



Математическая статистика

и основы анализа данных

Основы машинного обучения



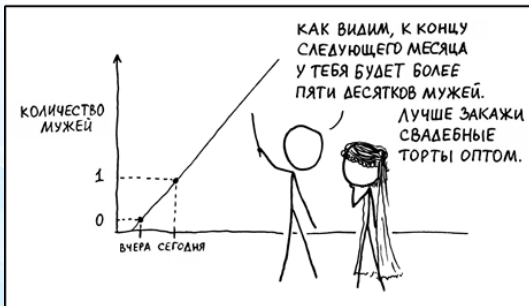
План занятия

- ▶ Задача регрессии
 - ▶ Линейная регрессия
 - ▶ Метрики качества для задачи регрессии
- ▶ Задача классификации
 - ▶ Классификатор kNN
 - ▶ Метрики качества для задачи классификации

Линейная регрессия

Метод наименьших квадратов

МОЁ ХОББИ: ЭКСТРАПОЛИРОВАТЬ





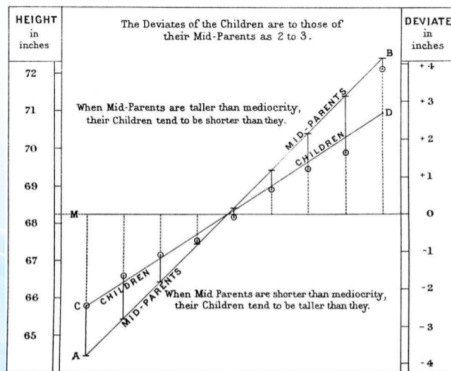
Первое упоминание регрессии

Впервые регрессия упоминается в работе Гальтона

"Регрессия к середине в наследственности роста", 1885 г.

x — рост родителей, y — рост детей

Установлена зависимость $y - \bar{y} \approx \frac{2}{3}(x - \bar{x})$, т.е. регрессия к середине.





Задача регрессии: интуиция

Есть объект, обладающий признаками x .

Примеры признаков: рост песика, экспрессия белка, энергия частицы.

Мы предполагаем, что есть зависимость какой-то численной характеристики объекта y от его признаков:

$$y \approx f(x)$$

Пример: зависимость уровня когнитивных способностей от параметров поражения мозга при рассеянном склерозе.

Однако мы не знаем, какова эта зависимость на самом деле.

На основании **данных** – набора объектов, для которых известны x и y , мы пытаемся "восстановить" зависимость:

$$y \approx \hat{f}(x)$$



Пример

Пусть x — рост песика, а y — его вес.

Что мы знаем?

- ▶ чем крупнее песик, тем больший вес он имеет;
- ▶ песики одинакового роста могут иметь разный вес.

Выводы:

- ▶ для фиксированного роста песика x его вес $y = f(x)$ является случайной величиной;
- ▶ в среднем вес $f(x)$ возрастает при увеличении роста песика x .



Пример

Простая зависимость:

$$y = \theta_0 + \theta_1 x + \varepsilon,$$

x — рост песика,

y — вес песика,

θ_0, θ_1 — неизвестные параметры,

ε — случайная составляющая с нулевым средним (погрешность).

Зависимость **линейна по параметрам**, линейна по аргументу.



Пример

Более сложная зависимость:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_2^2 + \varepsilon,$$

x_1 — рост песика,

x_2 — обхват туловища песика,

y — вес песика,

$\theta_0, \theta_1, \theta_2, \theta_3$ — неизвестные параметры,

ε — случайная составляющая с нулевым средним.

Зависимость **линейна по параметрам**, квадратична по аргументам.



Модель линейной регрессии

Рассматриваем функциональную зависимость вида

$$y = y(x) = \theta_1 x_1 + \dots + \theta_d x_d$$

x_1, \dots, x_d — признаки ,

$\theta = (\theta_1, \dots, \theta_d)^T$ — вектор параметров.

Для оценки θ производится n испытаний вида

$$Y_i = \theta_1 x_{i1} + \dots + \theta_d x_{id} + \varepsilon_i, \quad i = 1, \dots, n,$$

$x_i = (x_{i1}, \dots, x_{id})$ — признаковые описания объекта i
(обычно неслучайные),

ε_i — случайная ошибка измерений.



Модель линейной регрессии

Введем обозначения

$$Y = \begin{pmatrix} Y_1 \\ \dots \\ Y_n \end{pmatrix}, \quad X = \begin{pmatrix} x_{11} & \dots & x_{1d} \\ \dots & & \\ x_{n1} & \dots & x_{nd} \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \dots \\ \varepsilon_n \end{pmatrix}.$$

Матричная форма записи проведенных испытаний

$$Y = X\theta + \varepsilon.$$

$X \in \mathbb{R}^{n \times d}$ — регрессоры (или матрица плана эксперимента),

$Y \in \mathbb{R}^n$ — отклик.

Матричный вид зависимости: $y(x) = x^T \theta$.



Замечание

Зависимость $y = y(x)$ должна быть **линейна по параметрам**, но не обязана быть линейной по признакам.

Пусть z_1, \dots, z_k — набор "независимых" переменных.

Можно рассматривать модель

$$y(x) = \theta_1 x_1(z_1, \dots, z_k) + \dots + \theta_d x_d(z_1, \dots, z_k),$$

где $x_j(z_1, \dots, z_k)$ — некоторые функции (м.б. нелинейные).

Примеры:

▶ $x(z_1, \dots, z_k) = 1;$

▶ $x(z_1, \dots, z_k) = z_1;$

▶ $x(z_1, \dots, z_k) = \ln z_1;$

▶ $x(z_1, \dots, z_k) = z_1^2 z_2.$

Определение моментов инерции твёрдых тел с помощью трифилярного подвеса

- ▶ На платформу помещается тело — диск, разрезанный по диаметру;
 - ▶ I — момент инерции тела;
 - ▶ m — масса тела;
 - ▶ h — расстояние от половинок до оси вращения;
 - ▶ I_0 — момент инерции нераздвинутого диска.
- ▶ Половинки диска постепенно раздвигаются;
- ▶ Снимается зависимость момента инерции системы I от h .

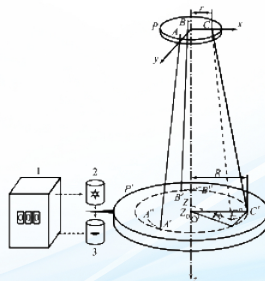


Рис. 2. Трифилярный подвес

По материалам "Модели и концепции физики: механика. Лабораторный практикум"



Пример: Момент инерции

Согласно теореме Гюйгенса-Штейнера должно выполняться:

$$I = I_0 + mh^2$$

Итого, предполагается линейная зависимость момента инерции I от квадрата расстояния h^2 . Мы хотим найти неизвестные m и I_0 .

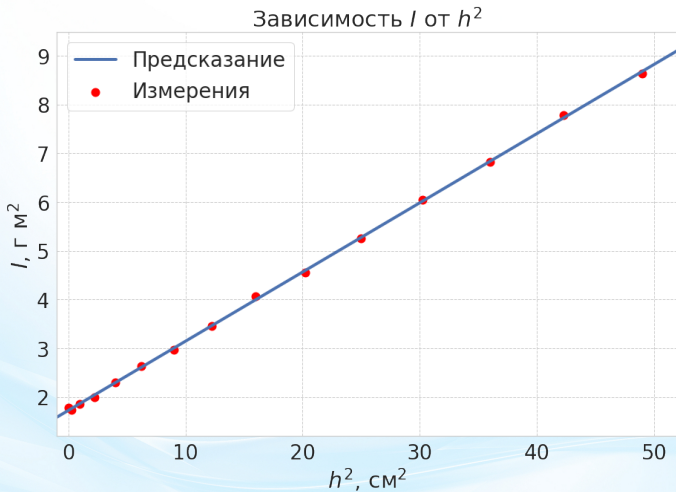
Наблюдения: $I_i = I_0 + mh_i^2 + \varepsilon_i$, где ε_i — погрешность.

В данном примере $x_1(t) = 1$, $x_2(h) = h^2$,

$$X = \begin{pmatrix} 1 & h_1^2 \\ \dots & \dots \\ 1 & h_n^2 \end{pmatrix}, Y = \begin{pmatrix} I_1 \\ \dots \\ I_n \end{pmatrix}, \theta = \begin{pmatrix} I_0 \\ m \end{pmatrix}.$$

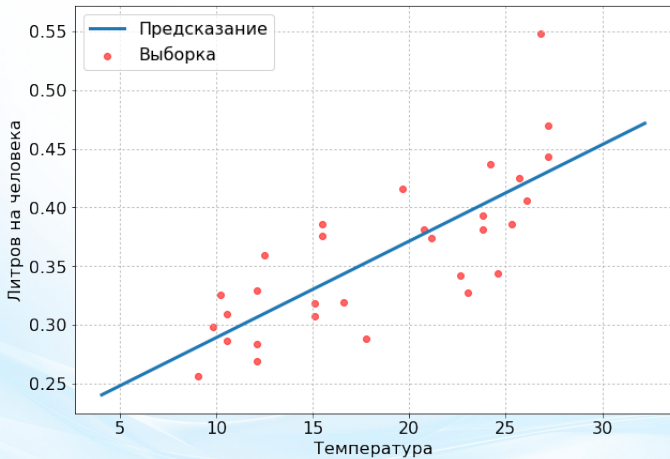


Пример: Момент инерции





Пример: Потребление мороженого





Метод наименьших квадратов

Зависимость: $y(x) = x^T \theta$, $\theta \in \mathbb{R}^d$.

Испытания: $Y = X\theta + \varepsilon$, $X \in \mathbb{R}^{n \times d}$, $Y \in \mathbb{R}^n$.

Хотим как-то **оценить** параметр θ на основании полученных данных.

Пусть $\hat{\theta} = \hat{\theta}(X, Y)$ — наша оценка θ .

Как понять, что она хорошая?

Метрика MSE:

$$MSE(\hat{\theta}) = \left\| Y - X\hat{\theta} \right\|^2$$

Оценка $\hat{\theta} = \arg \min_{\theta} MSE(\hat{\theta})$ называется **оценкой по методу наименьших квадратов** параметра θ .



Метод наименьших квадратов

Теорема. Если матрица $X^T X$ невырождена, то $\hat{\theta} = (X^T X)^{-1} X^T Y$.

$$MSE(\theta) = \|Y - X\theta\|^2 = (Y - X\theta)^T (Y - X\theta) = Y^T Y - 2Y^T X\theta + \theta^T X^T X\theta$$

Берем производную по θ и приравниваем ее к нулю.

$$\frac{\partial MSE(\theta)}{\partial \theta} = -2Y^T X + 2\theta^T X^T X = 0$$

Отсюда получается утверждение теоремы. □

Предсказанием отклика на новом объекте x будет величина $\hat{y}(x) = x^T \hat{\theta}$.



Реализация в sklearn

```
m = sklearn.linear_model.LinearRegression(fit_intercept=True)
```

Обучение модели:

```
m.fit(X, Y)
```

Вектор коэффициентов:

```
m.coef_
```

Свободный коэффициент:

```
m.intercept_
```

Предсказания:

```
m.predict(X)
```



Метрики качества в задаче регрессии



Обозначения

Пусть x_1, \dots, x_n — признаковые описания объектов;

$Y = (Y_1, \dots, Y_n)^T$ — наблюдения.

Пусть $\hat{f}(x)$ — оцененная нами зависимость.

В случае линейной регрессии $\hat{f}(x) = x^T \hat{\theta}$.

Пусть $\hat{Y}_i = \hat{f}(x_i)$ — предсказание нашей модели на i -м объекте;

$\hat{Y} = (\hat{Y}_1, \dots, \hat{Y}_n)^T$.

Метрики качества в задаче регрессии

Y — реальные наблюдения, \hat{Y} — предсказания.

- ▶ MSE (Mean Squared Error):

$$MSE(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

- ▶ MAE (Mean Absolute Error):

$$MAE(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^n |Y_i - \hat{Y}_i|$$

- ▶ MAPE (Mean Absolute Percentage Error):

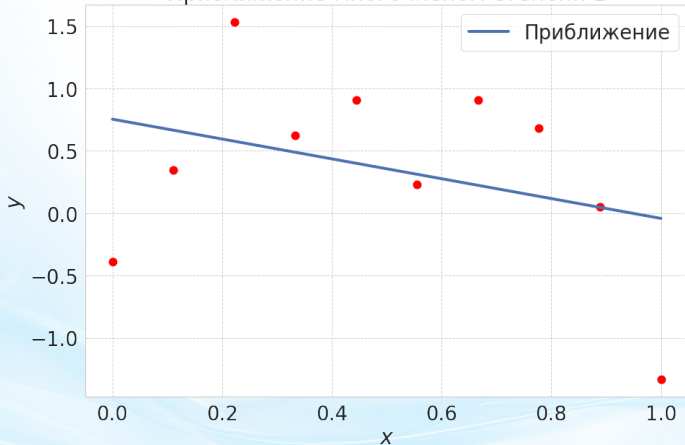
$$MAPE(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^n \left| \frac{Y_i - \hat{Y}_i}{Y_i} \right| * 100\%$$



Недообучение vs Переобучение

Зависимость: $y = 5x - 6x^2$, имеется погрешность

Приближение многочленом степени 1



Недообучение



Недообучение vs Переобучение

Зависимость: $y = 5x - 6x^2$, имеется погрешность

Приближение многочленом степени 10



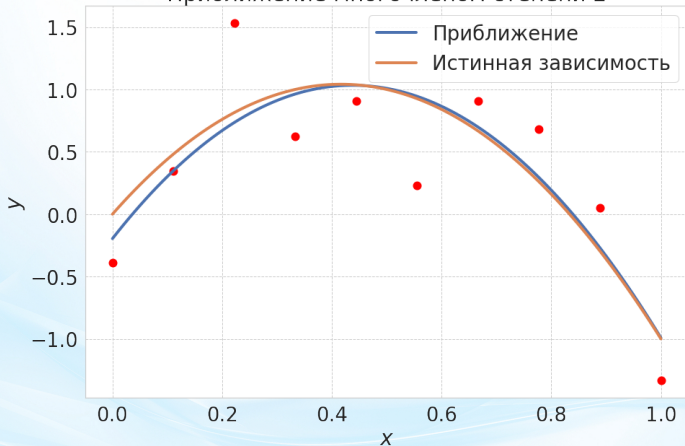
Переобучение



Недообучение vs Переобучение

Зависимость: $y = 5x - 6x^2$, имеется погрешность

Приближение многочленом степени 2



Нормально!



Тренировочная и тестовая выборки

Если все время работать с одной и той же выборкой (это жаргон, корректно понимать "реализацией выборки") и все больше улучшать модель, "подгонять" ее под выборку, может возникнуть переобучение.

Предсказание на **новом** объекте может быть неадекватным.

Поэтому перед началом работы имеющиеся данные делят на две части:

тренировочную (обучающую) и **тестовую** выборки.



На тренировочной выборке происходит **обучение** моделей (например, оценка коэффициентов в линейной регрессии).

На тестовой выборке происходит **оценка качества** итоговой модели с использованием метрик качества.



Задача классификации

Классификатор KNN



Задача классификации: интуиция

У нас есть объекты, заданные своими признаковыми описаниями x .

Пример: рост, вес, уровни экспрессии генов, фотография куска неба.

У каждого объекта есть **класс** y — значение из конечного множества классов C . *Пример: порода пса, тип клетки, индикатор наличия звезды.*

На основании **данных** — объектов, для которых известны x и y , хотим построить зависимость $y \approx \hat{f}(x)$.

Тогда мы сможем предсказывать класс на новом объекте по его признаковому описанию x : $\hat{y} = \hat{f}(x)$

Бинарная классификация: классов всего два.

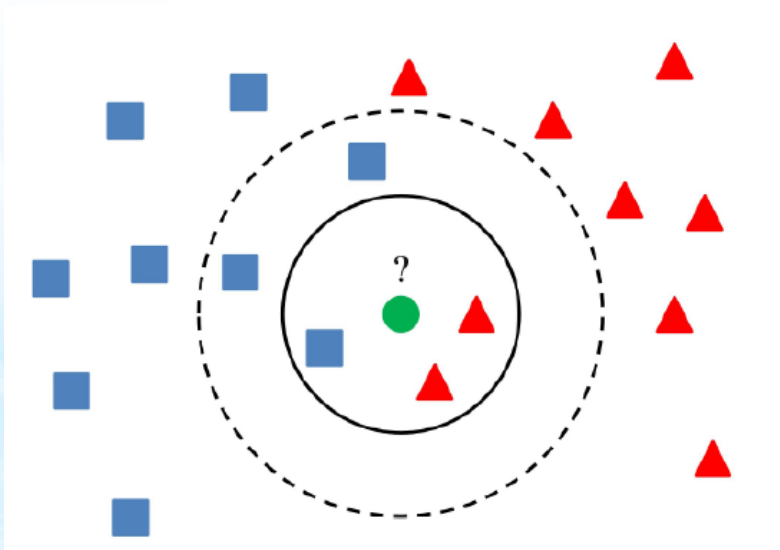
Пример: является ли клетка раковой.

Многоклассовая классификация: несколько классов.

Пример: тип клетки (нейрон, глия, эпителий, etc).



Метод ближайших соседей (kNN)





Метод ближайших соседей (kNN)

Пусть \mathcal{X} — метрическое пространство.

$x_1, \dots, x_n \in \mathcal{X}$ — обучающая выборка (признаки объектов).

Y_1, \dots, Y_n — соответствующая целевая переменная (например, класс).

Предположение:

Свойства объекта меняются не сильно в его окрестности.

Тогда давайте смотреть на свойства k ближайших соседей, где k — параметр.

Пусть $x \in \mathcal{X}$ — исследуемый объект.

$x_{(1)}, \dots, x_{(k)}$ — k его соседей в порядке удаления от x .

1. Классификация.

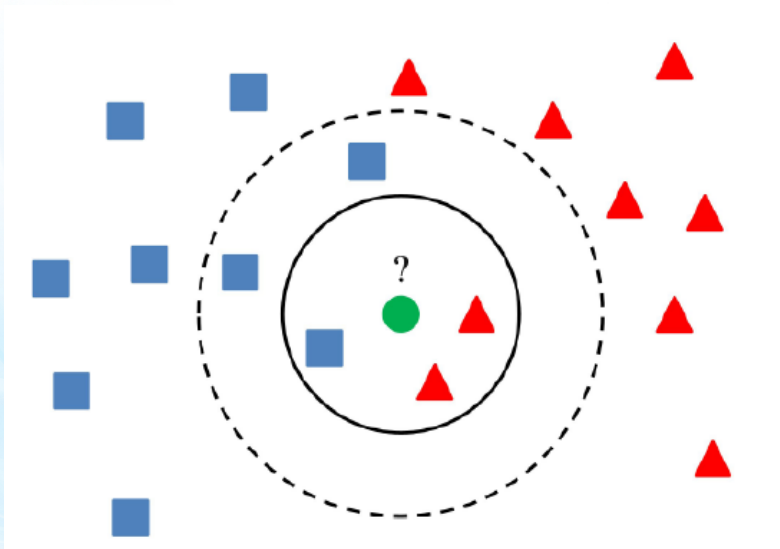
Предсказание — наиболее часто встречаемый класс.

2. Регрессия.

Предсказание — усреднение отклика по соседям.



А одинаково ли важны соседи?





Взвешенный kNN для задачи классификации

Пусть $x \in \mathcal{X}$ — исследуемый объект.

$x_{(1)}, \dots, x_{(k)}$ — k его соседей в порядке удаления от x .

$Y_{(1)}, \dots, Y_{(k)}$ — соответствующие классы, C — множество классов.

w_1, \dots, w_k — **вес/вклад** j -го соседа, определяемый пользователем.

Способы определения веса:

1. $w_j = 1 - j/k$ — зависящий от **номера** соседа;
2. $w_j = \|x - x_{(j)}\|^{-1}$ — зависящий от **расстояния** до соседа.

Тогда "оптимальный" класс в задаче классификации можно определить так:

$$\hat{y}(x) = \arg \max_{c \in C} \sum_{j=1}^k w_j I\{Y_{(j)} = c\}$$

Смысл формулы:

- ▶ Складываем веса соседей для каждого класса;
- ▶ Берем класс с наибольшим суммарным весом.



Метрика accuracy в задаче классификации

Данная метрика является самой тривиальной метрикой для задачи классификации и представляет собой **долю правильных ответов**.

Y_1, \dots, Y_n — истинные классы объектов.

$\hat{Y}_1, \dots, \hat{Y}_n$ — предсказанные классы объектов, т.е. $\hat{Y}_i = \hat{f}(x_i)$.

$$\text{accuracy}(Y, \hat{Y}) = \frac{1}{n} \sum_{i=1}^n I\{Y_i = \hat{Y}_i\}$$



Применения kNN на практике

kNN biological data



Все



Картинки



Новости



Видео



Карты



Ещё

Настройки

Инструменты

Результатов: примерно 1 130 000 (0,49 сек.)

link.springer.com › chapter · [Перевести эту страницу](#)

A KNN-Based Learning Method for Biology Species ...

KNN has been successfully used in natural language processing (NLP). Our work extends the learning method for **biological data**. We view the DNA or RNA ...

ieeexplore.ieee.org › document · [Перевести эту страницу](#)

Applied biological data mining based on improved K-means ...

Applied **biological data** mining based on improved K-means clustering algorithm and KNN classifier in **protein sub-cellular localization**. ... The experimental results based on protein sub-cellular localization prediction show that the methods proposed newly better work than the traditional methods.

www.nature.com › tpj201056 · [Перевести эту страницу](#)

k -Nearest neighbor models for microarray gene expression ...

30 июл. 2010 г. — The kDAP produces consistent KNN prediction models on a newly generated **data** set created by a different microarray technology. The resulting KNN model parameters reveal the underlying **biological** and practical characteristics of the end points.





Применения kNN на практике

A KNN-Based Learning Method for Biology Species Categorization

Dang et al., https://link.springer.com/chapter/10.1007/11539087_127

- ▶ Категоризация видов;
- ▶ Работают с ДНК и РНК разных видов;
- ▶ Строят признаки по последовательности ДНК/РНК методами анализа текстов (след. семестр);
- ▶ Для полученных признаков строят классификатор kNN;
- ▶ Категоризовали 43 бактерии.



Применения kNN на практике

Applied biological data mining based on improved K-means clustering algorithm and KNN classifier in protein sub-cellular localization

Zhenfeng Lei, Shunfang Wang, <https://ieeexplore.ieee.org/document/7883060>

- ▶ Локализация белков в клетках;
- ▶ Строят признаковое описание белков;
- ▶ Придумали различные улучшения kNN;

k-Nearest neighbor models for microarray gene expression analysis and clinical outcome prediction

Parry et al., <https://www.nature.com/articles/tpj201056>

- ▶ Датасеты с тремя типами рака;
breast cancer, neuroblastoma, multiple myeloma
- ▶ Экспериментировали с признаками, метриками, etc;
- ▶ Построили 463320 моделей KNN...



Применения kNN на практике

A Regression-based K nearest neighbor algorithm for gene function prediction from heterogeneous data

Zizhen Yao, Walter L Ruzzo,

<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-7-S1-S11>

- ▶ Предсказывают функцию генов;
- ▶ Оптимизируют метрику с помощью регрессии;
- ▶ Метрика — "взвешенная комбинация базовых метрик" (напоминает линейную регрессию!);
- ▶ Предсказывают функции генов *Escherichia coli* и сравнивают с известными.



BCĚ!