUDY
COGNITIVE

DEVELOPMENT
ELEMENTARY
PARTICLES
THE ARCTIC REGIONS
DARK
MATTER

QUANTUM
TECHNOLOGIES
BIOMEDICINE
APPLIED
STUDIES
PHOTONICS
ASTRONOMY
GLOBAL PRIORITY
ASTROPHYSICS
BIOINFORMATICS

LASER
PHYSICS
KNOWLEDGE
ECONOMY
GEOLOGY
ARCHEOLOGY
TECHNOLOGIES

N* Novosibirsk
State
University
***THE REAL SCIENCE**

Базовые методы ИИ

Семинар 2

Глушенко Андрей Валерьевич
ФФ НГУ

Регрессия

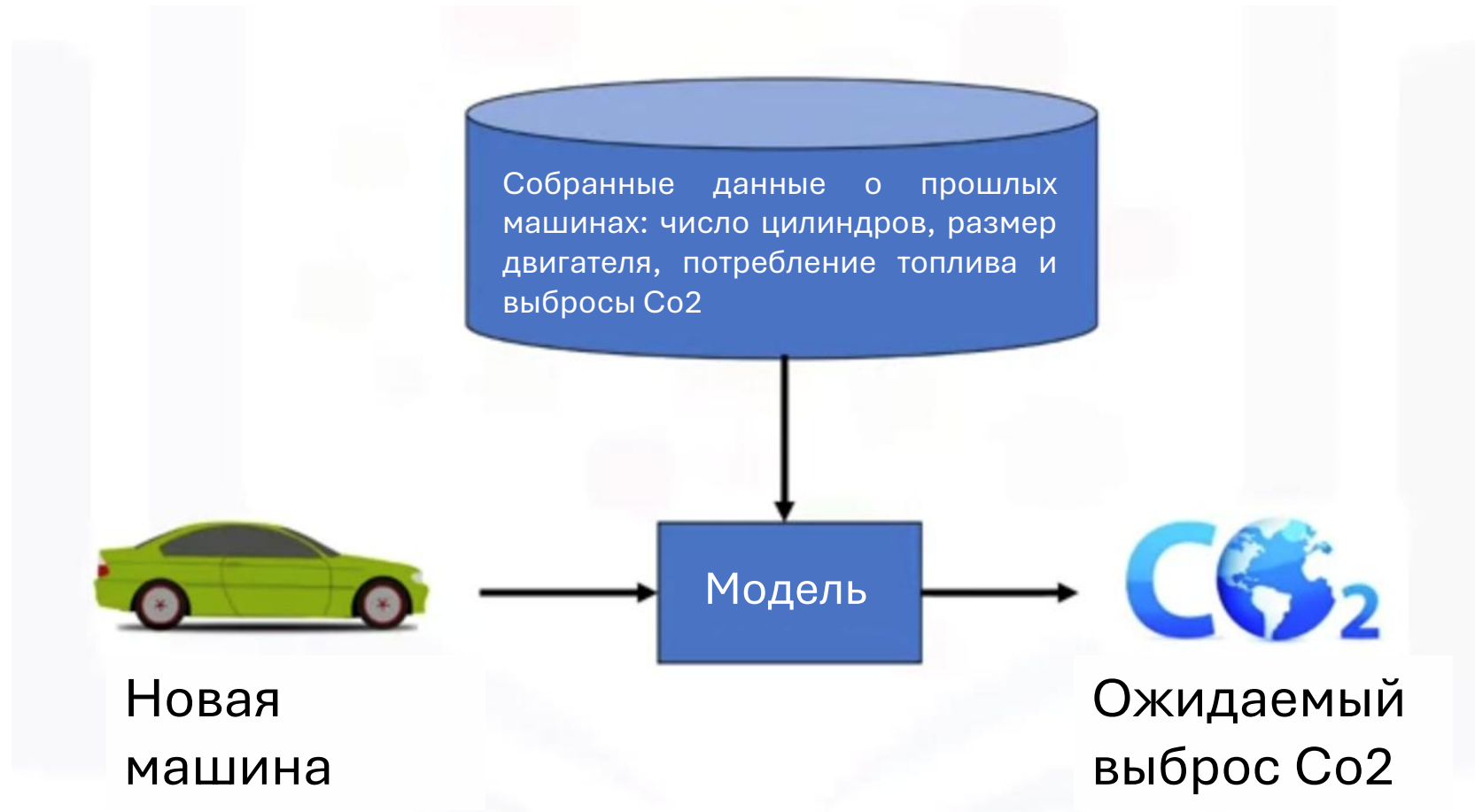
Введение

Simple Linear Regression


Практическое применение

Оценка модели и диагностика

Регрессия



Регрессия



	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Регрессия

- Задача предсказания непрерывного числа

X: Independent variable

Y: Dependent variable

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Типы регрессии

Простая регрессия

- ❑ Использование одной независимой переменной для прогнозирования зависимой переменной.
- ❑ Например: прогнозирование выбросов CO₂ на основе размера двигателя.
- ❑ Простая линейная регрессия, простая нелинейная регрессия

Множественная регрессия

- ❑ Использование более чем одной независимой переменной для прогнозирования зависимой переменной.
- ❑ Например: прогнозирование выбросов CO₂ на основе размера двигателя и количества цилиндров.
- ❑ Множественная линейная регрессия, множественная нелинейная регрессия

Применение регрессии

- Прогнозирование годового объема продаж человека на основе возраста, стажа работы и т. д
- Определение индивидуальной удовлетворенности на основе демографических и психологических факторов
- Прогнозирование цены дома на основе его размера, количества комнат и т. д
- Прогнозирование дохода от работы с учетом независимых переменных, таких как количество рабочих часов, образование, профессия, пол, возраст, стаж работы



Вопрос

Какой из примеров является примером применения регрессионного анализа?


- Прогнозирование наличия или отсутствия рака у пациента.
- Группировка похожих домов в районе.
- Прогнозирование количества осадков на следующий день.
- Прогнозирование победы или поражения команды.

Алгоритмы для решения задачи регрессии

- ☐ Порядковая регрессия
- ☐ Регрессия Пуассона
- ✓ Линейная регрессия (Простая и множественная)
- ☐ Полиномиальная регрессия
- ☐ Регрессия Лассо
- ☐ Гребневая регрессия
- ☐ Регрессия на основе дерева решений
- ☐ Регрессия на основе бустированного дерева решений

Простая линейная регрессия

Линейная регрессия — это аппроксимация линейной модели, используемая для описания взаимосвязи между двумя или более переменными.

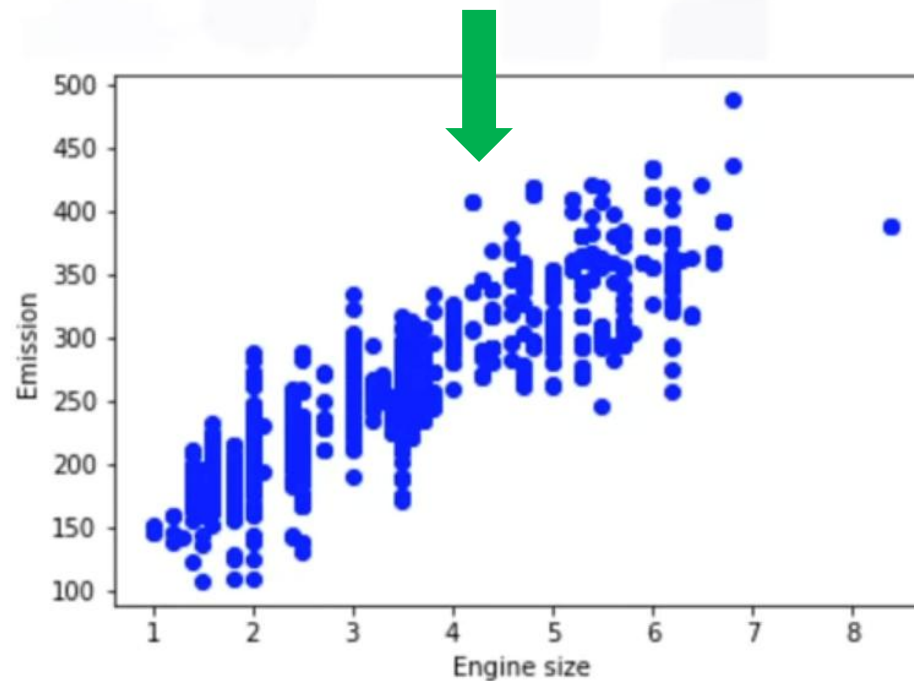


	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Простая линейная регрессия

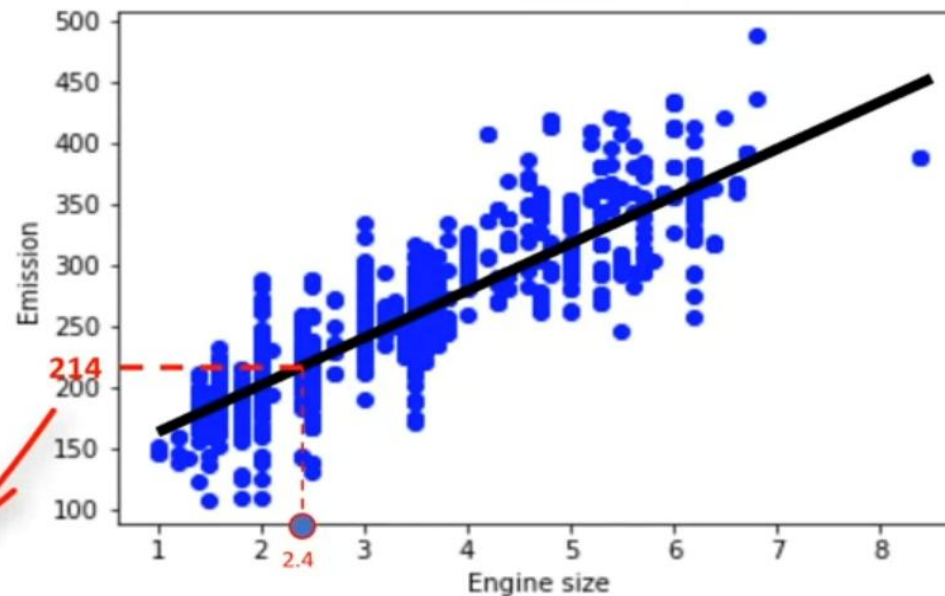
Изменение в одной
переменной объясняет
изменение в другой

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?



Простая линейная регрессия

	ENGINE SIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?



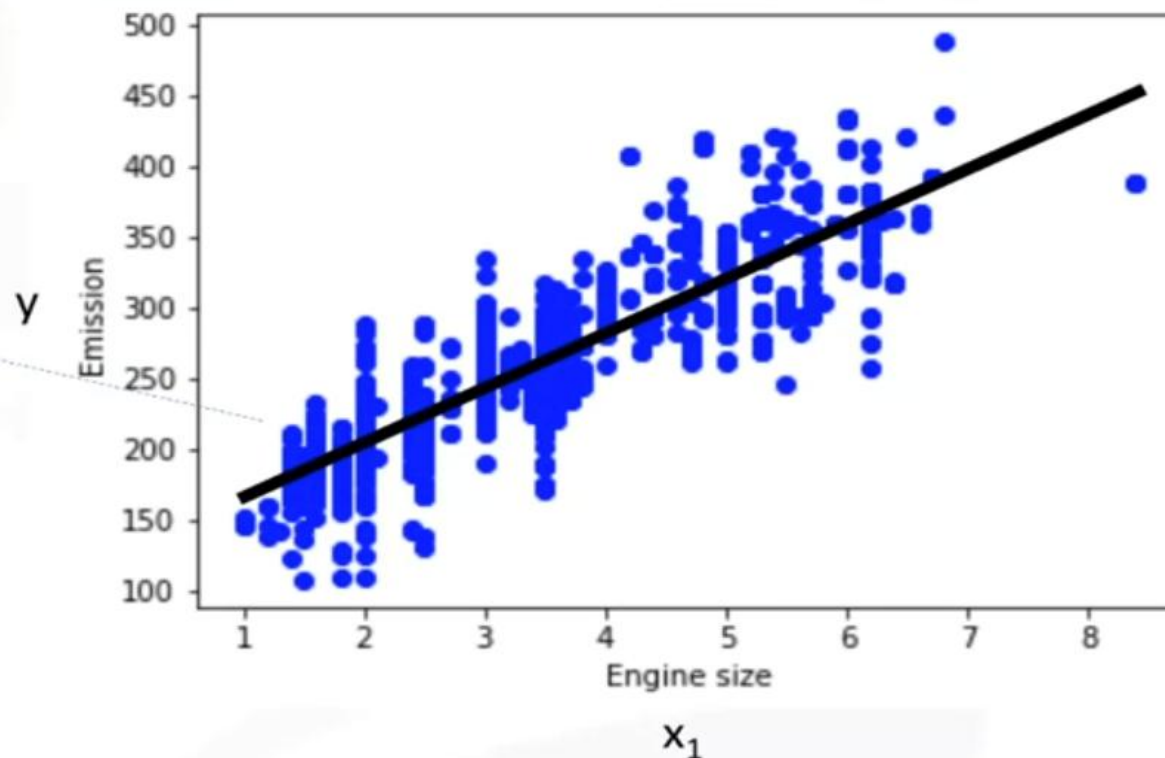
Простая линейная регрессия

Параметры, которые нужно найти

$$\hat{y} = \theta_0 + \theta_1 x_1$$

Независимая
переменная

Зависимая
переменная



Как найти оптимальный вариант?

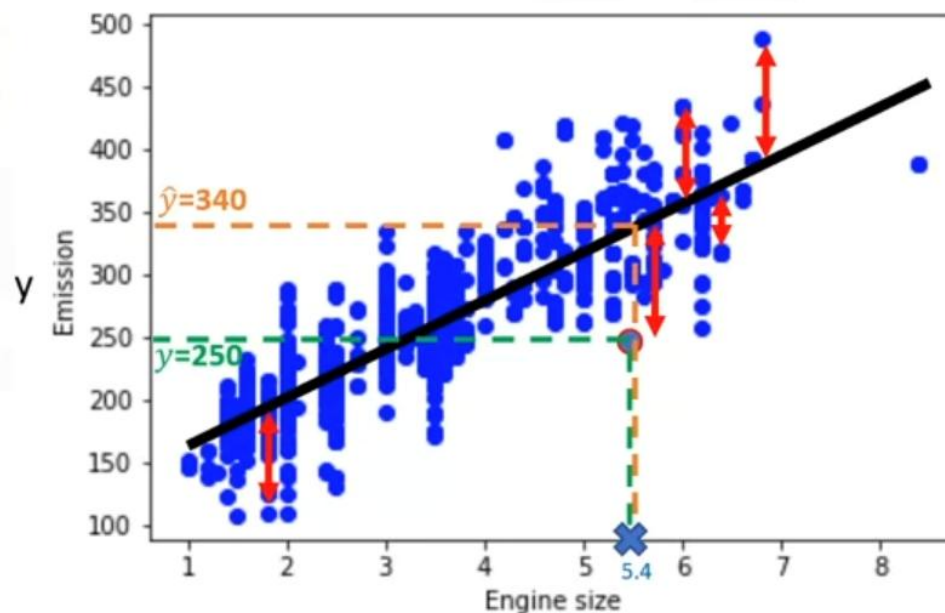
$x_1 = 5.4$ независимая переменная
 $y = 250$ истинный выброс Co2 для x_1

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$\hat{y} = 340$ предсказанный выброс Co2 для x_1

$$\begin{aligned}\text{Error} &= y - \hat{y} \\ &= 250 - 340 \\ &= -90\end{aligned}$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



Математический подход

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267

X_1 is indicated by a bracket on the left side of the table, spanning rows 4 through 8.

y is indicated by a bracket on the right side of the table, spanning rows 4 through 8.

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$\theta_1 = \frac{\sum_{i=1}^s (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^s (x_i - \bar{x})^2}$$

$$\bar{x} = (2.0 + 2.4 + 1.5 + \dots) / 9 = 3.03$$

$$\bar{y} = (196 + 221 + 136 + \dots) / 9 = 226.22$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

Минимизируя сумму квадратов ошибок (SSE)

Предсказание

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$Co2Emission = \theta_0 + \theta_1 EngineSize$$

$$Co2Emission = 125 + 39 EngineSize$$

$$Co2Emission = 125 + 39 \times 2.4$$

$$Co2Emission = 218.6$$



Преимущества

- Быстро
- Не нужна настройка гиперпараметров
- Интерпретируемая

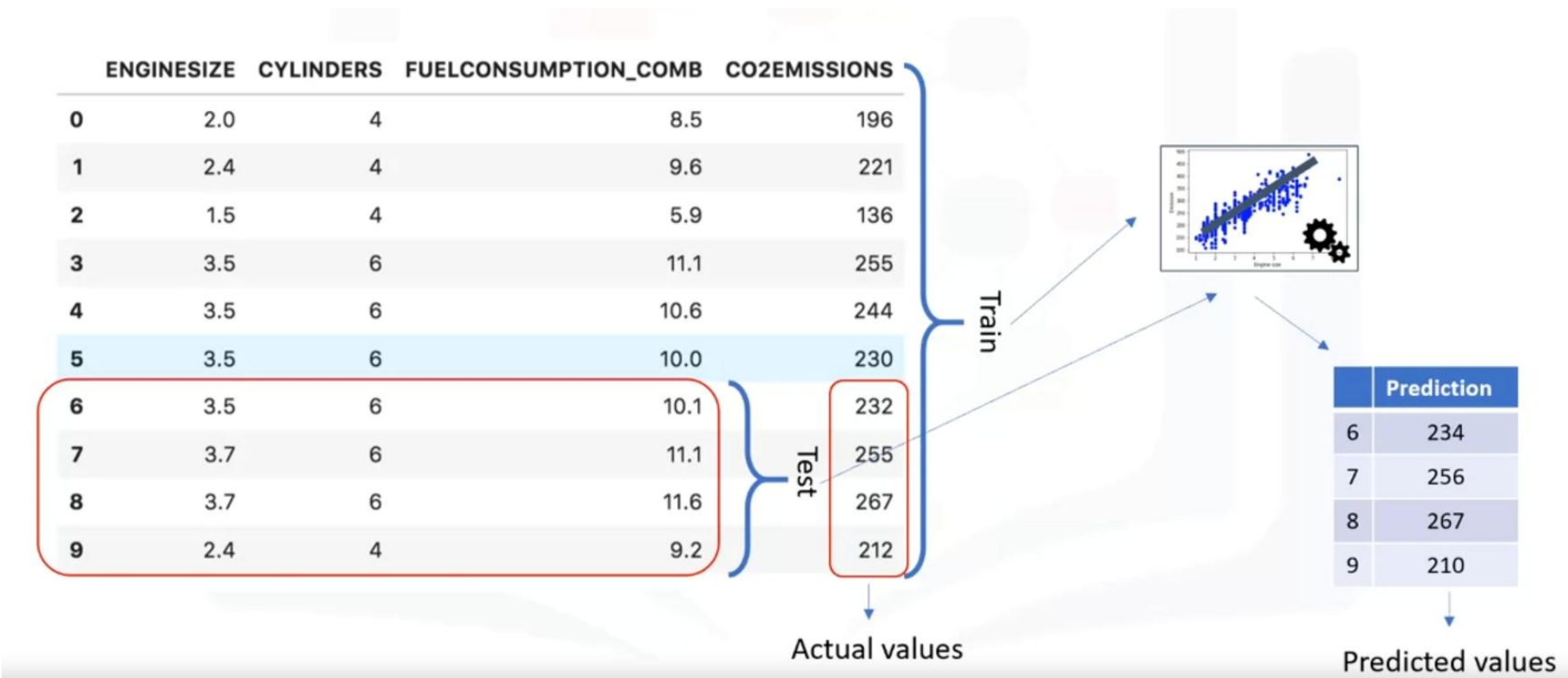


Disadvantages

- Чувствительная к выбросам
- Не поможет при нелинейной связи

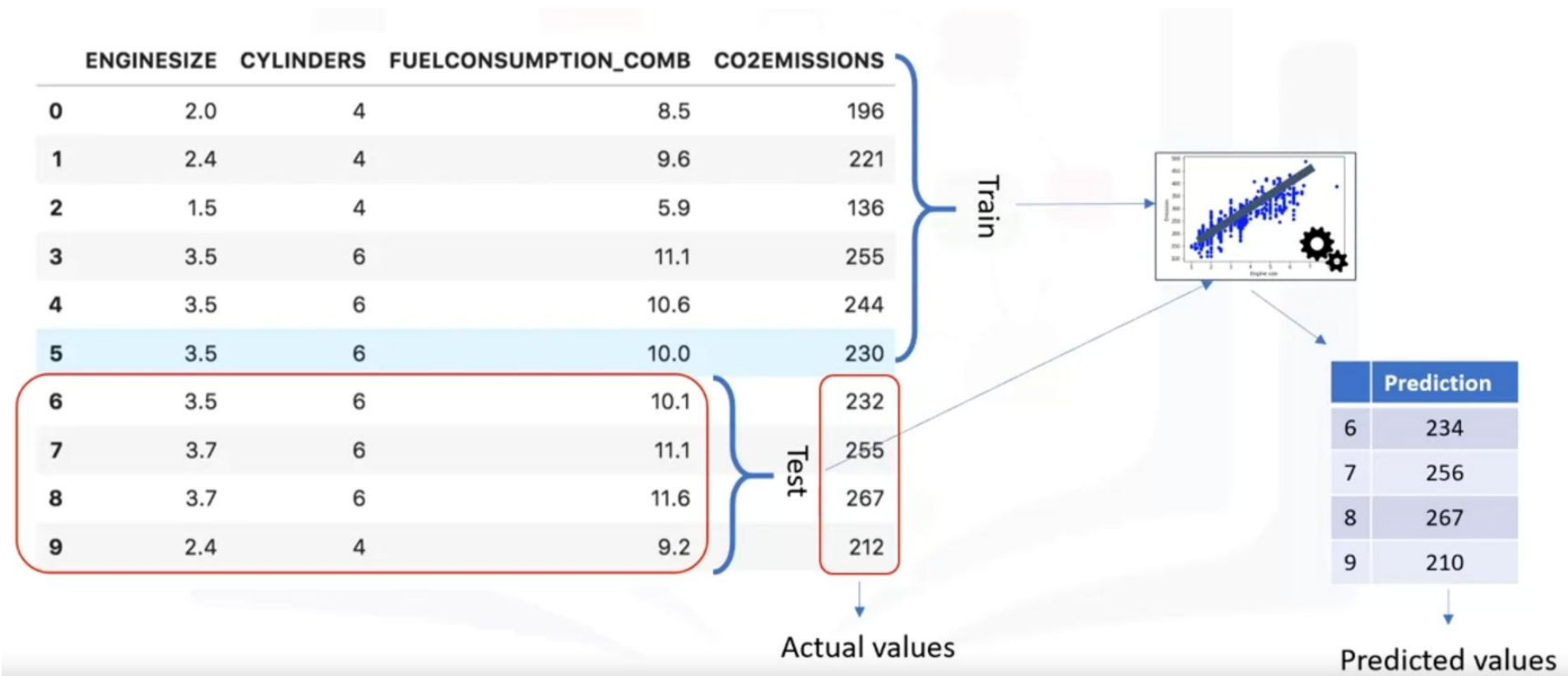
Как оценить точность (Accuracy) модели

Train/Test split



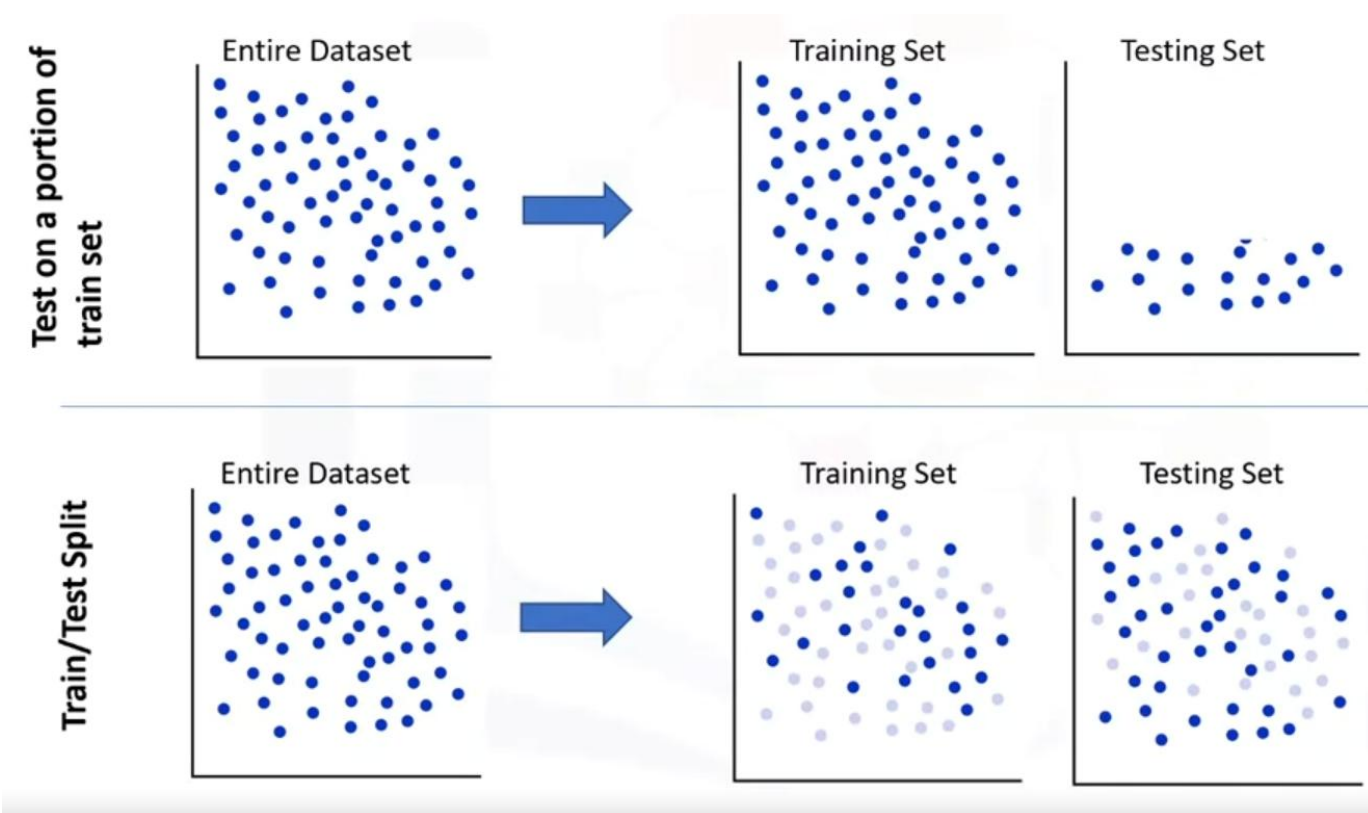
Как оценить точность (Accuracy) модели

Train/Test split



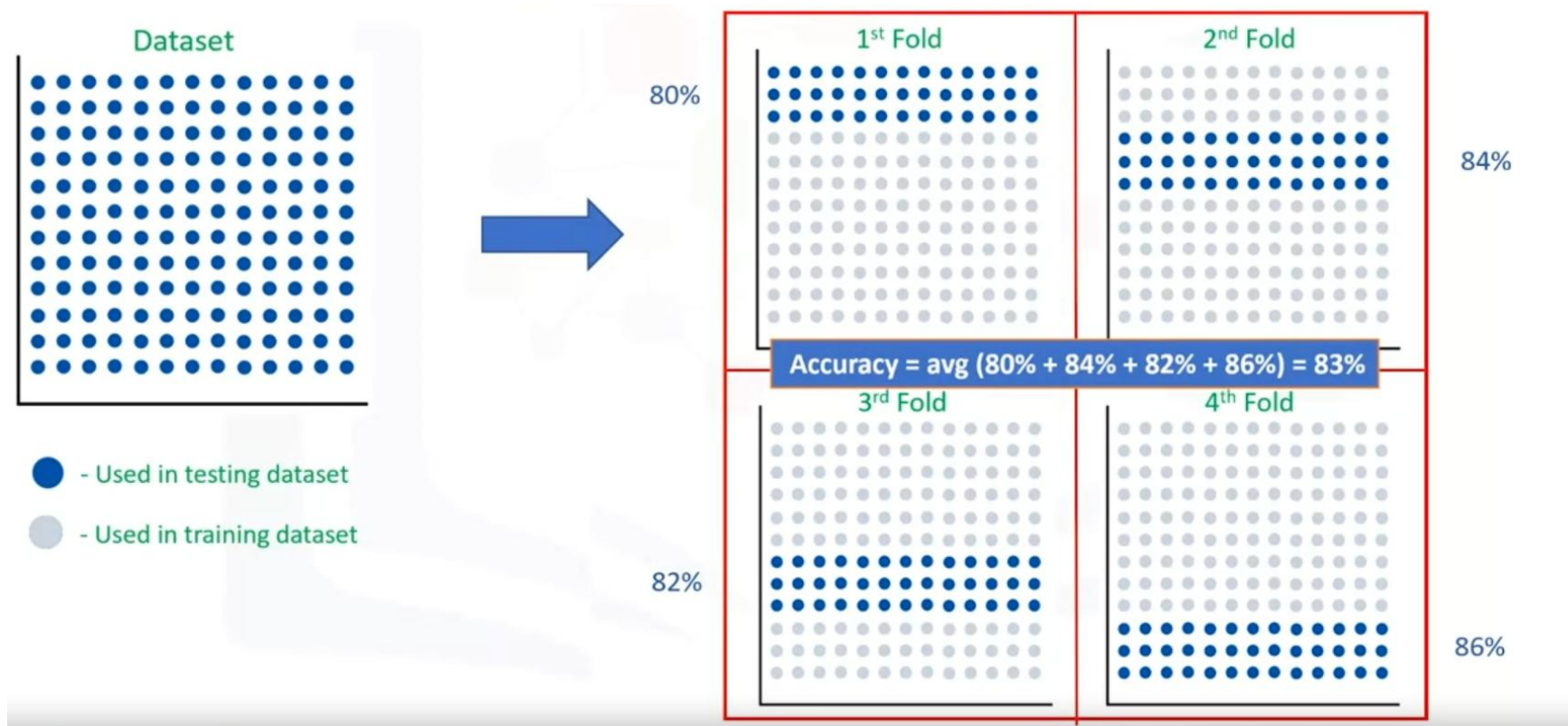
Как оценить точность (Accuracy) модели

Train/Test split

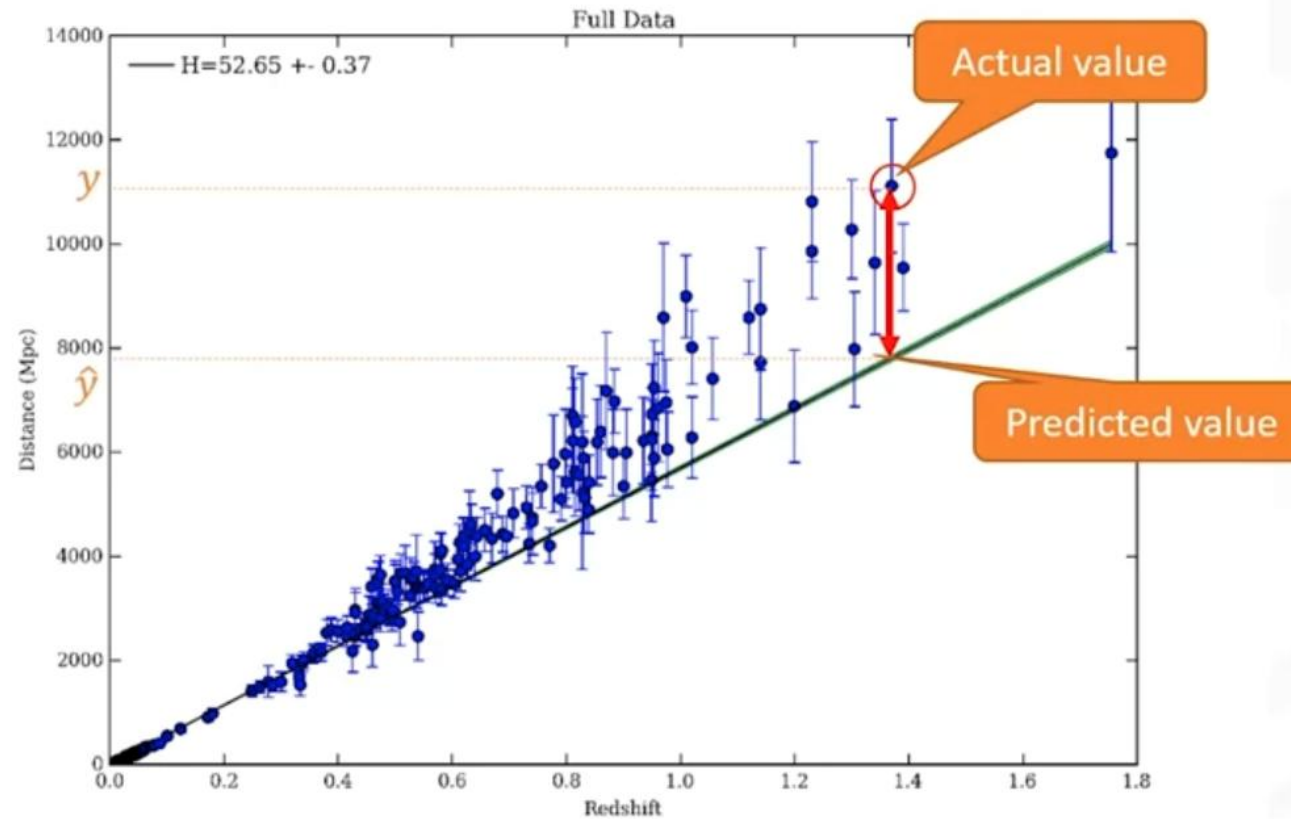


- Тестовая выборка часть обучающей
- Высокая точность на обученных данных
- Низкая точность на новых данных
- Разделение на независимые подвыборки
- Более высокая точность на новых данных
- Высоко зависит от того, на каких наборах данных обучалась и оценивалась

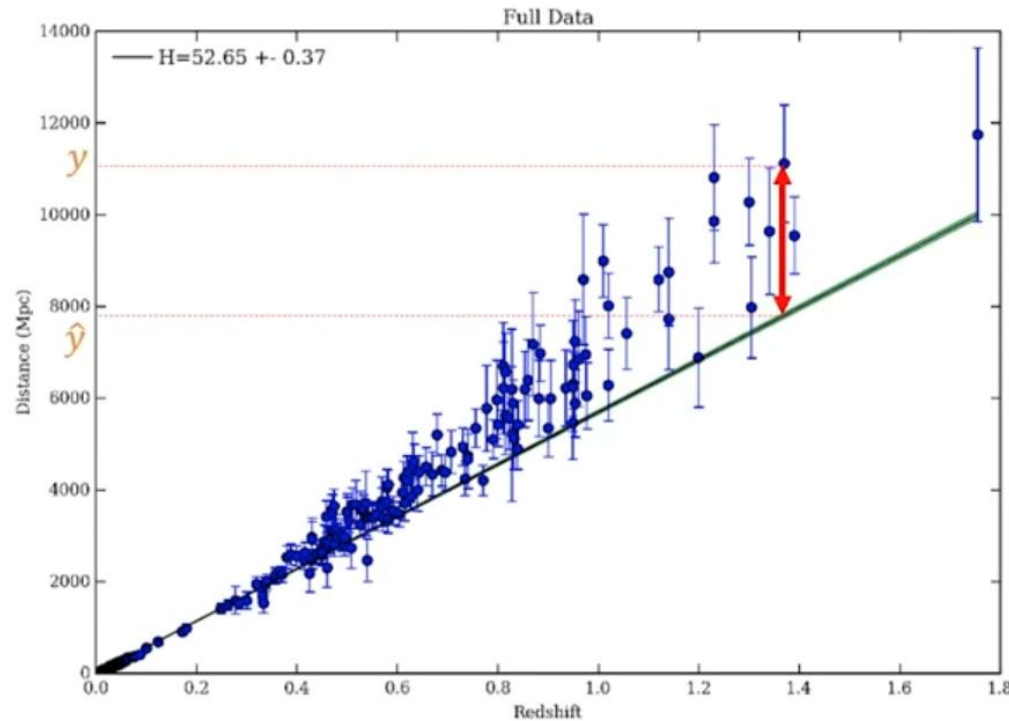
Кросс-валидация (k-fold cross-validation)



Оценка метрик в модели регрессии



Оценка метрик в модели регрессии



$$MAE = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

$$RAE = \frac{\sum_{j=1}^n |y_j - \hat{y}_j|}{\sum_{j=1}^n |y_j - \bar{y}|}$$

$$RSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

$$R^2 = 1 - RSE$$

Множественная линейная регрессия

$$Co2\ Em = \theta_0 + \theta_1 Engine\ size + \theta_2 Cylinders + \dots$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\hat{y} = \theta^T X$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots]$$
$$X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$

	X: Independent variable			Y: Dependent variable
	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Как найти лучшие параметры?

$$\hat{y} = \theta^T X$$

$$\hat{y}_i = 140$$

предсказанный x_i

$$y_i = 196$$

истинный x_i

$$y_i - \hat{y}_i = 196 - 140 = 56 \text{ остаточная ошибка}$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267

Как найти лучшие параметры?

- ❑ Математический подход: операции линейной алгебры (для небольших наборов данных)
- ❑ Оптимизационный подход: градиентный спуск (для больших наборов данных)

Множественная линейная регрессия

$$Co2\ Em = \theta_0 + \theta_1 Engine\ size + \theta_2 Cylinders + \dots$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\hat{y} = \theta^T X$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T y$$

X: Independent variable Y: Dependent variable

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Множественная линейная регрессия

$$Co2\ Em = \theta_0 + \theta_1 Engine\ size + \theta_2 Cylinders + \dots$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\hat{y} = \theta^T X$$

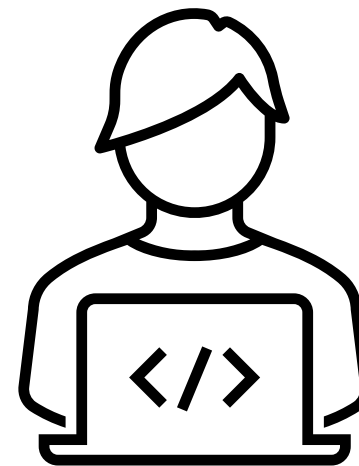
$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T y$$

Смещение (Bias): Чтобы модель могла предсказывать значения не из начала координат, мы добавляем фиктивный столбец из единиц к матрице признаков.

Метрика R^2 : Она показывает долю объясненной дисперсии. Значение 1.0 — идеальное предсказание, 0.0 — модель работает не лучше, чем простое предсказание среднего значения.

Практика



28

Семинар 2:

<https://colab.research.google.com/drive/1CIbOl4-sHVzBL8C4g9qEFI2OOIoWxsi-?usp=sharing>

Лабораторная работа 2:

<https://colab.research.google.com/drive/15x9TOZsFSvbrF-0bxKw8Wwk75N3voAM8?usp=sharing>

Датасет для lab02:

<https://www.kaggle.com/datasets/prokshitha/home-value-insights>