

Логистическая регрессия

[Все курсы](#) > [Оптимизация](#) > [Занятие 5](#)

Как мы уже знаем, несмотря на название, логистическая регрессия решает задачу классификации. Сегодня мы подробно разберем принцип работы и составные части алгоритма логистической регрессии, а также построим модели с одной и несколькими независимыми переменными.

Содержание занятия

- Бинарная логистическая регрессия
 - Задача бинарной классификации
 - Функция логистической регрессии
 - Сигмоида
 - Интерпретация коэффициентов
 - Решающая граница
 - Вычислительная устойчивость сигмоиды
 - Logistic loss или функция кросс-энтропии
 - График и формула функции логистической ошибки
 - Расчет логистической ошибки
 - Окончательный вариант
 - Производная функции логистической ошибки
 - Производная логарифмической функции
 - Производная сигмоиды

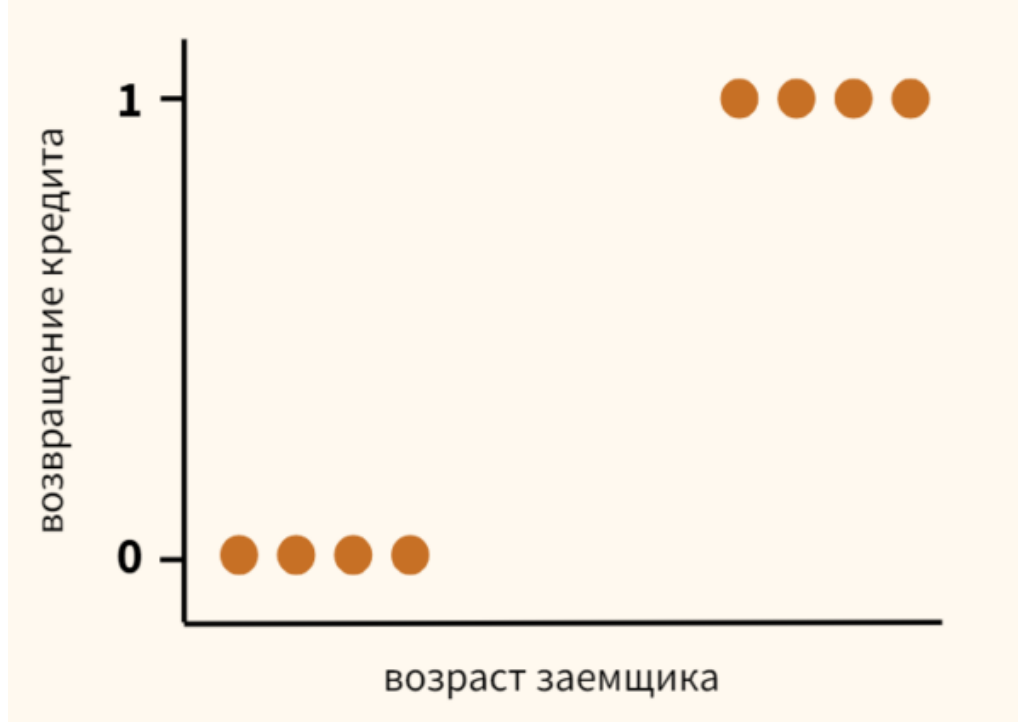
- Производная линейной функции
- Подготовка данных
 - Целевая переменная
 - Отбор признаков
 - Масштабирование признаков
- Обучение модели
- Прогноз и оценка качества
 - Прогноз модели
 - Метрика ассигасы и матрица ошибок
 - Решающая граница
- Написание класса
- Сравнение с sklearn
 - Обучение модели
 - Прогноз
 - Оценка качества
 - Решающая граница
- Бинарная полиномиальная регрессия
 - Полиномиальные признаки
 - Моделирование и оценка качества
- Мультиклассовая логистическая регрессия
 - Подготовка данных
 - Подход one-vs-rest
 - Подготовка датасетов
 - Обучение моделей
 - Прогноз и оценка качества
 - Написание класса
 - Сравнение с sklearn

- Мультиклассовая полиномиальная регрессия
- Softmax Regression
 - Подготовка признаков
 - Кодирование целевой переменной
 - Инициализация весов
 - Функция softmax
 - Функция потерь
 - Градиент
 - Обучение модели, прогноз и оценка качества
 - Написание класса
 - Сравнение с sklearn
- Подведем итог

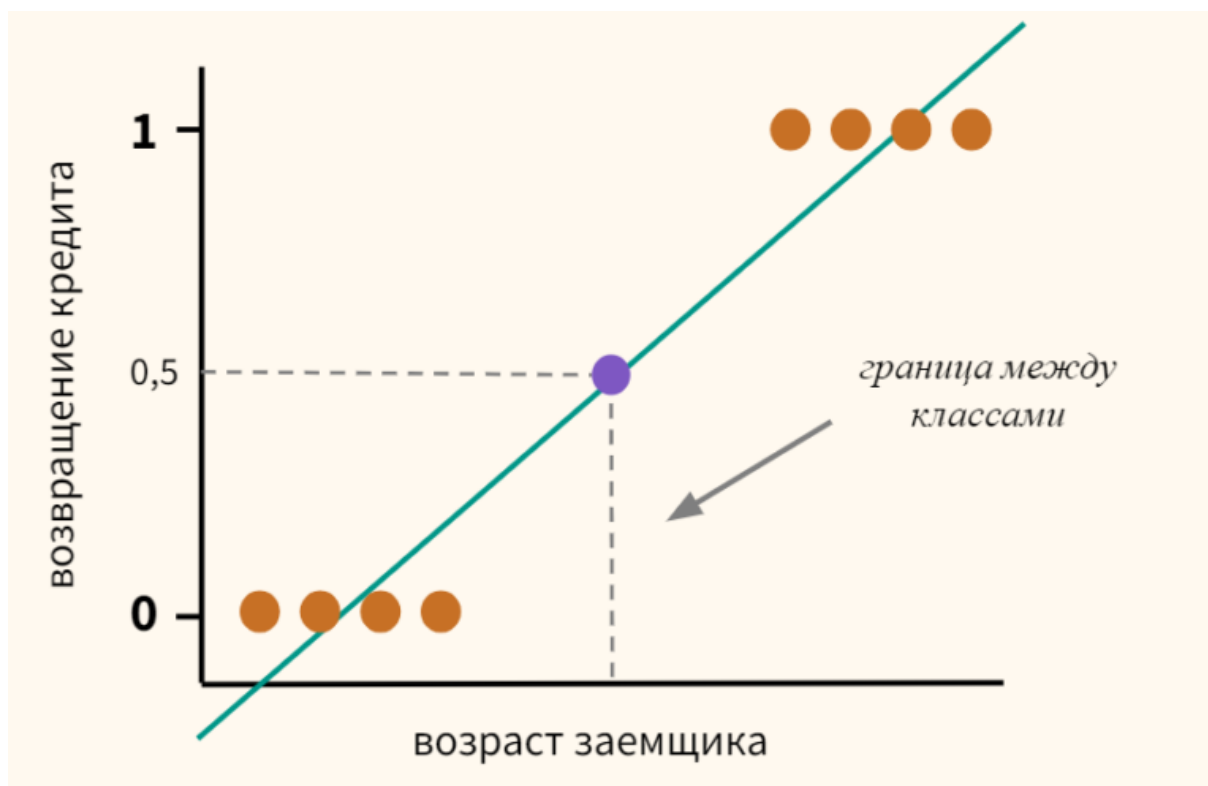
Бинарная логистическая регрессия

Задача бинарной классификации

Вернемся к задаче кредитного скоринга, про которую мы говорили, когда обсуждали принцип машинного обучения. Предположим, что мы собрали данные и выявили зависимость возвращения кредита (ось y) от возраста заемщика (ось x).

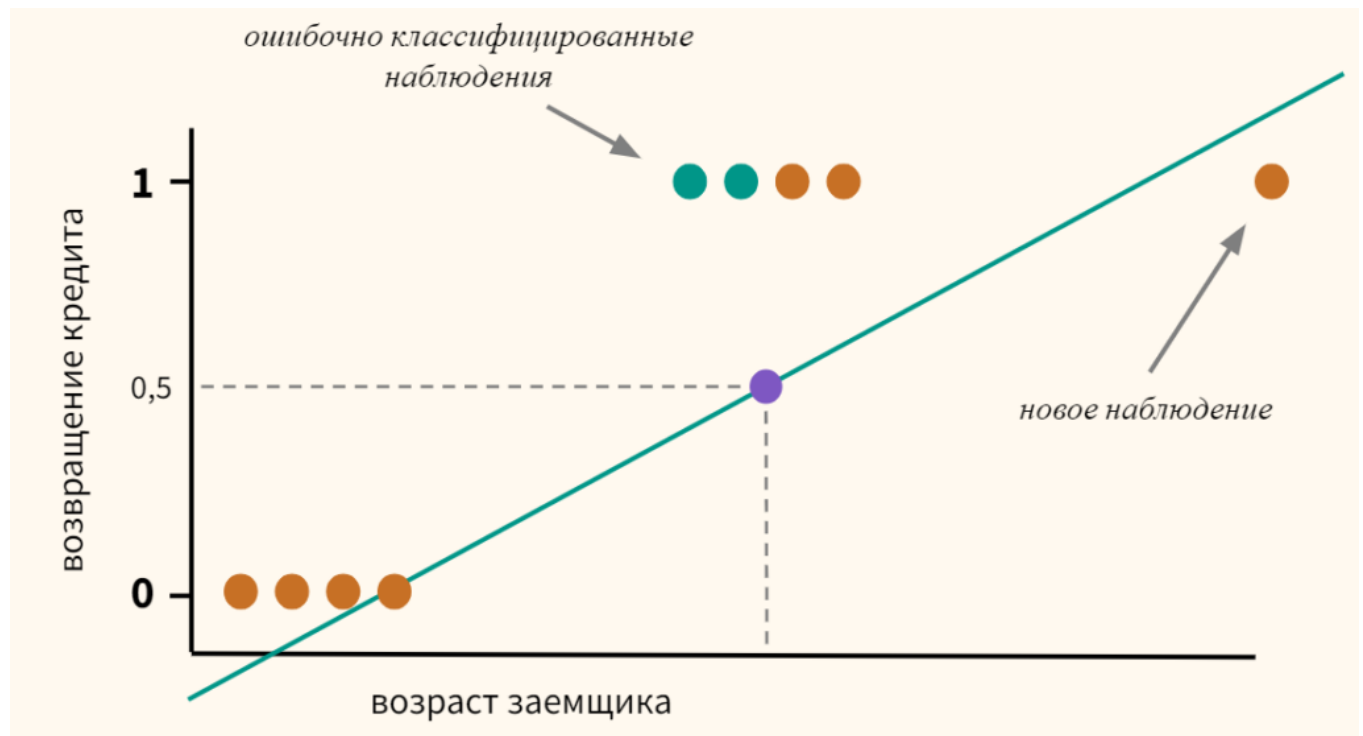


Как мы видим, в среднем более молодые заемщики реже возвращают кредит. Возникает вопрос, с помощью какой модели можно описать эту зависимость? Казалось бы, можно построить линейную регрессию таким образом, чтобы она выдавала некоторое значение y , если это значение окажется ниже 0,5 — отнести наблюдение к классу 0, если выше — к классу 1.



- Если $f_w(x) < 0,5 \rightarrow \hat{y} = 0$
- Если $f_w(x) \geq 0,5 \rightarrow \hat{y} = 1$

Однако, даже если предположить, что мы удачно провели линию регрессии (а на графике выше мы действительно провели ее вполне удачно), и наша модель может делать качественный прогноз, появление новых данных сместит эту границу, и, как следствие, ничего не добавит, а только ухудшит точность модели.



Теперь часть наблюдений, принадлежащих к классу 1, будет ошибочно отнесено моделью к классу 0.

Кроме этого, линейная регрессия по оси y выдает значения, сильно выходящие за пределы интересующего нас интервала от нуля до единицы.

Откроем ноутбук к этому занятию

```
1 # помимо стандартных библиотек мы также импортируем библиотеку warnings
2 # она позволит скрыть предупреждения об ошибках
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 import warnings
8
9 # кроме того, импортируем датасеты библиотеки sklearn
10 from sklearn import datasets
11 # а также функции для расчета метрики accuracy и построения матрицы ошибок
12 from sklearn.metrics import accuracy_score, confusion_matrix
13
14 # построенные нами модели мы будем сравнивать с результатом
```

```

15 # класса LogisticRegression библиотеки sklearn
16 from sklearn.linear_model import LogisticRegression
17
18 # среди прочего, мы построим модели полиномиальной логистической регрессии
19 from sklearn.preprocessing import PolynomialFeatures

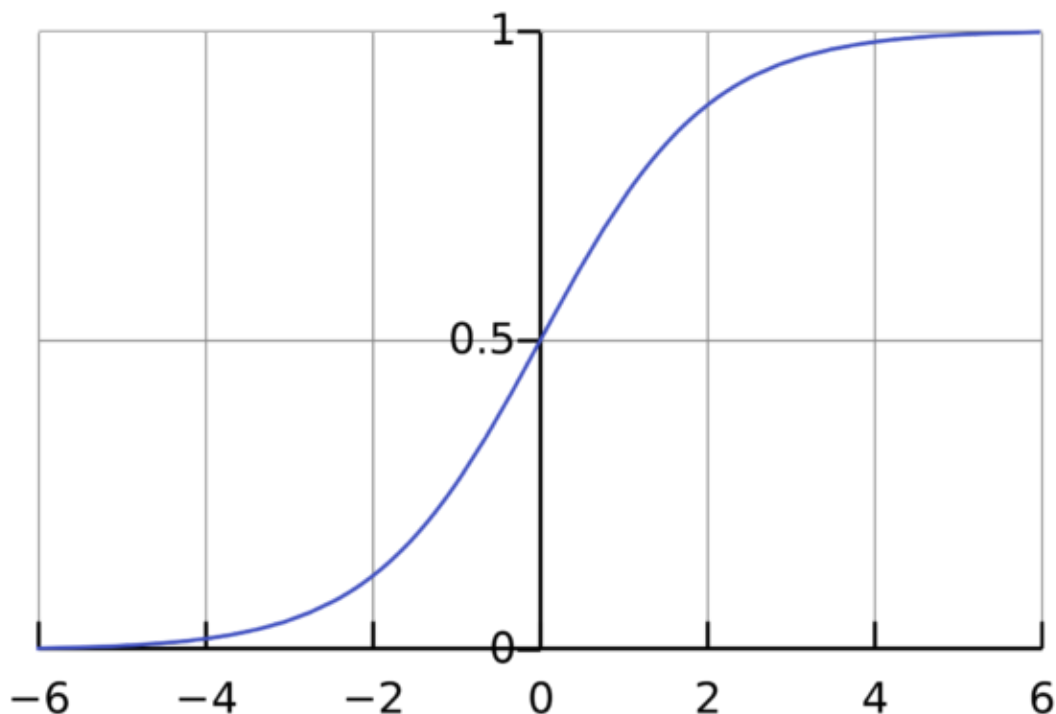
```

Функция логистической регрессии

Сигмоида

Возможное решение упомянутых выше сложностей — пропустить значение линейной регрессии через **сигмоиду** (sigmoid function), которая при любом значении x не выйдет из необходимого нам диапазона $0 \leq h(x) \leq 1$. Напомню формулу и график сигмоиды.

$$g(z) = \frac{1}{1 + e^{-z}}$$



Примечание: обратите внимание, когда z представляет собой большое отрицательное число, знаменатель становится очень большим

$1 + e^{-(-5)} \approx 148,413$ и значение сигмоиды стремится к нулю; когда z является большим положительным числом, знаменатель, а вместе с ним и все выражение стремятся к единице $1 + e^{-(5)} \approx 0,0067$.

Тогда мы можем построить линейную модель, значение которой будет подаваться в сигмоиду.

$$z = X\theta \rightarrow h_{\theta}(x) = \frac{1}{1 + e^{-(X\theta)}}$$

В этом смысле никакой ошибки в названии «логистическая регрессия» нет. Этот алгоритм решает задачу классификации через модель линейной регрессии.

Если вы не помните, почему мы записали множественную линейную функцию как θx , посмотрите [предыдущую лекцию](#).

Приведем код на Питоне.

```
1 def h(x, thetas):  
2     z = np.dot(x, thetas)  
3     return 1.0 / (1 + np.exp(-z))
```

Теперь посмотрим, как интерпретировать коэффициенты.

Интерпретация коэффициентов

Для любого значения x через $h_{\theta}(x)$ мы будем получать вероятность от 0 до 1, что объект принадлежит к классу $y = 1$. Например, если класс 1 означает, что заемщик вернул кредит, то $h_{\theta}(x) = 0,8$ говорит о том, что согласно нашей модели (с параметрами θ), для данного заемщика (x) вероятность возвращения кредита составляет 80 процентов.

В общем случае мы можем записать вероятность вот так.

$$h_{\theta}(x) = P(y = 1|x; \theta)$$

Это выражение можно прочесть как вероятность принадлежности к классу 1 при условии x с параметрами θ (probability of $y = 1$ given x , parameterized by θ).

Поскольку, как мы помним, сумма вероятностей событий, образующих полную группу, всегда равна единице, вероятность принадлежности к классу 0 будет равна

$$P(y = 0|x; \theta) = 1 - P(y = 1|x; \theta)$$

Решающая граница

Решающая граница (decision boundary) — это порог, который определяет к какому классу отнести то или иное наблюдение. Если выбрать порог на уровне 0,5, то все что выше или равно этому порогу мы отнесем к классу 1, все что ниже — к классу 0.

$$y = 1, h_{\theta}(x) \geq 0,5$$

$$y = 0, h_{\theta}(x) < 0,5$$

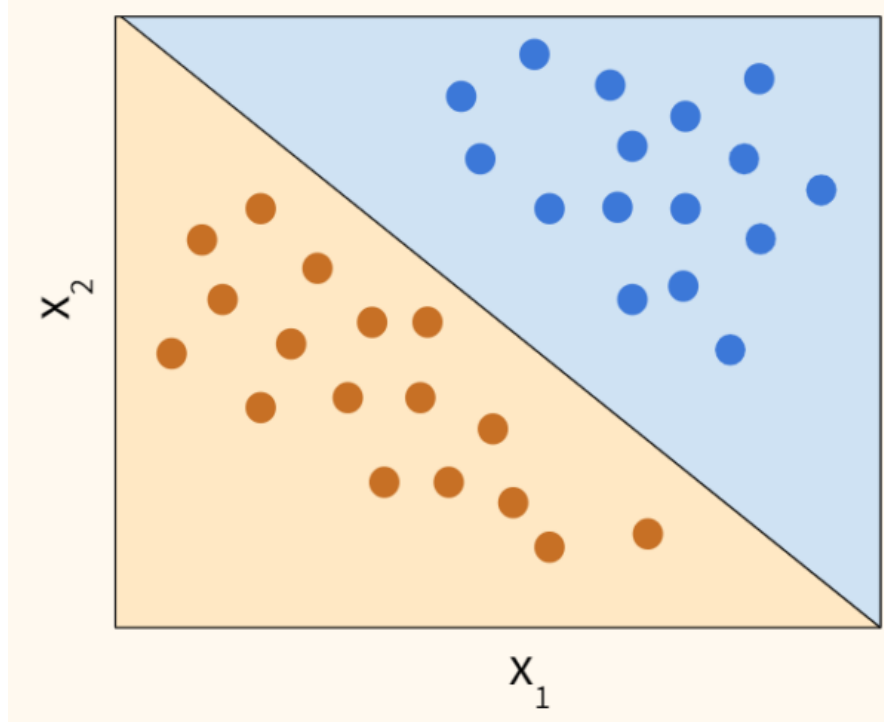
Теперь обратите внимание на сигмоиду. Сигмоида $g(z)$ принимает значения больше 0,5, если $z \geq 0$, а так как $z = X\theta$, то можно сказать, что

- $h_{\theta}(x) \geq 0,5$ и $y = 1$, когда $X\theta \geq 0$, и соответственно
- $h_{\theta}(x) < 0,5$ и $y = 0$, когда $X\theta < 0$.

Уравнение решающей границы

Предположим, что у нас есть два признака x_1 и x_2 . Вместе они образуют так называемое пространство ввода (input space), то есть все имеющиеся у нас наблюдения. Мы можем представить это пространство на координатной плоскости, дополнительно выделив цветом наблюдения, относящиеся к разным классам.

Кроме того, представим, что мы уже построили модель логистической регрессии, и она провела для нас соответствующую границу между двумя классами.



Возникает вопрос. Как, зная коэффициенты θ_0 , θ_1 и θ_2 модели, найти уравнение линии решающей границы? Для начала договоримся, что уравнение решающей границы будет иметь вид $x_2 = mx_1 + c$, где m — наклон прямой, а c — сдвиг.

Теперь вспомним, что модель с двумя признаками (до подачи в сигмоиду) имеет вид

$$z = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

Также не забудем, что граница проходит там, где $h_\theta(x) = 0,5$, а значит $z = 0$. Значит,

$$0 = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

Чтобы найти c (то есть сдвиг линии решающей границы вдоль оси x_2) приравняем x_1 к нулю и решим для x_2 (именно эта точка и будет сдвигом c).

$$0 = \theta_0 + 0 + \theta_2 x_2 \rightarrow x_2 = -\frac{\theta_0}{\theta_2} \rightarrow c = -\frac{\theta_0}{\theta_2}$$

Теперь займемся наклоном m . Возьмем некоторую точку на линии решающей границы с координатами (x_1^a, x_2^a) , (x_1^b, x_2^b) . Тогда наклон m будет равен

$$m = \frac{x_2^b - x_2^a}{x_1^b - x_1^a}$$

Так как эти точки расположены на решающей границе, то справедливо, что

$$0 = \theta_1 x_1^b + \theta_2 x_2^b + \theta_0 - (\theta_1 x_1^a + \theta_2 x_2^a + \theta_0)$$
$$-\theta_2(x_2^b - x_2^a) = \theta_1(x_1^b - x_1^a)$$

А значит,

$$\frac{x_2^b - x_2^a}{x_1^b - x_1^a} = -\frac{\theta_1}{\theta_2} \rightarrow m = -\frac{\theta_1}{\theta_2}$$

Вычислительная устойчивость сигмоиды

При очень больших отрицательных или положительных значениях z может возникнуть переполнение памяти (overflow).

```
1 # возьмем большое отрицательное значение
2 z = -999
3 1 / (1 + np.exp(-z))
```

```
1 RuntimeWarning: overflow encountered in exp
2 0.0
```

Преодолеть это ограничение и добиться **вычислительной устойчивости** (numerical stability) алгоритма можно с помощью следующего тождества.

$$g(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-z}} \times \frac{e^z}{e^z} = \frac{e^z}{e^z(1 + e^{-z})} = \frac{e^z}{e^z + 1}$$

Что интересно, первая часть тождества устойчива при очень больших положительных значениях z .

```
1 z = 999
2 1 / (1 + np.exp(-z))
```

```
1 1.0
```

При этом вторая стабильна при очень больших отрицательных значениях.

```
1 z = -999
2 np.exp(z) / (np.exp(z) + 1)
```

```
1 0.0
```

Объединим обе части с помощью условия.

```
1 def stable_sigmoid(z):
2     if z >= 0:
3         return 1 / (1 + np.exp(-z))
4     else:
5         return np.exp(z) / (np.exp(z) + 1)
```

Примечание. Мы не использовали более лаконичный код, например, функцию `np.where()`, потому что эта функция прежде чем применить условие рассчитывает оба сценария (в данном случае обе части тождества), а это ровно то, чего мы хотим избежать, чтобы не возникло ошибки. Простое условие с `if` препятствует выполнению той части кода, которая нам не нужна.

Можно также использовать **функцию `expit()`** библиотеки `scipy`.

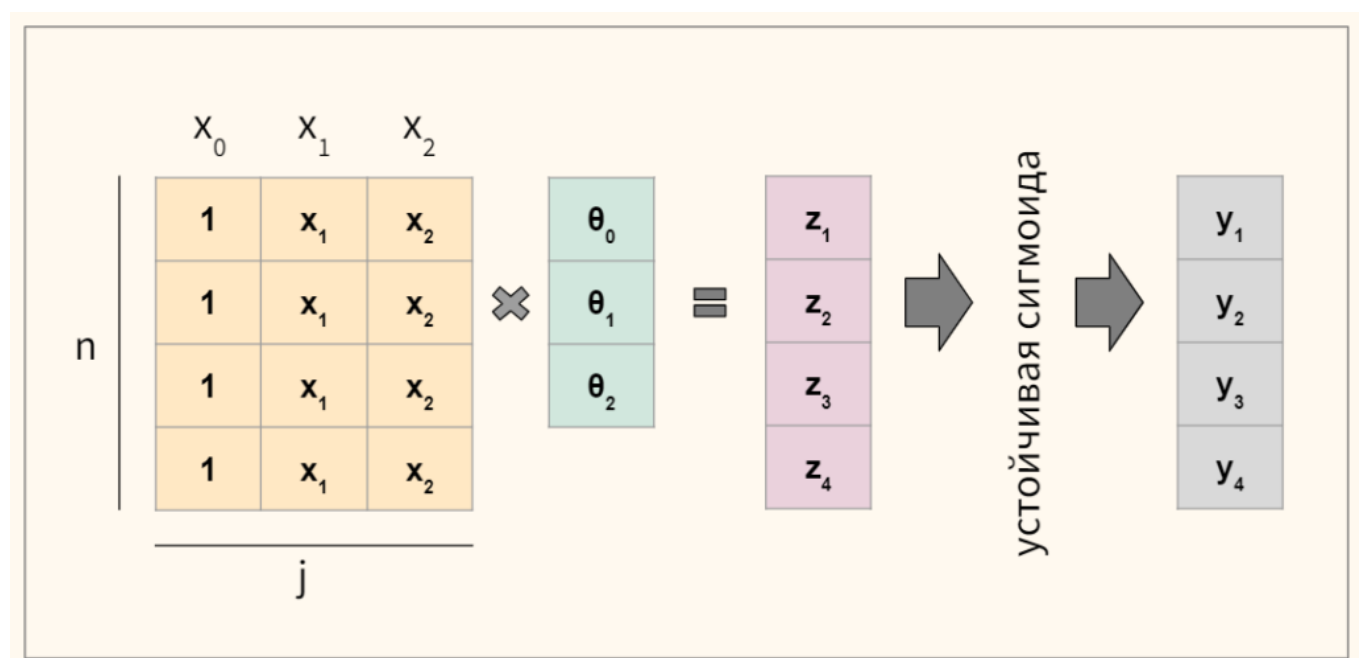
```
1 from scipy.special import expit
2 expit(999), expit(-999)
```

```
1 (1.0, 0.0)
```

Остается написать линейную функцию и подать ее результат в сигмоиду.

```
1 def h(x, thetas):
2     z = np.dot(x, thetas)
3
4     return np.array([stable_sigmoid(value) for value in z])
```

Протестируем код. Предположим, что в нашем датасете четыре наблюдения и три коэффициента. Схематично расчеты будут выглядеть следующим образом.



Пропишем это на Питоне.

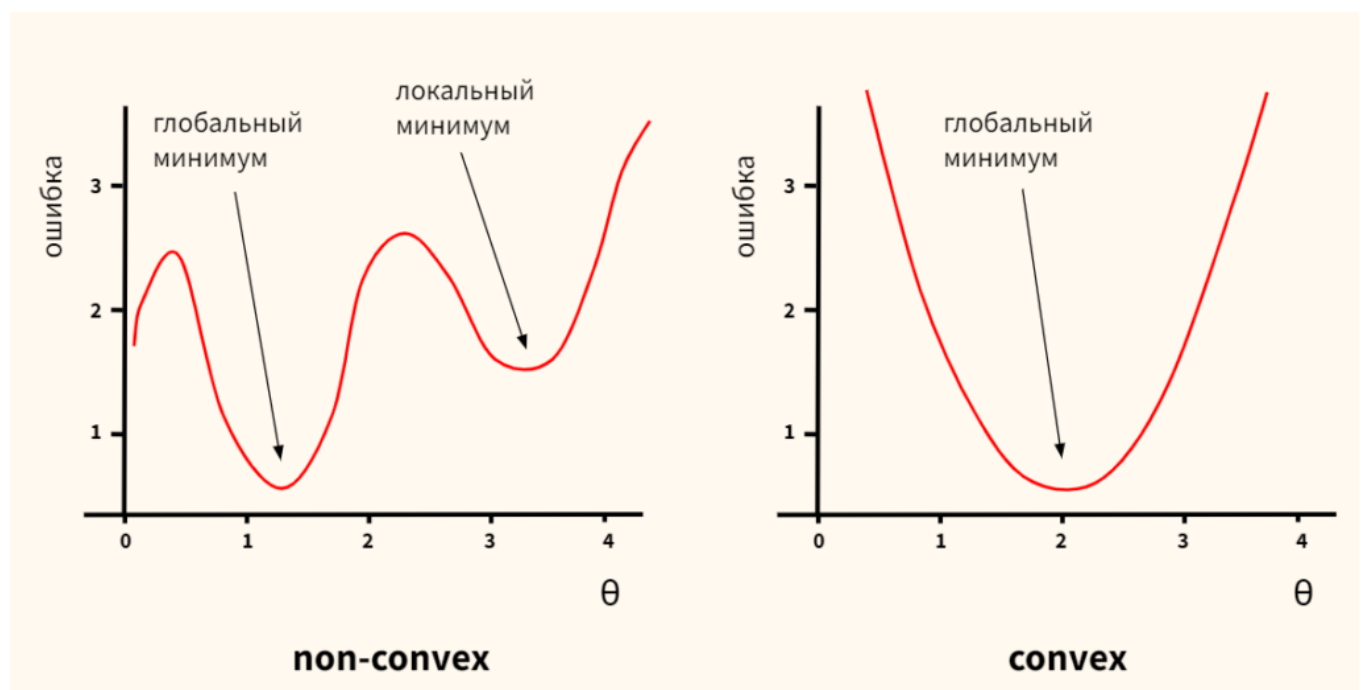
```
1 # возьмем массив наблюдений 4 x 3 с числами от 1 до 12
2 x = np.arange(1, 13).reshape(4, 3)
3
4 # и трехмерный вектор коэффициентов
5 thetas = np.array([-3, 1, 1])
6
7 # подадим их в модель
8 h(x, thetas)
```

```
1 array([0.88079708, 0.26894142, 0.01798621, 0.00091105])
```

Модель работает корректно. Теперь обсудим, как ее обучать, то есть какую функцию потерь использовать для оптимизации параметров θ .

Logistic loss или функция кросс-энтропии

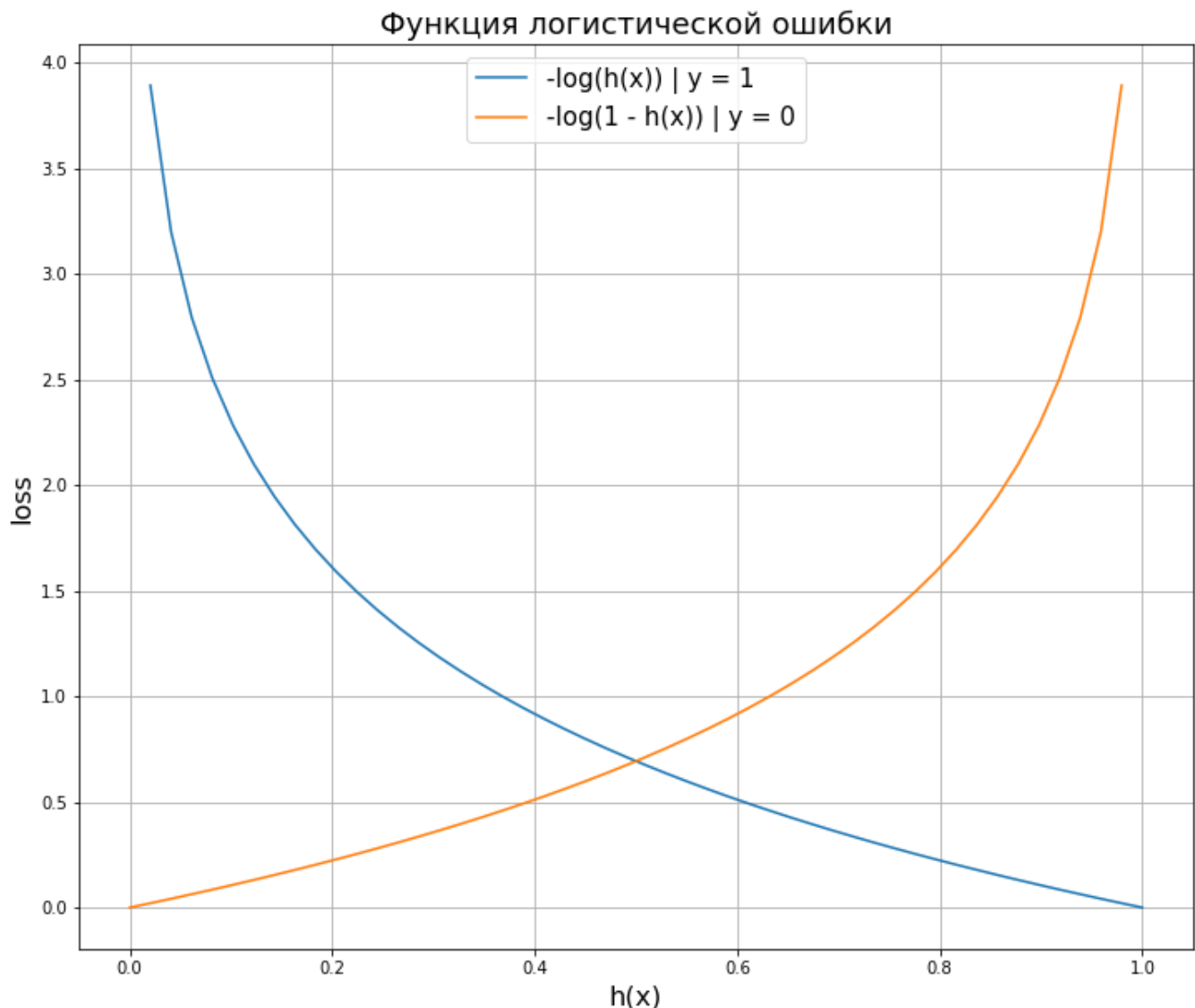
В модели логистической регрессии мы не можем использовать MSE. Дело в том, что если мы поместим результат сигмоиды (представляющей собою нелинейную функцию) в MSE, то на выходе получим невыпуклую функцию (non-convex), глобальный минимум которой довольно сложно найти.



Вместо MSE мы будем использовать функцию логистической ошибки, которую еще называют функцией бинарной кросс-энтропии (log loss, binary cross-entropy loss).

График и формула функции логистической ошибки

Вначале посмотрим на нее на графике.



Разберемся, как она работает. Наша модель $h_{\theta}(x)$ может выдавать *вероятность* от 0 до 1, фактические значения y только 0 и 1.

Сценарий 1. Предположим, что для конкретного заемщика в обучающем датасете истинное значение/ целевой класс записан как 1 (то есть заемщик вернул кредит). Тогда «срабатывает» синяя ветвь графика и ошибка измеряется по ней. Соответственно, чем ближе выдаваемая моделью вероятность к единице, тем меньше ошибка.

$$-\log(P(y = 1|x; \theta)) = -\log(h_{\theta}(x)), y = 1$$

Сценарий 2. Заемщик не вернул кредит и его целевая переменная записана как 0. Тогда срабатывает оранжевая ветвь. Ошибка модели будет минимальна при значениях, близких к нулю.

$$-\log(1 - P(y = 1|x; \theta)) = -\log(1 - h_{\theta}(x)), y = 0$$

Добавлю, что **минус логарифм** в данном случае очень удачно отвечает нашему желанию иметь нулевую ошибку при правильном прогнозе и наказать алгоритм высокой ошибкой (асимптотически стремящейся к бесконечности) в случае неправильного прогноза.

В итоге нам нужно будет найти **сумму вероятностей** принадлежности к классу 1 для сценария 1 и сценария 2.

$$J(\theta) = \begin{cases} -\log(h_{\theta}(x))|y = 1 \\ -\log(1 - h_{\theta}(x))|y = 0 \end{cases}$$

Однако, для каждого наблюдения нам нужно учитывать только одну из вероятностей (либо $y = 1$, либо $y = 0$). Как нам переключаться между ними? На самом деле очень просто.

В качестве **переключателя** можно использовать целевую переменную. В частности, умножим левую часть функции на y , а правую на $1 - y$. Тогда, если речь идет о классе 1, первая часть умножится на единицу, вторая на ноль и исчезнет. Если речь идет о классе 0, произойдет обратное, исчезнет левая часть, а правая останется. Получается

$$J(\theta) = -\frac{1}{n} \sum y \cdot \log(h_{\theta}(x)) + (1 - y) \cdot \log(1 - h_{\theta}(x))$$

Рассмотрим ее работу на учебном примере.

Расчет логистической ошибки

Предположим, мы построили модель и для каждого наблюдения получили некоторый прогноз (вероятность).

```
1 # выведем результат работы модели (вероятности) y_pred и целевую переменную y
2 output = pd.DataFrame({
3     'y'           : [1, 1, 1, 0, 0, 1, 1, 0],
4     'y_pred'      : [0.93, 0.81, 0.78, 0.43, 0.54, 0.49, 0.22, 0.1]
5 })
6
7 output
```

| | y | y_pred |
|---|---|--------|
| 0 | 1 | 0.93 |
| 1 | 1 | 0.81 |
| 2 | 1 | 0.78 |
| 3 | 0 | 0.43 |
| 4 | 0 | 0.54 |
| 5 | 1 | 0.49 |
| 6 | 1 | 0.22 |
| 7 | 0 | 0.10 |

Найдем вероятность принадлежности к классу 1.

```
1 # оставим вероятность, если y = 1, и вычтем вероятность из единицы, если
2 output['y=1 prob'] = np.where(output['y'] == 0, 1 - output['y_pred'], ou
3 output
```

| | y | y_pred | y=1 prob |
|---|---|--------|----------|
| 0 | 1 | 0.93 | 0.93 |
| 1 | 1 | 0.81 | 0.81 |
| 2 | 1 | 0.78 | 0.78 |
| 3 | 0 | 0.43 | 0.57 |
| 4 | 0 | 0.54 | 0.46 |
| 5 | 1 | 0.49 | 0.49 |
| 6 | 1 | 0.22 | 0.22 |
| 7 | 0 | 0.10 | 0.90 |

Возьмем отрицательный логарифм из каждой вероятности.

```
1 output['-log'] = -np.log(output['y=1 prob'])
2 output
```

| | y | y_pred | y=1 prob | -log |
|---|---|--------|----------|----------|
| 0 | 1 | 0.93 | 0.93 | 0.072571 |
| 1 | 1 | 0.81 | 0.81 | 0.210721 |
| 2 | 1 | 0.78 | 0.78 | 0.248461 |
| 3 | 0 | 0.43 | 0.57 | 0.562119 |
| 4 | 0 | 0.54 | 0.46 | 0.776529 |
| 5 | 1 | 0.49 | 0.49 | 0.713350 |
| 6 | 1 | 0.22 | 0.22 | 1.514128 |
| 7 | 0 | 0.10 | 0.90 | 0.105361 |

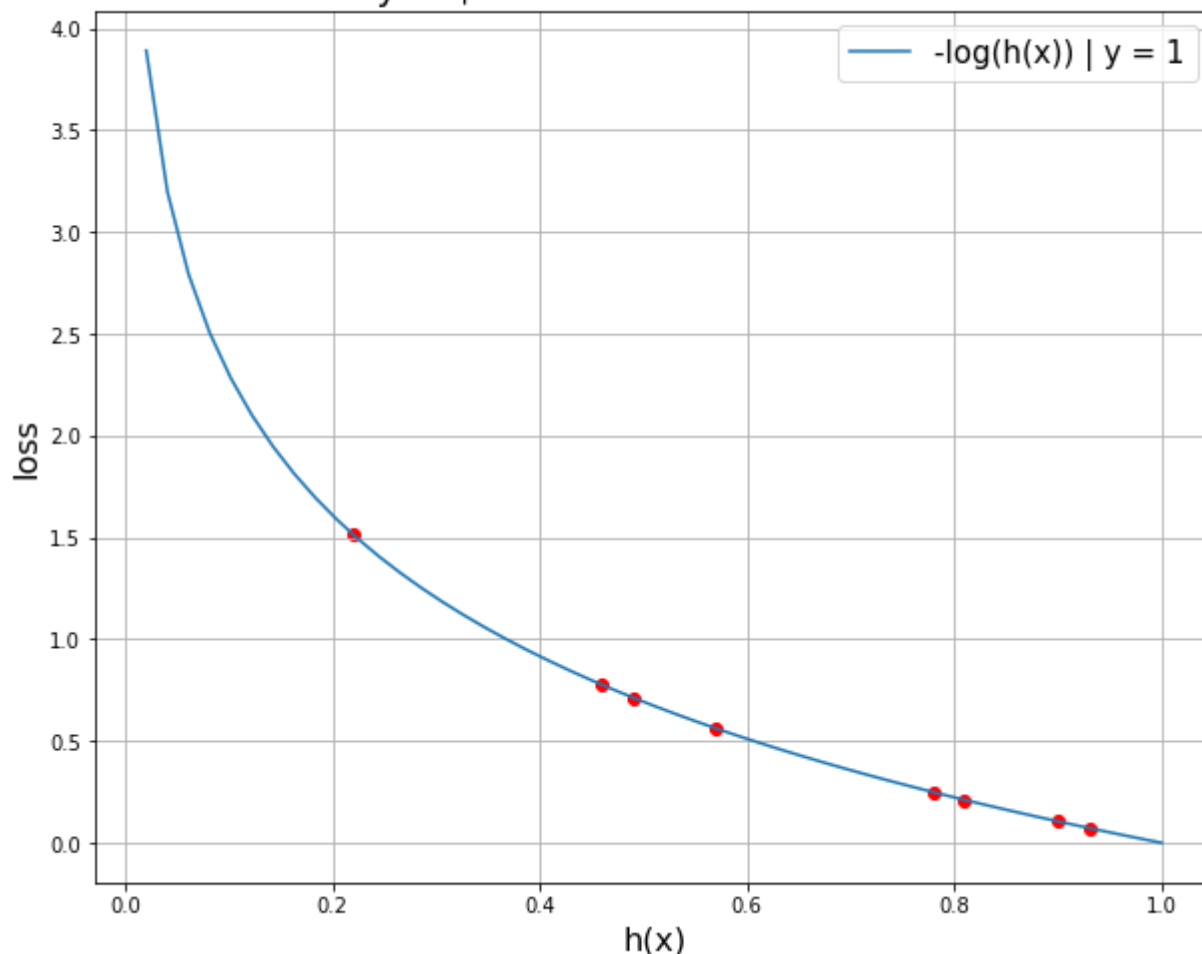
Выведем каждое из получившихся значений на графике.

```

1 plt.figure(figsize = (10, 8))
2
3 # создадим точки по оси x в промежутке от 0 до 1
4 x_vals = np.linspace(0, 1)
5
6 # выведем кривую функции логистической ошибки
7 plt.plot(x_vals, -np.log(x_vals), label = '-log(h(x)) | y = 1')
8
9 # выведем каждое из значений отрицательного логарифма
10 plt.scatter(output['y=1 prob'], output['-log'], color = 'r')
11
12 # зададим заголовок, подписи к осям, легенду и сетку
13 plt.xlabel('h(x)', fontsize = 16)
14 plt.ylabel('loss', fontsize = 16)
15 plt.title('Функция логистической ошибки', fontsize = 18)
16 plt.legend(loc = 'upper right', prop = {'size': 15})
17 plt.grid()
18
19 plt.show()

```


Функция логистической ошибки



Как мы видим, так как мы всегда выражаем вероятность принадлежности к классу 1, графически нам будет достаточно одной ветви. Остается сложить результаты и разделить на количество наблюдений.

```
1 | output['-log'].mean()
```

```
1 | 0.5254048659083239
```

Окончательный вариант

Напишем функцию логистической ошибки, которую будем использовать в нашем алгоритме.

```
1 | def objective(y, y_pred):
2 |
3 |     # рассчитаем функцию потерь для y = 1, добавив 1e-9, чтобы избежать о
4 |     y_one_loss = y * np.log(y_pred + 1e-9)
5 |
6 |     # также рассчитаем функцию потерь для y = 0
7 |     y_zero_loss = (1 - y) * np.log(1 - y_pred + 1e-9)
8 |
9 |     # сложим и разделим на количество наблюдений
```

```
10 | return -np.mean(y_zero_loss + y_one_loss)
```

Проверим ее работу на учебных данных.

```
1 | objective(output['y'], output['y_pred'])
```

```
1 | 0.525404864006128
```

Теперь займемся поиском производной.

Производная функции логистической ошибки

Предположим, что $G(\theta)$ — одна из частных производных описанной выше функции логистической ошибки $J(\theta)$,

$$G = y \cdot \log(h) + (1 - y) \cdot \log(1 - h)$$

где h — это сигмоида $1/(1 + e^{-z})$, а $z(\theta)$ — линейная функция $x\theta$. Тогда по chain rule нам нужно найти производные следующих функций

$$\frac{\partial G}{\partial \theta} = \frac{\partial G}{\partial h} \cdot \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial \theta}$$

Производная логарифмической функции

Начнем с производной логарифмической функции.

$$\frac{\partial}{\partial x} \ln f(x) = \frac{1}{f(x)}$$

Теперь, помня, что x и y — это константы, найдем первую производную.

$$\begin{aligned} & \frac{\partial G}{\partial h} [y \cdot \log(h) + (1 - y) \cdot \log(1 - h)] \\ &= y \cdot \frac{\partial G}{\partial h} [\log(h)] + (1 - y) \cdot \frac{\partial G}{\partial h} [\log(1 - h)] \\ &= \frac{1}{h} y + \frac{1}{1 - h} \cdot \frac{\partial G}{\partial h} [1 - h] \cdot (1 - y) \end{aligned}$$

Упростим выражение (не забыв про производную разности).

$$\begin{aligned}
&= \frac{h}{y} + \frac{\frac{\partial G}{\partial h}(1-h)(1-y)}{1-h} = \frac{h}{y} + \frac{(0-1)(1-y)}{1-h} \\
&= \frac{y}{h} - \frac{1-y}{1-h} = \frac{y-h}{h(1-h)}
\end{aligned}$$

Теперь займемся производной сигмоиды.

Производная сигмоиды

Вначале упростим выражение.

$$\frac{\partial h}{\partial z} \left[\frac{1}{1 + e^{-z}} \right] = \frac{\partial h}{\partial z} [(1 + e^{-z})^{-1}]$$

Теперь перейдем к нахождению производной

$$\begin{aligned}
&= -(1 + e^{-z})^{-2}) \cdot (-e^{-z}) = \frac{e^{-z}}{(1 + e^{-z})^2} \\
&= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}} = \frac{1}{1 + e^{-z}} \cdot \frac{(1 + e^{-z}) - 1}{1 + e^{-z}} \\
&= \frac{1}{1 + e^{-z}} \cdot \left(\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \right) \\
&= \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}} \right)
\end{aligned}$$

В терминах предложенной выше нотации получается

$$h(1 - h)$$

Производная линейной функции

Наконец найдем производную линейной функции.

$$\frac{\partial z}{\partial \theta} = x$$

Перемножим производные и найдем градиент по каждому из признаков j для n наблюдений.

$$\frac{\partial J}{\partial \theta} = \frac{y - h}{h(1 - h)} \cdot h(1 - h) \cdot x_j \cdot \frac{1}{n} = x_j \cdot (y - h) \cdot \frac{1}{n}$$

Замечу, что хотя производная похожа на градиент функции линейной регрессии, на самом деле это разные функции, h в данном случае сигмоида.

Для нахождения градиента (всех частных производных одновременно) перепишем формулу в векторной нотации.

$$\nabla_{\theta} J = X^T (h(X\theta) - y) \times \frac{1}{n}$$

Схематично для четырех наблюдений и трех коэффициентов нахождение градиента будет выглядеть следующим образом.

Объявим соответствующую функцию.

```
1 def gradient(x, y, y_pred, n):  
2     return np.dot(x.T, (y_pred - y)) / n
```

На всякий случай напомним, что прогнозные значения (y_{pred}) мы получаем с помощью объявленной ранее функции $h(x, \theta)$.

Подготовка данных

В качестве примера возьмем встроенный в sklearn датасет, в котором нам предлагается определить класс вина по его характеристикам.

```
1  # импортируем датасет о вине из модуля datasets
2  data = datasets.load_wine()
3
4  # превратим его в датафрейм
5  df = pd.DataFrame(data.data, columns = data.feature_names)
6
7  # добавим целевую переменную
8  df['target'] = data.target
9
10 # посмотрим на первые три строки
11 df.head(3)
```

Выше представлена только часть датасета. Полностью его можно посмотреть в [ноутбуке](#).

Целевая переменная

Посмотрим на количество наблюдений и признаков (размерность матрицы), а также уникальные значения (классы) в целевой переменной.

```
1  df.shape, np.unique(df.target)
```

```
1  ((178, 14), array([0, 1, 2]))
```

Как мы видим, у нас три класса, а должно быть два, потому что пока что мы создаем алгоритм бинарной классификации. Отфильтруем значения так, чтобы осталось только два класса.

```
1  # применим маску датафрейма и удалим класс 2
```

```
2 df = df[df.target != 2]
3
4 # посмотрим на результат
5 df.shape, df.target.unique()
```

```
1 ((130, 14), array([0, 1]))
```

Отбор признаков

Наша целевая переменная выражена бинарной категорией или, как еще говорят, находится на дихотомической шкале (dichotomous variable). В этом случае применять коэффициент корреляции Пирсона не стоит и можно использовать точечно-бисериальную корреляцию (point-biserial correlation). Рассчитаем корреляцию признаков и целевой переменной нашего датасета.

```
1 # импортируем модуль stats из библиотеки scipy
2 from scipy import stats
3
4 # создадим два списка, один для названий признаков, второй для корреляций
5 columns, correlations = [], []
6
7 # пройдемся по всем столбцам датафрейма кроме целевой переменной
8 for col in df.drop('target', axis = 1).columns:
9     # поместим название признака в список columns
10    columns.append(col)
11    # рассчитаем корреляцию этого признака с целевой переменной
12    # и поместим результат в список корреляций
13    correlations.append(stats.pointbiserialr(df[col], df['target'])[0])
14
15 # создадим датафрейм на основе заполненных списков
16 # и применим градиентную цветовую схему
17 pd.DataFrame({'column': columns, 'correlation': correlations}).style.background
```

Наиболее коррелирующим с целевой переменной признаком является пролин (proline). Визуально оценим насколько сильно отличается этот показатель для классов вина 0 и 1.

```
1 # зададим размер графика
2 plt.figure(figsize = (10, 8))
3 # на точечной диаграмме выведем пролин по оси x, а класс вина по оси y
4 sns.scatterplot(x = df.proline, y = df.target, s = 80);
```

Теперь посмотрим на зависимость двух признаков (спирт и пролин) от целевой переменной.

```
1 # зададим размер графика
2 plt.figure(figsize = (10, 8))
3 # на точечной диаграмме по осям x и y выведем признаки,
4 # с помощью параметра hue разделим соответствующие классы целевой переменной
5 sns.scatterplot(x = df.alcohol, y = df.proline, hue = df.target, s = 80)
6 # добавим легенду, зададим ее расположение и размер
7 plt.legend(loc = 'upper left', prop = {'size': 15})
8 # выведем результат
9 plt.show()
```


В целом можно сказать, что классы линейно разделимы (другими словами, мы можем провести прямую между ними). Поместим признаки в переменную X , а целевую переменную — в y .

```
1 X = df[['alcohol', 'proline']]
2 y = df['target']
```

Масштабирование признаков

Как и в случае с линейной регрессией, для алгоритма логистической регрессии важно, чтобы признаки были приведены к одному масштабу. Для этого используем стандартизацию.

```
1 # т.е. приведем данные к нулевому среднему и единичному СКО
2 X = (X - X.mean()) / X.std()
3 X.head()
```

Проверим результат.

```
1 | X.alcohol.mean(), X.alcohol.std(), X.proline.mean(), X.proline.std()
```

```
1 | (6.8321416900009635e-15, 1.0, -5.465713352000771e-17, 1.0)
```

Теперь мы готовы к созданию и обучению модели.

Обучение модели

Вначале объявим уже знакомую нам функцию, которая добавит в датафрейм столбец под названием `x0`, заполненный единицами.

```
1 | def add_ones(x):  
2 |     # важно! метод .insert() изменяет исходный датафрейм  
3 |     return x.insert(0, 'x0', np.ones(x.shape[0]))
```

Применим ее к нашему датафрейму с признаками.

```
1 | # добавим столбец с единицами  
2 | add_ones(X)  
3 |  
4 | # и посмотрим на результат  
5 | X.head()
```

Создадим вектор начальных весов (он будет состоять из нулей), а также переменную `n`, в которой будет храниться количество наблюдений.

```
1 | thetas, n = np.zeros(X.shape[1]), X.shape[0]
2 | thetas, n
```

```
1 | (array([0., 0., 0.]), 130)
```

Кроме того, создадим список, в который будем записывать размер ошибки функции потерь.

```
1 | loss_history = []
```

Теперь выполним основную работу по минимизации функции потерь и поиску оптимальных весов (выполнение кода ниже у меня заняло около 30 секунд).

```
1 | # в цикле из 20000 итераций
2 | for i in range(20000):
3 |     # рассчитаем прогнозное значение с текущими весами
4 |     y_pred = h(X, thetas)
5 |     # посчитаем уровень ошибки при текущем прогнозе
6 |     loss_history.append(objective(y, y_pred))
7 |     # рассчитаем градиент
8 |     grad = gradient(X, y, y_pred, n)
9 |     # используем градиент для улучшения весов модели
10 |    # коэффициент скорости обучения будет равен 0,001
11 |    thetas = thetas - 0.001 * grad
```

Посмотрим на получившиеся веса и финальный уровень ошибки.

```
1 | # чтобы посмотреть финальный уровень ошибки,
2 | # достаточно взять последний элемент списка loss_history
3 | thetas, loss_history[-1]
```

```
1 | (array([ 0.23234188, -1.73394252, -1.89350543]), 0.12282503517421262)
```

Модель обучена. Теперь мы можем сделать прогноз и оценить результат.

Прогноз и оценка качества

Прогноз модели

Объявим функцию `predict()`, которая будет предсказывать к какому классу относится то или иное наблюдение. От функции $h(x, thetas)$ эта функция будет

отличаться тем, что выдаст не только вероятность принадлежности к тому или иному классу, но и непосредственно сам предполагаемый класс (0 или 1).

```
1 def predict(x, thetas):
2     # найдем значение линейной функции
3     z = np.dot(x, thetas)
4     # проведем его через устойчивую сигмоиду
5     probs = np.array([stable_sigmoid(value) for value in z])
6     # если вероятность больше или равна 0,5 - отнесем наблюдение к классу 1
7     # в противном случае к классу 0
8     # дополнительно выведем значение вероятности
9     return np.where(probs >= 0.5, 1, 0), probs
```

Вызовем функцию `predict()` и запишем прогноз класса и вероятность принадлежности к этому классу в переменные `y_pred` и `probs` соответственно.

```
1 # запишем прогноз класса и вероятность этого прогноза в переменные y_pred и probs
2 y_pred, probs = predict(X, thetas)
3
4 # посмотрим на прогноз и вероятность для первого наблюдения
5 y_pred[0], probs[0]
```

```
1 (0, 0.022908352078195617)
```

Здесь важно напомнить, что вероятность, близкая к нулю, говорит о принадлежности к классу 0. В качестве упражнения выведите класс последнего наблюдения и соответствующую вероятность.

Метрика ассурасу и матрица ошибок

Оценим результат с помощью метрики ассурасу и матрицы ошибок.

```
1 # функцию accuracy_score() мы импортировали в начале ноутбука
2 accuracy_score(y, y_pred)
```

```
1 0.9615384615384616
```

```
1 # функцию confusion_matrix() мы импортировали в начале ноутбука
2 # столбцами будут прогнозные значения (Forecast),
3 # строками - фактические (Actual)
4 pd.DataFrame(confusion_matrix(y, y_pred),
5               columns = ['Forecast 0', 'Forecast 1'],
6               index = ['Actual 0', 'Actual 1'])
```

Как мы видим, алгоритм ошибся пять раз. Дважды он посчитал, что наблюдение относится к классу 1, хотя на самом деле это был класс 0, и трижды, наоборот, неверно отнес класс 1 к классу 0.

Решающая граница

Выше мы уже вывели уравнение решающей границы. Воспользуемся им, чтобы визуально оценить насколько удачно классификатор справился с поставленной задачей.

```
1 # рассчитаем сдвиг (с) и наклон (m) линии границы
2 c, m = -thetas[0]/thetas[2], -thetas[1]/thetas[2]
3 c, m
```

```
1 (0.1227046263531282, -0.915731474695505)
```

```
1 # найдем минимальное и максимальное значения для спирта (ось x)
2 xmin, xmax = min(X['alcohol']), max(X['alcohol'])
3 # найдем минимальное и максимальное значения для пролина (ось y)
4 ymin, ymax = min(X['proline']), max(X['proline'])
5
6 # запишем значения оси x в переменную xd
7 xd = np.array([xmin, xmax])
8 xd
```

```
1 array([-2.15362589,  2.12194856])
```

```
1 # подставим эти значения, а также значения сдвига и наклона в уравнение
2 yd = m * xd + c
3
4 # в результате мы получим координаты двух точек, через которые проходит
5 (xd[0], yd[0]), (xd[1], yd[1])
```

```
1 ((-2.1536258890738247, 2.0948476376971197),
2  (2.1219485561396647, -1.8204304541886445))
```

```
1 # зададим размер графика
2 plt.figure(figsize = (11, 9))
3
4 # построим пунктирную линию по двум точкам, найденным выше
5 plt.plot(xd, yd, 'k', lw = 1, ls = '--')
```

```

6
7 # дополнительно отобразим наши данные,
8 sns.scatterplot(x = X['alcohol'], y = X['proline'], hue = y, s = 70)
9
10 # которые снова снабдим легендой
11 plt.legend(loc = 'upper left', prop = {'size': 15})
12
13 # минимальные и максимальные значения по обеим осям будут границами графика
14 plt.xlim(xmin, xmax)
15 plt.ylim(ymin, ymax)
16
17 # по желанию, разделенные границей половинки можно закрасить
18 # tab: означает, что цвета берутся из палитры Tableau
19 # plt.fill_between(xd, yd, ymin, color='tab:blue', alpha = 0.2)
20 # plt.fill_between(xd, yd, ymax, color='tab:orange', alpha = 0.2)
21
22 # а также добавить обозначения переменных в качестве подписей к осям
23 # plt.xlabel('x_1')
24 # plt.ylabel('x_2')
25
26 plt.show()

```

На графике хорошо видны те пять значений, в которых ошибся наш классификатор.

Написание класса

Остается написать класс бинарной логистической регрессии.

```
1 class LogReg():
2
3     # в методе .__init__() объявим переменные для весов и уровня ошибки
4     def __init__(self):
5         self.thetas = None
6         self.loss_history = []
7
8     # метод .fit() необходим для обучения модели
9     # этому методу мы передадим признаки и целевую переменную
10    # кроме того, мы зададим значения по умолчанию
11    # для количества итераций и скорости обучения
12    def fit(self, x, y, iter = 20000, learning_rate = 0.001):
13
14        # метод создаст "правильные" копии датафрейма
15        x, y = x.copy(), y.copy()
16
17        # добавит столбец из единиц
18        self.add_ones(x)
19
20        # инициализирует веса и запишет в переменную n количество наблюдений
21        thetas, n = np.zeros(x.shape[1]), x.shape[0]
22
23        # создадим список для записи уровня ошибки
24        loss_history = []
25
26        # в цикле равном количеству итераций
27        for i in range(iter):
28            # метод сделает прогноз с текущими весами
29            y_pred = self.h(x, thetas)
30            # найдет и запишет уровень ошибки
31            loss_history.append(self.objective(y, y_pred))
32            # рассчитает градиент
33            grad = self.gradient(x, y, y_pred, n)
34            # и обновит веса
35            thetas -= learning_rate * grad
36
37        # метод выдаст веса и список с историей ошибок
38        self.thetas = thetas
39        self.loss_history = loss_history
40
41    # метод .predict() делает прогноз с помощью обученной модели
42    def predict(self, x):
43
44        # метод создаст "правильную" копию модели
45        x = x.copy()
```

```

46     # добавит столбец из единиц
47     self.add_ones(x)
48     # рассчитает значения линейной функции
49     z = np.dot(x, self.thetas)
50     # передаст эти значения в сигмоиду
51     probs = np.array([self.stable_sigmoid(value) for value in z])
52     # выдаст принадлежность к определенному классу и соответствующую ве
53     return np.where(probs >= 0.5, 1, 0), probs
54
55     # ниже приводятся служебные методы, смысл которых был разобран ранее
56     def add_ones(self, x):
57         return x.insert(0, 'x0', np.ones(x.shape[0]))
58
59     def h(self, x, thetas):
60         z = np.dot(x, thetas)
61         return np.array([self.stable_sigmoid(value) for value in z])
62
63     def objective(self, y, y_pred):
64         y_one_loss = y * np.log(y_pred + 1e-9)
65         y_zero_loss = (1 - y) * np.log(1 - y_pred + 1e-9)
66         return -np.mean(y_zero_loss + y_one_loss)
67
68     def gradient(self, x, y, y_pred, n):
69         return np.dot(x.T, (y_pred - y)) / n
70
71     def stable_sigmoid(self, z):
72         if z >= 0:
73             return 1 / (1 + np.exp(-z))
74         else:
75             return np.exp(z) / (np.exp(z) + 1)

```

Проверим работу написанного нами класса. Вначале подготовим данные и обучим модель.

```

1     # проверим работу написанного нами класса
2     # поместим признаки и целевую переменную в X и y
3     X = df[['alcohol', 'proline']]
4     y = df['target']
5
6     # приведем признаки к одному масштабу
7     X = (X - X.mean())/X.std()
8
9     # создадим объект класса LogReg
10    model = LogReg()
11
12    # и обучим модель
13    model.fit(X, y)
14
15    # посмотрим на атрибуты весов и финального уровня ошибки
16    model.thetas, model.loss_history[-1]

```



```
1 | (array([ 0.23234188, -1.73394252, -1.89350543]), 0.12282503517421262)
```

Затем сделаем прогноз и оценим качество модели.

```
1 | # сделаем прогноз
2 | y_pred, probs = model.predict(X)
3 |
4 | # и посмотрим на класс первого наблюдения и вероятность
5 | y_pred[0], probs[0]
```

```
1 | (0, 0.022908352078195617)
```

```
1 | # рассчитаем accuracy
2 | accuracy_score(y, y_pred)
```

```
1 | 0.9615384615384616
```

```
1 | # создадим матрицу ошибок
2 | pd.DataFrame(confusion_matrix(y, y_pred),
3 |               columns = ['Forecast 0', 'Forecast 1'],
4 |               index = ['Actual 0', 'Actual 1'])
```

Модель показала точно такой же результат. Методы класса LogReg работают. Теперь давайте сравним работу нашего класса с классом LogisticRegression библиотеки sklearn.

Сравнение с sklearn

Обучение модели

Вначале обучим модель.

```
1 | # подготовим данные
2 | X = df[['alcohol', 'proline']]
3 | y = df['target']
4 |
5 | X = (X - X.mean())/X.std()
6 |
7 | # создадим объект класса LogisticRegression и запишем его в переменную r
8 | model = LogisticRegression()
9 |
```

```
10 # обучим модель
11 model.fit(X, y)
12
13 # посмотрим на получившиеся веса модели
14 model.intercept . model.coef
```

```
1 (array([0.30838852]), array([[ -2.09622008,  -2.45991159]]))
```

Прогноз

Теперь необходимо сделать прогноз и найти соответствующие вероятности. В классе **LogisticRegression** библиотеки **sklearn** метод **.predict()** отвечает за предсказание принадлежности к определенному классу, а метод **.predict_proba()** отвечает за вероятность такого прогноза.

```
1 # выполним предсказание класса
2 y_pred = model.predict(X)
3
4 # и найдем вероятности
5 probs = model.predict_proba(X)
6
7 # посмотрим на класс и вероятность первого наблюдения
8 y_pred[0], probs[0]
```

```
1 (0, array([0.9904622, 0.0095378]))
```

Модель предсказала для первого наблюдения класс 0. При этом, обратите внимание, что метод **.predict_proba()** для каждого наблюдения выдает две вероятности, первая — это вероятность принадлежности к классу 0, вторая — к классу 1.

Оценка качества

Рассчитаем метрику ассурасу.

```
1 accuracy_score(y, y_pred)
```

```
1 0.9615384615384616
```

И построим матрицу ошибок.

```
1 pd.DataFrame(confusion_matrix(y, y_pred),
2               columns = ['Forecast 0', 'Forecast 1'],
```

```
3 index = ['Actual 0', 'Actual 1'])
```

Как мы видим, хотя веса модели и предсказанные вероятности немного отличаются, ее точность осталась неизменной.

Решающая граница

Построим решающую границу.

```
1 # найдем сдвиг и наклон для уравнения решающей границы
2 c, m = -model.intercept_ / model.coef_[0][1], -model.coef_[0][0] / model
3 c, m
```

```
1 (array([0.12536569]), -0.8521526076691505)
```

```
1 # посмотрим на линию решающей границы
2 plt.figure(figsize = (11, 9))
3
4 xmin, xmax = min(X['alcohol']), max(X['alcohol'])
5 ymin, ymax = min(X['proline']), max(X['proline'])
6 xd = np.array([xmin, xmax])
7 yd = m*xd + c
8 plt.plot(xd, yd, 'k', lw=1, ls='--')
9 sns.scatterplot(x = X['alcohol'], y = X['proline'], hue = y, s = 70)
10 plt.legend(loc = 'upper left', prop = {'size': 15})
11
12 plt.xlim(xmin, xmax)
13 plt.ylim(ymin, ymax)
14
15 plt.show()
```

Бинарная полиномиальная регрессия

Идея **бинарной полиномиальной логистической регрессии** (binary polynomial logistic regression) заключается в том, чтобы использовать полином внутри сигмоиды и соответственно создать нелинейную границу между двумя классами.

Полиномиальные признаки

Уравнение полинома на основе двух признаков будет выглядеть следующим образом.

$$y = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1 x_2 + \theta_5 x_2^2$$

Реализуем этот алгоритм на практике и посмотрим, улучшатся ли результаты. Вначале, подготовим и масштабируем данные.

```
1 X = df[['alcohol', 'proline']]
2 y = df['target']
3
4 X = (X - X.mean())/X.std()
```

Теперь преобразуем наши данные так, как если бы мы использовали полином второй степени.

Смысл создания полиномиальных признаков мы детально разобрали на занятии по множественной линейной регрессии.

```
1 # создадим объект класса PolynomialFeatures
2 # укажем, что мы хотим создать полином второй степени
3 polynomial_features = PolynomialFeatures(degree = 2)
4
5 # преобразуем данные с помощью метода .fit_transform()
6 X_poly = polynomial_features.fit_transform(X)
```

Сравним исходные признаки с полиномиальными.

```
1 # посмотрим на первое наблюдение
2 X.head(1)
```

```
1 # должно получиться шесть признаков
2 X_poly[:1]
```

```
1 array([[1.          , 1.44685785, 0.77985116, 2.09339765, 1.12833378,
2         0.60816783]])
```

Моделирование и оценка качества

Обучим модель, сделаем прогноз и оценим результат.

```
1 # создадим объект класса LogisticRegression
2 poly_model = LogisticRegression()
3
4 # обучим модель на полиномиальных признаках
5 poly_model = poly_model.fit(X_poly, y)
6
7 # сделаем прогноз
8 y_pred = poly_model.predict(X_poly)
9
10 # рассчитаем accuracy
11 accuracy_score(y_pred, y)
```

```
1 0.9615384615384616
```

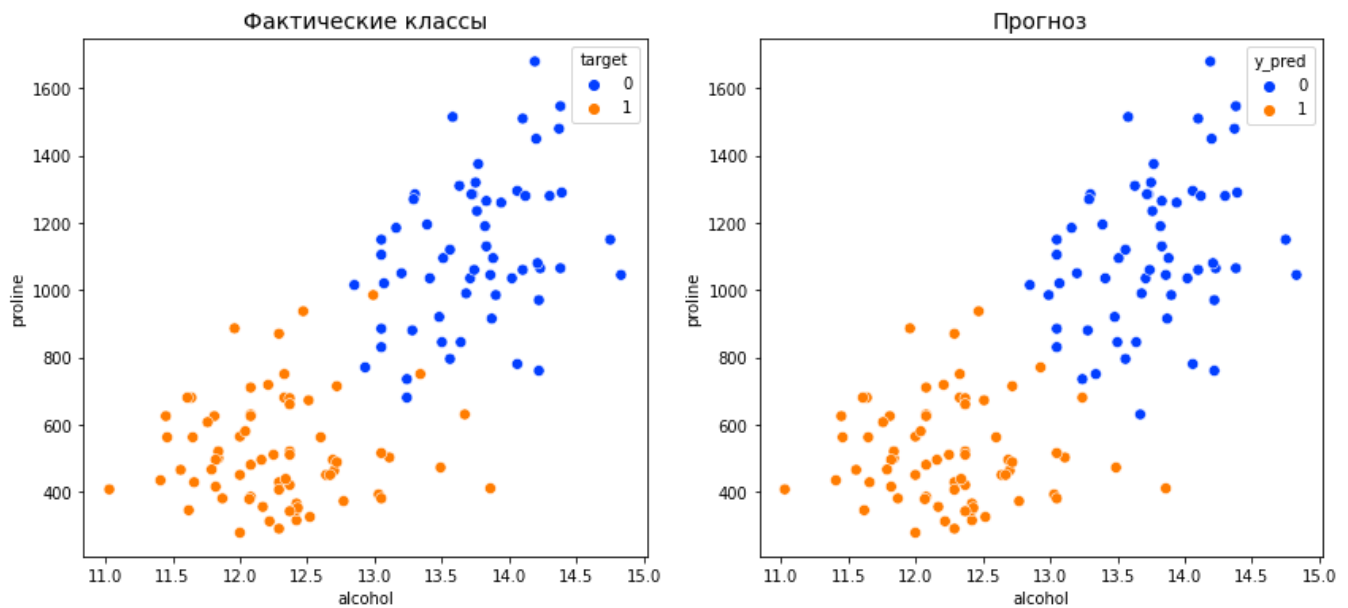
Построим матрицу ошибок.

```
1 pd.DataFrame(confusion_matrix(y, y_pred),
2               columns = ['Forecast 0', 'Forecast 1'],
3               index = ['Actual 0', 'Actual 1'])
```

Для того чтобы визуально оценить качество модели, построим два графика: фактических классов и прогнозных. Вначале создадим датасет, в котором будут исходные признаки (alcohol, proline) и прогнозные значения (y_pred).

```
1 # сделаем копию исходного датафрейма с нужными признаками
2 predictions = df[['alcohol', 'proline']].copy()
3
4 # и добавим новый столбец с прогнозными значениями
5 predictions['y_pred'] = y_pred
6
7 # посмотрим на результат
8 predictions.head(3)
```

```
1 # создадим два подграфика с помощью функции plt.subplots()
2 # расположим подграфики на одной строке
3 fig, (ax1, ax2) = plt.subplots(1, 2,
4                                # пропишем размер,
5                                figsize = (14, 6),
6                                # а также расстояние между подграфиками
7                                gridspec_kw = {'wspace' : 0.2})
8
9 # на левом подграфике выведем фактические классы
10 sns.scatterplot(data = df, x = 'alcohol', y = 'proline', hue = 'target',
11                ax1.set_title('Фактические классы', fontsize = 14))
12
13 # на правом - прогнозные
14 sns.scatterplot(data = predictions, x = 'alcohol', y = 'proline', hue =
15                ax2.set_title('Прогноз', fontsize = 14))
16
17 # зададим общий заголовок
18 fig.suptitle('Бинарная полиномиальная регрессия', fontsize = 16)
19
20 plt.show()
```



Как вы видите, нам не удалось добиться улучшения по сравнению с обычной полиномиальной регрессией.

Напомню, что создание подграфиков мы подробно разобрали на занятии по исследовательскому анализу данных.

В качестве упражнения предлагаю вам выяснить, какая степень полинома позволит улучшить результат прогноза на этих данных и насколько, таким образом, улучшится качество предсказаний.

Перейдем ко второй части нашего занятия.

Мультиклассовая логистическая регрессия

Как поступить, если нужно предсказать не два класса, а больше? Сегодня мы рассмотрим два подхода: one-vs-rest и кросс-энтропию. Начнем с того, что подготовим данные.

Подготовка данных

Вернем исходный датасет с тремя классами.


```

1 # вновь импортируем датасет о вине
2 data = datasets.load_wine()
3
4 # превратим его в датафрейм
5 df = pd.DataFrame(data.data, columns = data.feature_names)
6
7 # приведем признаки к одному масштабу
8 df = (df - df.mean())/df.std()
9
10 # добавим целевую переменную
11 df['target'] = data.target
12
13 # убедимся, что у нас присутствуют все три класса
14 df.target.value_counts()

```

```

1 1    71
2 0    59
3 2    48
4 Name: target, dtype: int64

```

В целевой переменной большое двух классов, а значит точечно-бисериальный коэффициент корреляции мы использовать не можем. Воспользуемся корреляционным отношением (correlation ratio).

```

1 # код ниже был подробно разобран на предыдущем занятии
2 def correlation_ratio(numerical, categorical):
3
4     values = np.array(numerical)
5     ss_total = np.sum((values.mean() - values) ** 2)
6
7     cats = np.unique(categorical, return_inverse = True)[1]
8
9     ss_betweengroups = 0
10
11     for c in np.unique(cats):
12
13         group = values[np.argwhere(cats == c).flatten()]
14         ss_betweengroups += len(group) * (group.mean() - values.mean()) ** 2
15
16     return np.sqrt(ss_betweengroups/ss_total)

```

```

1 # создадим два списка, один для названий признаков, второй для значений
2 columns, correlations = [], []
3
4 # пройдемся по всем столбцам датафрейма кроме целевой переменной
5 for col in df.drop('target', axis = 1).columns:
6     # поместим название признака в список columns
7     columns.append(col)
8     # рассчитаем взаимосвязь этого признака с целевой переменной
9     # и поместим результат в список значений корреляционного отношения

```

```

10     correlations.append(correlation_ratio(df[col], df['target']))
11
12     # создадим датафрейм на основе заполненных списков
13     # и применим градиентную цветовую схему
14     pd.DataFrame({'column': columns, 'correlation': correlations}).style.back

```

| | column | correlation |
|----|------------------------------|-------------|
| 0 | alcohol | 0.779024 |
| 1 | malic_acid | 0.544857 |
| 2 | ash | 0.363394 |
| 3 | alcalinity_of_ash | 0.538689 |
| 4 | magnesium | 0.352680 |
| 5 | total_phenols | 0.719163 |
| 6 | flavanoids | 0.853098 |
| 7 | nonflavanoid_phenols | 0.489519 |
| 8 | proanthocyanins | 0.506986 |
| 9 | color_intensity | 0.761353 |
| 10 | hue | 0.732522 |
| 11 | od280/od315_of_diluted_wines | 0.827438 |
| 12 | proline | 0.838935 |

Теперь наибольшую корреляцию с целевой переменной показывают флавоноиды (flavanoids) и пролин (proline). Их и оставим.

```

1 df = df[['flavanoids', 'proline', 'target']].copy()
2 df.head(3)

```

| | flavanoids | proline | target |
|---|------------|----------|--------|
| 0 | 1.031908 | 1.010159 | 0 |
| 1 | 0.731565 | 0.962526 | 0 |
| 2 | 1.212114 | 1.391224 | 0 |

Посмотрим, насколько легко можно разделить эти классы.

```

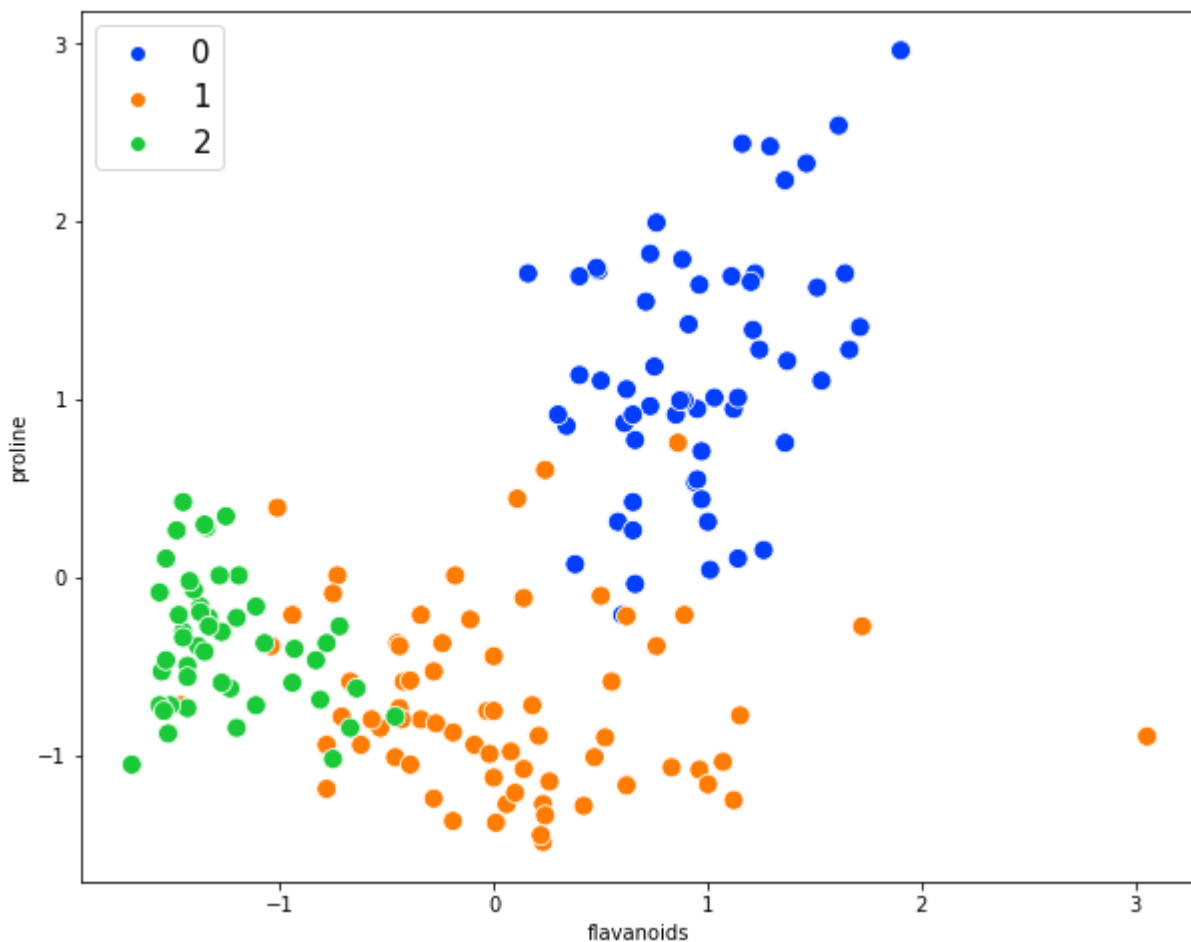
1 # зададим размер графика
2 plt.figure(figsize = (10, 8))
3 # построим точечную диаграмму с двумя признаками, разделяющей категориаль

```

```

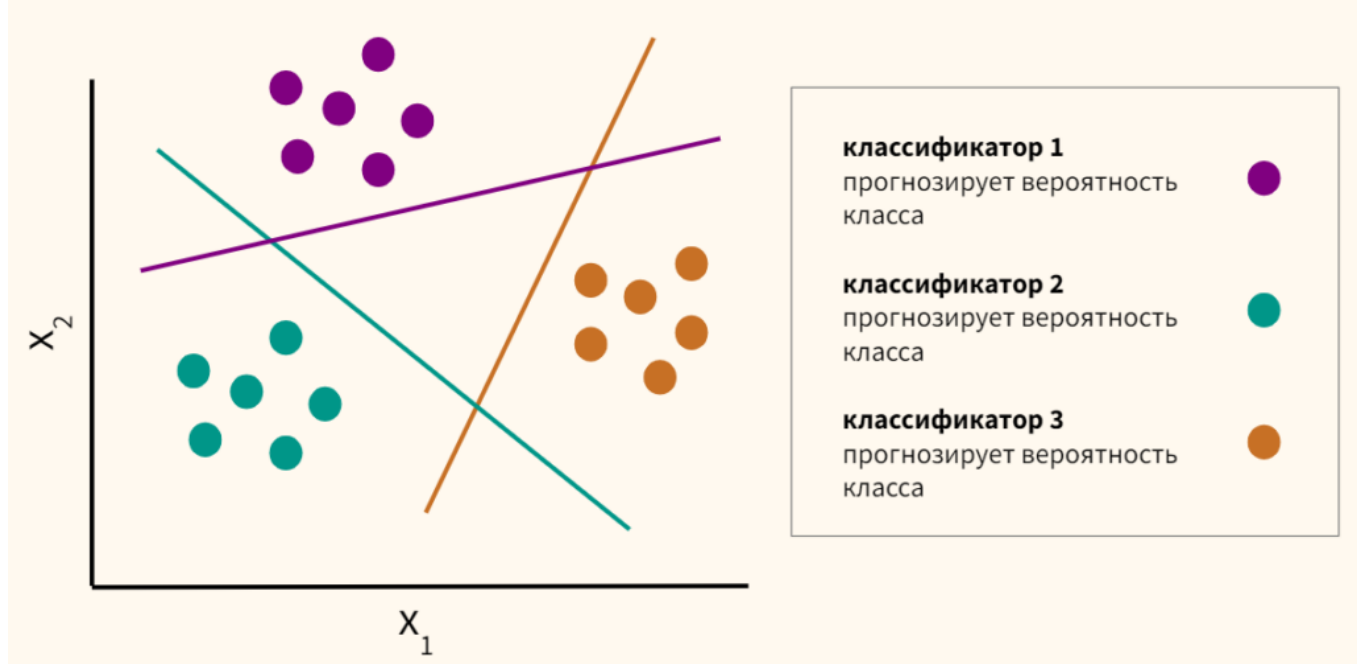
4 sns.scatterplot(x = df.flavanoids, y = df.proline, hue = df.target, palette='magma')
5 # добавим легенду
6 plt.legend(loc = 'upper left', prop = {'size': 15})
7 plt.show()

```



Перейдем непосредственно к алгоритмам мультиклассовой логистической регрессии. Начнем с подхода one-vs-rest.

Подход one-vs-rest



Подход one-vs-rest или one-vs-all предполагает, что мы отделяем один класс, а остальные наоборот объединяем. Так мы поступаем с каждым классом и строим по одной модели логистической регрессии относительно каждого из класса. Например, если у нас три класса, то у нас будет три модели логистической регрессии. Далее мы смотрим на получившиеся вероятности и выбираем наибольшую.

$$h_{\theta}^{(i)}(x) = P(y = i|x; \theta), i \in \{0, 1, 2\}$$

При таком подходе сам по себе алгоритм логистической регрессии претерпевает лишь несущественные изменения, главное правильно подготовить данные для обучения модели.

Подготовка датасетов

```
1 # поместим признаки и данные в соответствующие переменные
2 x1, x2 = df.columns[0], df.columns[1]
3
4 target = df.target.unique()
5 target
```

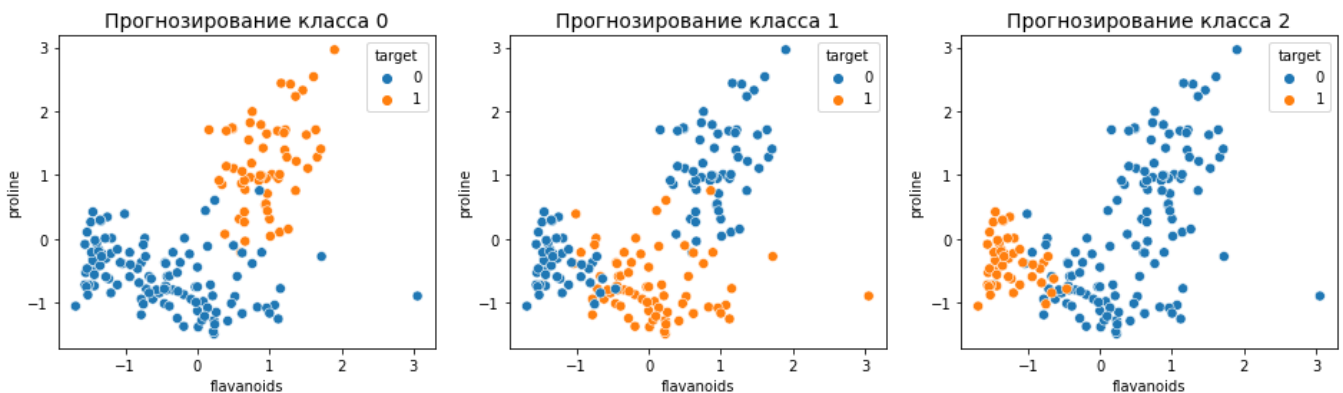
```
1 array([0, 1, 2])
```

```
1 # сделаем копии датафреймов
2 ovr_0, ovr_1, ovr_2 = df.copy(), df.copy(), df.copy()
3
4 # в каждом из них сделаем целевым классом 0-й, 1-й или 2-й классы
5 # например, в ovr_0 первым классом будет класс 0, а классы 1 и 2 - нулевые
6 ovr_0['target'] = np.where(df['target'] == target[0], 1, 0)
```

```

7  ovr_1['target'] = np.where(df['target'] == target[1], 1, 0)
8  ovr_2['target'] = np.where(df['target'] == target[2], 1, 0)
9
10 # выведем разделение на классы на графике
11 fig, (ax1, ax2, ax3) = plt.subplots(1, 3,
12                                     figsize = (16, 4),
13                                     gridspec_kw = {'wspace': 0.2, 'hspace': 0.2})
14
15 sns.scatterplot(data = ovr_0, x = x1, y = x2, hue = 'target', s = 50, ax=ax1)
16 ax1.set_title('Прогнозирование класса 0', fontsize = 14)
17
18 sns.scatterplot(data = ovr_1, x = x1, y = x2, hue = 'target', s = 50, ax=ax2)
19 ax2.set_title('Прогнозирование класса 1', fontsize = 14)
20
21 sns.scatterplot(data = ovr_2, x = x1, y = x2, hue = 'target', s = 50, ax=ax3)
22 ax3.set_title('Прогнозирование класса 2', fontsize = 14)
23
24 plt.show()

```



Обучение моделей

```

1  models = []
2
3  # поочередно обучим каждую из моделей
4  for ova_n in [ovr_0, ovr_1, ovr_2]:
5      X = ova_n[['flavanoids', 'proline']]
6      y = ova_n['target']
7
8      model = LogReg()
9      model.fit(X, y)
10
11     # каждую обученную модель поместим в список
12     models.append(model)

```

```

1  # убедимся, что все работает
2  # например, выведем коэффициенты модели 1
3  models[0].thetas

```

```
1 array([-0.99971466,  1.280398   ,  2.04834457])
```

Прогноз и оценка качества

```
1 # вновь перенесем данные из исходного датафрейма
2 X = df[['flavanoids', 'proline']]
3 y = df['target']
4
5 # в список probs будем записывать результат каждой модели
6 # для каждого наблюдения
7 probs = []
8
9 for model in models:
10     _, prob = model.predict(X)
11     probs.append(prob)
```

```
1 # очевидно, для каждого наблюдения у нас будет три вероятности
2 # принадлежности к целевому классу
3 probs[0][0], probs[1][0], probs[2][0]
```

```
1 (0.9161148288779738, 0.1540913395345091, 0.026621132600103174)
```

```
1 # склеим и изменим размерность массива таким образом, чтобы
2 # строки были наблюдениями, а столбцы вероятностями
3 all_probs = np.concatenate(probs, axis = 0).reshape(len(probs), -1).T
4 all_probs.shape
```

```
1 # каждая из 178 строк - это вероятность одного наблюдения
2 # принадлежать к классу 0, 1, 2
3 all_probs[0]
```

```
1 array([0.91611483, 0.15409134, 0.02662113])
```

Обратите внимание, при использовании подхода one-vs-rest вероятности в сумме не дают единицу.

```
1 # например, первое наблюдение вероятнее всего принадлежит к классу 0
2 np.argmax(all_probs[0])
```

```
1
```

```
1 # найдем максимальную вероятность в каждой строке,
2 # индекс вероятности [0, 1, 2] и будет прогнозом
3 y_pred = np.argmax(all_probs, axis = 1)
4
5 # рассчитаем accuracy
6 accuracy_score(y, y_pred)
```

```
1 0.9157303370786517
```

```

1 # выведем матрицу ошибок
2 pd.DataFrame(confusion_matrix(y, y_pred),
3               columns = ['Forecast 0', 'Forecast 1', 'Forecast 2'],
4               index = ['Actual 0', 'Actual 1', 'Actual 2'])

```

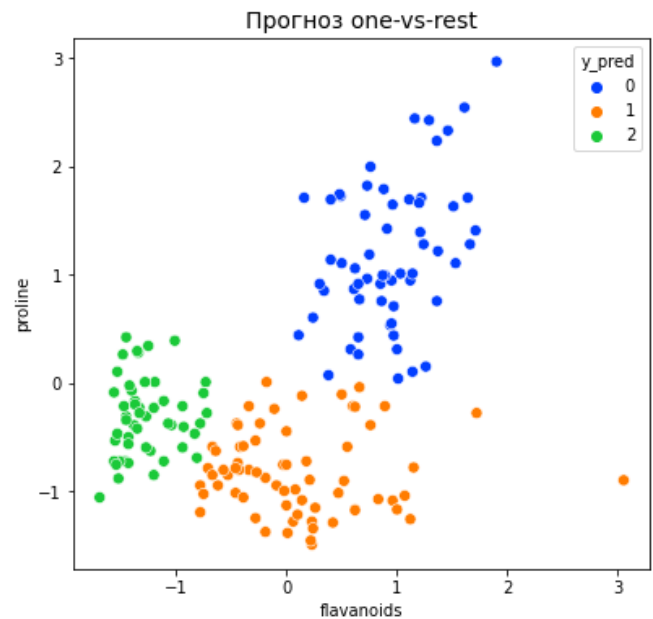
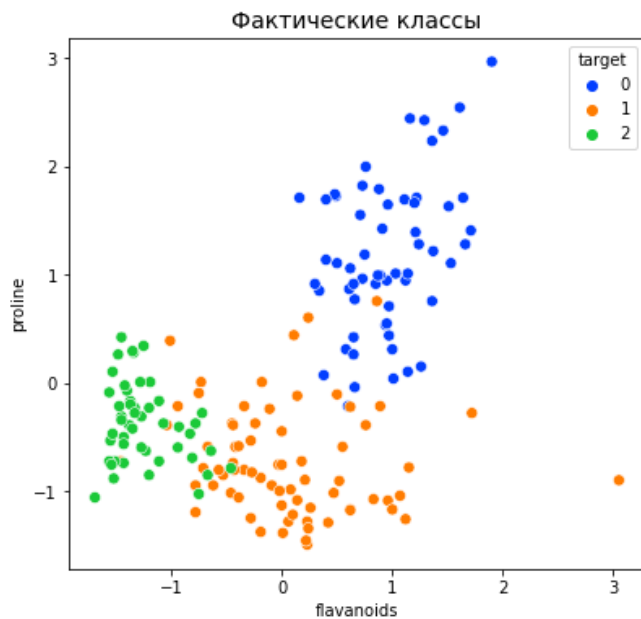
| | Forecast 0 | Forecast 1 | Forecast 2 |
|----------|------------|------------|------------|
| Actual 0 | 57 | 2 | 0 |
| Actual 1 | 3 | 62 | 6 |
| Actual 2 | 0 | 4 | 44 |

Сравним фактическое и прогнозное распределение классов на точечной диаграмме.

```

1 predictions = df[['flavanoids', 'proline']].copy()
2 predictions['y_pred'] = y_pred
3
4 fig, (ax1, ax2) = plt.subplots(1, 2,
5                               figsize = (14, 6),
6                               gridspec_kw = {'wspace': 0.2, 'hspace': 0.2})
7
8 sns.scatterplot(data = df, x = 'flavanoids', y = 'proline', hue = 'target',
9                 palette = 'bright', s = 50, ax = ax1)
10
11 ax1.set_title('Фактические классы', fontsize = 14)
12
13 sns.scatterplot(data = predictions, x = 'flavanoids', y = 'proline', hue = 'y_pred',
14                 palette = 'bright', s = 50, ax = ax2)
15
16 ax2.set_title('Прогноз one-vs-rest', fontsize = 14)
17
18 plt.show()

```



Написание класса

Поместим достигнутый выше результат в класс.

```

1 class OVR_LogReg():
2
3     def __init__(self):
4         self.models_thetas = []
5         self.models_loss = []
6
7     def fit(self, x, y, iter = 20000, learning_rate = 0.001):
8
9         dfs = self.preprocess(x, y)
10
11         models_thetas, models_loss = [], []
12
13         for ovr_df in dfs:
14
15             x = ovr_df.drop('target', axis = 1).copy()
16             y = ovr_df.target.copy()
17
18             self.add_ones(x)
19
20             loss_history = []
21             thetas, n = np.zeros(x.shape[1]), x.shape[0]
22
23             for i in range(iter):
24                 y_pred = self.h(x, thetas)
25                 loss_history.append(self.objective(y, y_pred))
26                 grad = self.gradient(x, y, y_pred, n)
27                 thetas -= learning_rate * grad
28

```



```

29     models_thetas.append(thetas)
30     models_loss.append(loss_history)
31
32     self.models_thetas = models_thetas
33     self.models_loss = models_loss
34
35 def predict(self, x):
36     x = x.copy()
37     probs = []
38     self.add_ones(x)
39     for t in self.models_thetas:
40         z = np.dot(x, t)
41         prob = np.array([self.stable_sigmoid(value) for value in z])
42         probs.append(prob)
43
44     all_probs = np.concatenate(probs, axis = 0).reshape(len(probs), -1)
45     y_pred = np.argmax(all_probs, axis = 1)
46
47     return y_pred, all_probs
48
49 def preprocess(self, x, y):
50
51     x, y = x.copy(), y.copy()
52
53     x['target'] = y
54     classes = x.target.unique()
55
56     dfs = []
57     ovr_df = None
58
59     for c in classes:
60         ovr_df = x.drop('target', axis = 1).copy()
61         ovr_df['target'] = np.where(x['target'] == classes[c], 1, 0)
62         dfs.append(ovr_df)
63
64     return dfs
65
66 def add_ones(self, x):
67     return x.insert(0, 'x0', np.ones(x.shape[0]))
68
69 def h(self, x, thetas):
70     z = np.dot(x, thetas)
71     return np.array([self.stable_sigmoid(value) for value in z])
72
73 def objective(self, y, y_pred):
74     y_one_loss = y * np.log(y_pred + 1e-9)
75     y_zero_loss = (1 - y) * np.log(1 - y_pred + 1e-9)
76     return -np.mean(y_zero_loss + y_one_loss)
77
78 def gradient(self, x, y, y_pred, n):
79     return np.dot(x.T, (y_pred - y)) / n

```

```

80
81     def stable_sigmoid(self, z):
82         if z >= 0:
83             return 1 / (1 + np.exp(-z))
84         else:
85             return np.exp(z) / (np.exp(z) + 1)

```

Проверим класс в работе.

```

1  X = df[['flavanoids', 'proline']]
2  y = df['target']
3
4  model = OVR_LogReg()
5  model.fit(X, y)
6
7  y_pred, probs = model.predict(X)
8  accuracy_score(y_pred, y)

```

```

1  0.9157303370786517

```

Сравнение с sklearn

Для того чтобы применить подход one-vs-rest в классе LogisticRegression, необходимо использовать значение параметра multi_class = 'ovr'.

```

1  X = df[['flavanoids', 'proline']]
2  y = df['target']
3
4  ovr_model = LogisticRegression(multi_class = 'ovr')
5  ovr_model = ovr_model.fit(X, y)
6  y_pred = ovr_model.predict(X)
7
8  accuracy_score(y_pred, y)

```

```

1  0.898876404494382

```

Мультиклассовая полиномиальная регрессия

Как мы увидели в предыдущем разделе, линейная решающая граница допустила некоторое количество ошибок. Попробуем улучшить результат, применив мультиклассовую полиномиальную логистическую регрессию.

```

1  X = df[['flavanoids', 'proline']]
2  y = df['target']

```

```

3
4 polynomial_features = PolynomialFeatures(degree = 7)
5
6 X_poly = polynomial_features.fit_transform(X)
7
8 poly_ovr_model = LogisticRegression(multi_class = 'ovr')
9 poly_ovr_model = poly_ovr_model.fit(X_poly, y)
10 y_pred = poly_ovr_model.predict(X_poly)
11
12 accuracy_score(y_pred, y)

```

```

1 0.9157303370786517

```

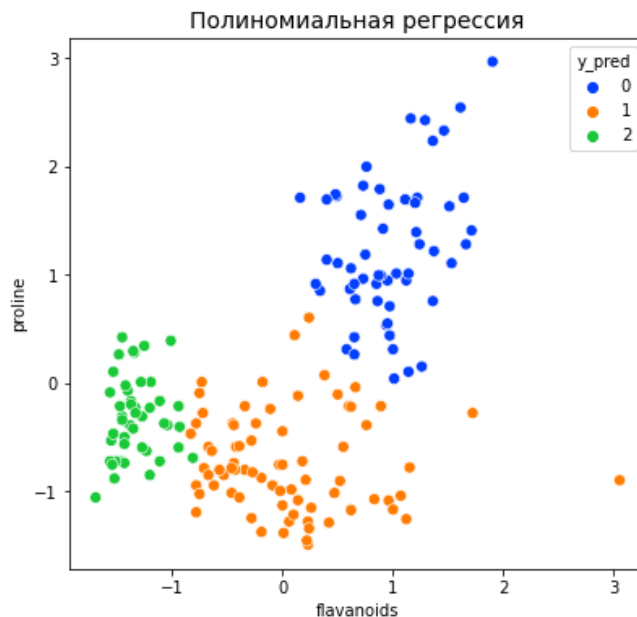
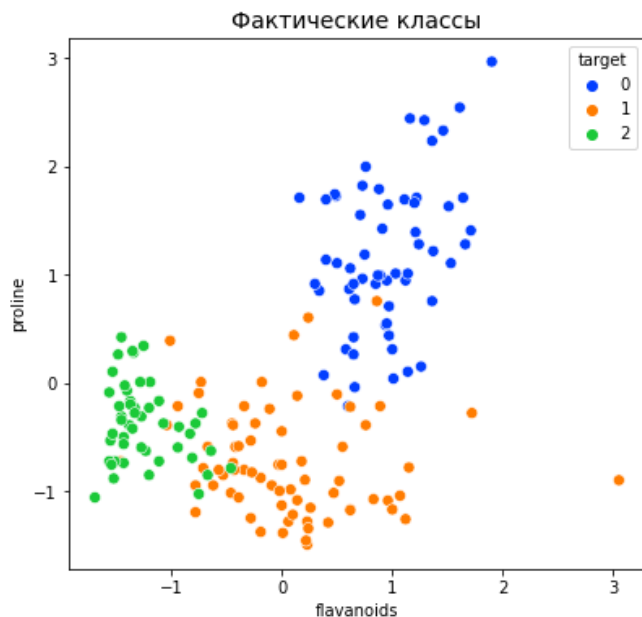
Результат, по сравнению с моделью sklearn без полиномиальных признаков, стал чуть лучше. Однако это было достигнуто за счет полинома достаточно высокой степени (degree = 7), что неэффективно с точки зрения временной сложности алгоритма.

Посмотрим, какие нелинейные решающие границы удалось построить алгоритму.

```

1 predictions = df[['flavanoids', 'proline']].copy()
2 predictions['y_pred'] = y_pred
3
4 fig, (ax1, ax2) = plt.subplots(1, 2,
5                               figsize = (14, 6),
6                               gridspec_kw = {'wspace': 0.2, 'hspace': 0.2})
7
8 sns.scatterplot(data = df, x = 'flavanoids', y = 'proline', hue = 'target')
9 ax1.set_title('Фактические классы', fontsize = 14)
10
11 sns.scatterplot(data = predictions, x = 'flavanoids', y = 'proline', hue = 'y_pred')
12 ax2.set_title('Полиномиальная регрессия', fontsize = 14)
13
14 plt.show()

```



Softmax Regression

Еще один подход при создании мультиклассовой логистической регрессии заключается в том, чтобы не разбивать многоклассовые данные на несколько датасетов и использовать бинарный классификатор, а сразу применять функции, которые подходят для работы с множеством классов.

Такую регрессию часто называют **Softmax Regression** из-за того, что в ней используется уже знакомая нам по занятию об основах нейросетей функция softmax. Вначале подготовим данные.

Подготовка признаков

Возьмем признаки flavanoids и proline и добавим столбец из единиц.

```
1 def add_ones(x):
2     # важно! метод .insert() изменяет исходный датафрейм
3     return x.insert(0, 'x0', np.ones(x.shape[0]))
```

```
1 X = df[['flavanoids', 'proline']]
2
3 add_ones(X)
4 X.head(3)
```

| | x0 | flavanoids | proline |
|---|-----|------------|----------|
| 0 | 1.0 | 1.031908 | 1.010159 |
| 1 | 1.0 | 0.731565 | 0.962526 |
| 2 | 1.0 | 1.212114 | 1.391224 |

Кодирование целевой переменной

Напишем собственную функцию для one-hot encoding.

```

1  def ohe(y):
2      # количество примеров и количество классов
3      examples, features = y.shape[0], len(np.unique(y))
4      # нулевая матрица: количество наблюдений x количество признаков
5      zeros_matrix = np.zeros((examples, features))
6      # построчно проходимся по нулевой матрице и с помощью индекса заполняем
7      for i, (row, digit) in enumerate(zip(zeros_matrix, y)):
8          zeros_matrix[i][digit] = 1
9
10     return zeros_matrix

```

```

1  y = df['target']
2
3  y_enc = ohe(df['target'])
4  y_enc[:3]

```

```

1  array([[1., 0., 0.],
2         [1., 0., 0.],
3         [1., 0., 0.]])

```

Такой же результат можно получить с помощью класса LabelBinarizer.

```

1  lb = LabelBinarizer()
2  lb.fit(y)
3  lb.classes_

```

```

1  array([0, 1, 2])

```

```

1  y_lb = lb.transform(y)
2  y_lb[:5]

```

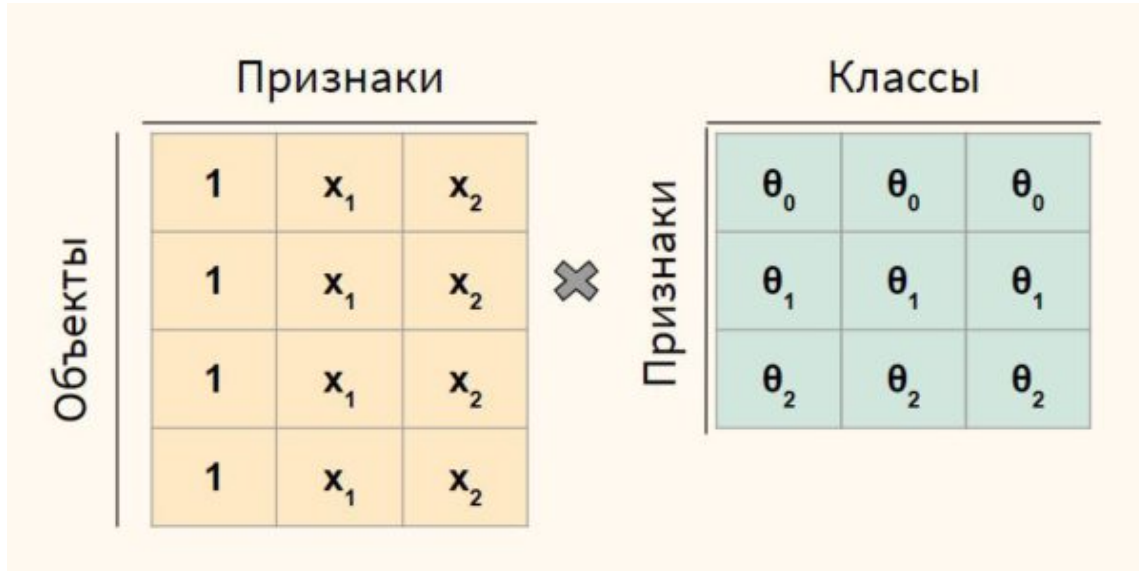
```

1  array([[1, 0, 0],
2         [1, 0, 0],
3         [1, 0, 0],
4         [1, 0, 0],
5         [1, 0, 0]])

```

Инициализация весов

Создадим нулевую матрицу весов. Она будет иметь размерность: количество признаков (строки) x количество классов (столбцы). Приведем схематичный пример для четырех наблюдений, трех признаков (включая сдвиг θ_0) и трех классов.



Инициализируем веса.

```
1 thetas = np.zeros((3, 3))
2 thetas
```

```
1 array([[0., 0., 0.],
2        [0., 0., 0.],
3        [0., 0., 0.]])
```

Функция softmax

Подробнее изучим функцию softmax. Приведем формулу.

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{k=1}^N e^{z_k}}$$

Рассмотрим ее реализацию на Питоне.

Напомню, что $z = (-X\theta)$. Соответственно в нашем случае мы будем умножать матрицу 178 x 3 на 3 x 3.

```
1 X.shape, thetas.shape
```

```
1 | ((178, 3), (3, 3))
```

В результате получим матрицу 178 x 3, где каждая строка — это прогнозные значения принадлежности одного наблюдения к каждому из трех классов.

```
1 | z = np.dot(-X, thetas)
2 | z.shape
```

```
1 | (178, 3)
```

Так как мы умножаем на ноль, при первой итерации эти значения будут равны нулю.

```
1 | z[:5]
```

```
1 | array([[0., 0., 0.],
2 |        [0., 0., 0.],
3 |        [0., 0., 0.],
4 |        [0., 0., 0.],
5 |        [0., 0., 0.]])
```

Для того чтобы обеспечить вычислительную устойчивость softmax мы можем вычесть из z максимальное значение в каждой из 178 строк (пока что, опять же на первой итерации, оно равно нулю).

$$\text{softmax}(z)_i = \frac{e^{z_i - \max(z)}}{\sum_{k=1}^N e^{z_k - \max(z)}}$$

```
1 | # axis = -1 - это последняя ось
2 | # keepdims = True сохраняет размерность (в данном случае двумерный массив)
3 | np.max(z, axis = -1, keepdims = True)[:5]
```

```
1 | array([[0.],
2 |        [0.],
3 |        [0.],
4 |        [0.],
5 |        [0.]])
```

```
1 | z = z - np.max(z, axis = -1, keepdims = True)
2 | z[:5]
```

```
1 | array([[0., 0., 0.],
2 |        [0., 0., 0.],
3 |        [0., 0., 0.],
4 |        [0., 0., 0.],
5 |        [0., 0., 0.]])
```

Смысл такого преобразования в том, что оно делает значения z нулевыми или отрицательными.

```
1 arr = np.array([-2, 3, 0, -7, 6])
2 arr - max(arr)
```

```
1 array([-8, -3, -6, -13,  0])
```

Далее, число возводимое в увеличивающуюся отрицательную степень стремится к нулю, а не к бесконечности и, таким образом, не вызывает переполнения памяти. Найдем числитель и знаменатель из формулы softmax.

```
1 numerator = np.exp(z)
2 numerator[:5]
```

```
1 array([[1., 1., 1.],
2        [1., 1., 1.],
3        [1., 1., 1.],
4        [1., 1., 1.],
5        [1., 1., 1.]])
```

```
1 denominator = np.sum(numerator, axis = -1, keepdims = True)
2 denominator[:5]
```

```
1 array([[3.],
2        [3.],
3        [3.],
4        [3.],
5        [3.]])
```

Разделим числитель и знаменатель и, таким образом, вычислим вероятность принадлежности каждого из наблюдений (строки результата) к одному из трех классов (столбцы).

```
1 softmax = numerator / denominator
2 softmax[:5]
```

```
1 array([[0.33333333, 0.33333333, 0.33333333],
2        [0.33333333, 0.33333333, 0.33333333],
3        [0.33333333, 0.33333333, 0.33333333],
4        [0.33333333, 0.33333333, 0.33333333],
5        [0.33333333, 0.33333333, 0.33333333]])
```

На первой итерации при одинаковых θ мы получаем, что логично, одинаковые вероятности принадлежности к каждому из классов. Напишем функцию.

```
1 def stable_softmax(x, thetas):
2     z = np.dot(-x, thetas)
3     z = z - np.max(z, axis = -1, keepdims = True)
```



```

4     numerator = np.exp(z)
5     denominator = np.sum(numerator, axis = -1, keepdims = True)
6     softmax = numerator / denominator
7     return softmax

```

```

1     probs = stable_softmax(X, thetas)
2     probs[:3]

```

```

1     array([[0.33333333, 0.33333333, 0.33333333],
2           [0.33333333, 0.33333333, 0.33333333],
3           [0.33333333, 0.33333333, 0.33333333]])

```

Примечание. Обратите внимание, что сигмоида — это частный случай функции softmax для двух классов $[z_1, 0]$. Вероятность класса z_1 будет равна

$$\text{softmax}(z_1) = \frac{e^{z_1}}{e^{z_1} + e^0} = \frac{e^{z_1}}{e^{z_1} + 1}$$

Если разделить и числитель, и знаменатель на e^{z_1} , то получим

$$\text{sigmoid}(z_1) = \frac{e^{z_1}}{1 + e^{-z_1}}$$

Вычислять вероятность принадлежности ко второму классу нет необходимости, достаточно вычесть результат сигмоиды из единицы.

Теперь нужно понять, насколько сильно при таких весах ошибается наш алгоритм.

Функция потерь

Вспомним функцию бинарной кросс-энтропии. То есть функции ошибки для двух классов.

$$L(y, \theta) = -\frac{1}{n} \sum y \cdot \log(h_{\theta}(x)) + (1 - y) \cdot \log(1 - h_{\theta}(x))$$

Напомню, что у выступает в роли своего рода переключателя, сохраняющего одну из частей выражения, и обнуляющего другую. Теперь посмотрите на функцию категориальной (многоклассовой) кросс-энтропии (categorical cross-entropy).

$$L(y_{\text{one}}, \text{softmax}) = - \sum y_{\text{one}} \log(\text{softmax})$$

Разберемся, что здесь происходит. y_{ohe} содержит закодированную целевую переменную, например, для наблюдения класса 0 [1, 0, 0], softmax содержит вектор вероятностей принадлежности наблюдения к каждому из классов, например, [0,3 0,4 0,3] (мы видим, что алгоритм ошибается).

В данном случае закодированная целевая переменная также выступает в виде переключателя. Здесь при умножении «срабатывает» только первая вероятность $1 \times 0,3 + 0 \times 0,4 + 0 \times 0,4$. Если подставить в формулу, то получаем (np.sum() добавлена для сохранения единообразия с формулой выше, в данном случае у нас одно наблюдение и сумма не нужна).

```
1 y_ohe = np.array([1, 0, 0])
2 softmax = np.array([0.3, 0.4, 0.4])
3
4 -np.sum(y_ohe * np.log(softmax))
```

```
1 1.2039728043259361
```

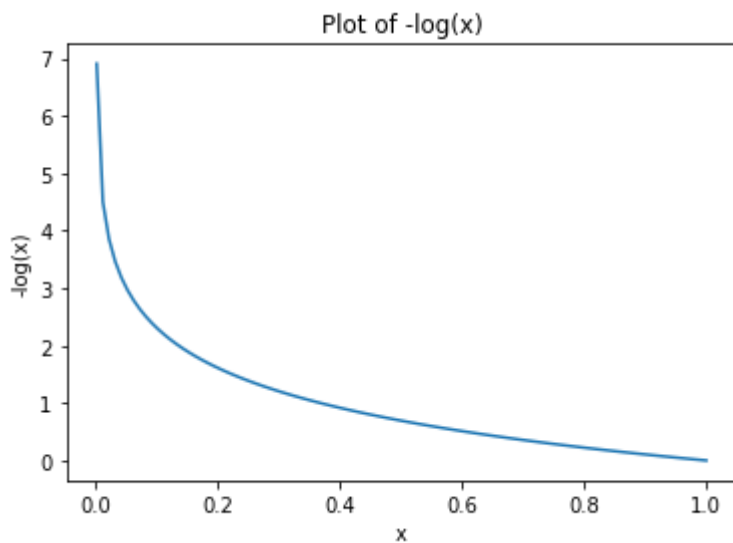
Если бы модель в своих вероятностях ошибалась меньше, то и общая ошибка была бы меньше.

```
1 y_ohe = np.array([1, 0, 0])
2 softmax = np.array([0.4, 0.3, 0.4])
3
4 -np.sum(y_ohe * np.log(softmax))
```

```
1 0.916290731874155
```

Функция $-\log$ позволяет снижать ошибку при увеличении вероятности верного (сохраненного переключателем) класса.

```
1 x_arr = np.linspace(0.001,1, 100)
2 sns.lineplot(x=x_arr,y=-np.log(x_arr))
3 plt.title('Plot of -log(x)')
4 plt.xlabel('x')
5 plt.ylabel('-log(x)');
```



Напишем функцию.

```
1 # добавим константу в логарифм для вычислительной устойчивости
2 def cross_entropy(probs, y_enc, epsilon = 1e-9):
3     n = probs.shape[0]
4     ce = -np.sum(y_enc * np.log(probs + epsilon)) / n
5     return ce
```

Рассчитаем ошибку для нулевых весов.

```
1 ce = cross_entropy(probs, y_enc)
2 ce
```

```
1 1.0986122856681098
```

Для снижения ошибки нужно найти градиент.

Градиент

Приведем формулу градиента без дифференцирования.

$$\nabla_{\theta} J = \frac{1}{n} \times X^T \cdot (y_{ohc} - softmax)$$

По сути, мы умножаем транспонированную матрицу признаков (3 x 178) на разницу между закодированной целевой переменной и вероятностями функции softmax (178 x 3).

```
1 def gradient_softmax(X, probs, y_enc):
2     # если не добавить функцию np.array(), будет выводиться датафрейм
3     return np.array(1 / probs.shape[0] * np.dot(X.T, (y_enc - probs)))
```

```
1 gradient_softmax(X, probs, y_enc)

1 array([[ -0.00187266,  0.06554307, -0.06367041],
2         [ 0.31627721,  0.02059572, -0.33687293],
3         [ 0.38820566, -0.28801792, -0.10018774]])
```

Обучение модели, прогноз и оценка качества

Выполним обучение модели.

```
1 loss_history = []
2
3 # в цикле
4 for i in range(30000):
5     # рассчитаем прогнозное значение с текущими весами
6     probs = stable_softmax(X, thetas)
7     # посчитаем уровень ошибки при текущем прогнозе
8     loss_history.append(cross_entropy(probs, y_enc, epsilon = 1e-9))
9     # рассчитаем градиент
10    grad = gradient_softmax(X, probs, y_enc)
11    # используем градиент для улучшения весов модели
12    thetas = thetas - 0.002 * grad
```

Посмотрим на получившиеся коэффициенты (напомню, что первая строка матрицы это сдвиг (intercept, θ_0)) и достигнутый уровень ошибки.

```
1 thetas

1 array([[ 0.11290134, -0.90399727,  0.79109593],
2        [-1.7550965 , -0.7857371 ,  2.5408336 ],
3        [-1.93839311,  1.77140542,  0.16698769]])
```

```
1 loss_history[0], loss_history[-1]
```

```
1 (1.0986122856681098, 0.2569641080523888)
```

Сделаем прогноз и оценим качество.

```
1 y_pred = np.argmax(stable_softmax(X, thetas), axis = 1)
```

```
1 accuracy_score(y, y_pred)
```

```
1 0.9044943820224719
```

```
1 pd.DataFrame(confusion_matrix(y, y_pred),
2               columns = ['Forecast 0', 'Forecast 1', 'Forecast 2'],
3               index = ['Actual 0', 'Actual 1', 'Actual 2'])
```

| | Forecast 0 | Forecast 1 | Forecast 2 |
|----------|------------|------------|------------|
| Actual 0 | 56 | 3 | 0 |
| Actual 1 | 3 | 62 | 6 |
| Actual 2 | 0 | 5 | 43 |

Написание класса

Объединим созданные выше компоненты в класс.

```

1  class SoftmaxLogReg():
2
3      def __init__(self):
4          self.loss_ = None
5          self.thetas_ = None
6
7      def fit(self, x, y, iter = 30000, learning_rate = 0.002):
8
9          loss_history = []
10
11         self.add_ones(x)
12         y_enc = self.ohe(y)
13         thetas = np.zeros((x.shape[1], y_enc.shape[1]))
14
15         for i in range(iter):
16             probs = self.stable_softmax(x, thetas)
17             loss_history.append(self.cross_entropy(probs, y_enc, epsilon = 1e-9))
18             grad = self.gradient_softmax(x, probs, y_enc)
19             thetas = thetas - 0.002 * grad
20
21             self.thetas_ = thetas
22             self.loss_ = loss_history
23
24         def predict(self, x, y):
25             return np.argmax(self.stable_softmax(x, thetas), axis = 1)
26
27         def stable_softmax(self, x, thetas):
28             z = np.dot(-x, thetas)
29             z = z - np.max(z, axis = -1, keepdims = True)
30             numerator = np.exp(z)
31             denominator = np.sum(numerator, axis = -1, keepdims = True)
32             softmax = numerator / denominator
33             return softmax
34
35         def cross_entropy(self, probs, y_enc, epsilon = 1e-9):

```

```

36     n = probs.shape[0]
37     ce = -np.sum(y_enc * np.log(probs + epsilon)) / n
38     return ce
39
40     def gradient_softmax(self, x, probs, y_enc):
41         return np.array(1 / probs.shape[0] * np.dot(x.T, (y_enc - probs)))
42
43     def add_ones(self, x):
44         return x.insert(0, 'x0', np.ones(x.shape[0]))
45
46     def ohe(self, y):
47         examples, features = y.shape[0], len(np.unique(y))
48         zeros_matrix = np.zeros((examples, features))
49         for i, (row, digit) in enumerate(zip(zeros_matrix, y)):
50             zeros_matrix[i][digit] = 1
51
52     return zeros_matrix

```

Обучим модель, сделаем прогноз и оценим качество.

```

1  X = df[['flavanoids', 'proline']]
2  y = df['target']
3
4  model = SoftmaxLogReg()
5
6  model.fit(X, y)
7  model.thetas_, model.loss_[-1]

```

```

1  (array([[ 0.11290134, -0.90399727,  0.79109593],
2          [-1.7550965 , -0.7857371 ,  2.5408336 ],
3          [-1.93839311,  1.77140542,  0.16698769]]), 0.2569641080523888)

```

```

1  y_pred = model.predict(X, y)
2
3  accuracy_score(y, y_pred)

```

```

1  0.9044943820224719

```

Сравнение с sklearn

Для того чтобы использовать softmax логистическую регрессию в sklearn, соответствующему классу нужно передать параметр `multi_class = 'multinomial'`.

```

1  X = df[['flavanoids', 'proline']]
2  y = df['target']
3
4  # создадим объект класса LogisticRegression и запишем его в переменную r
5  model = LogisticRegression(multi_class = 'multinomial')

```

```
6
7 # обучим нашу модель
8 model.fit(X, y)
9
10 # посмотрим на получившиеся веса модели
11 model.intercept , model.coef
```

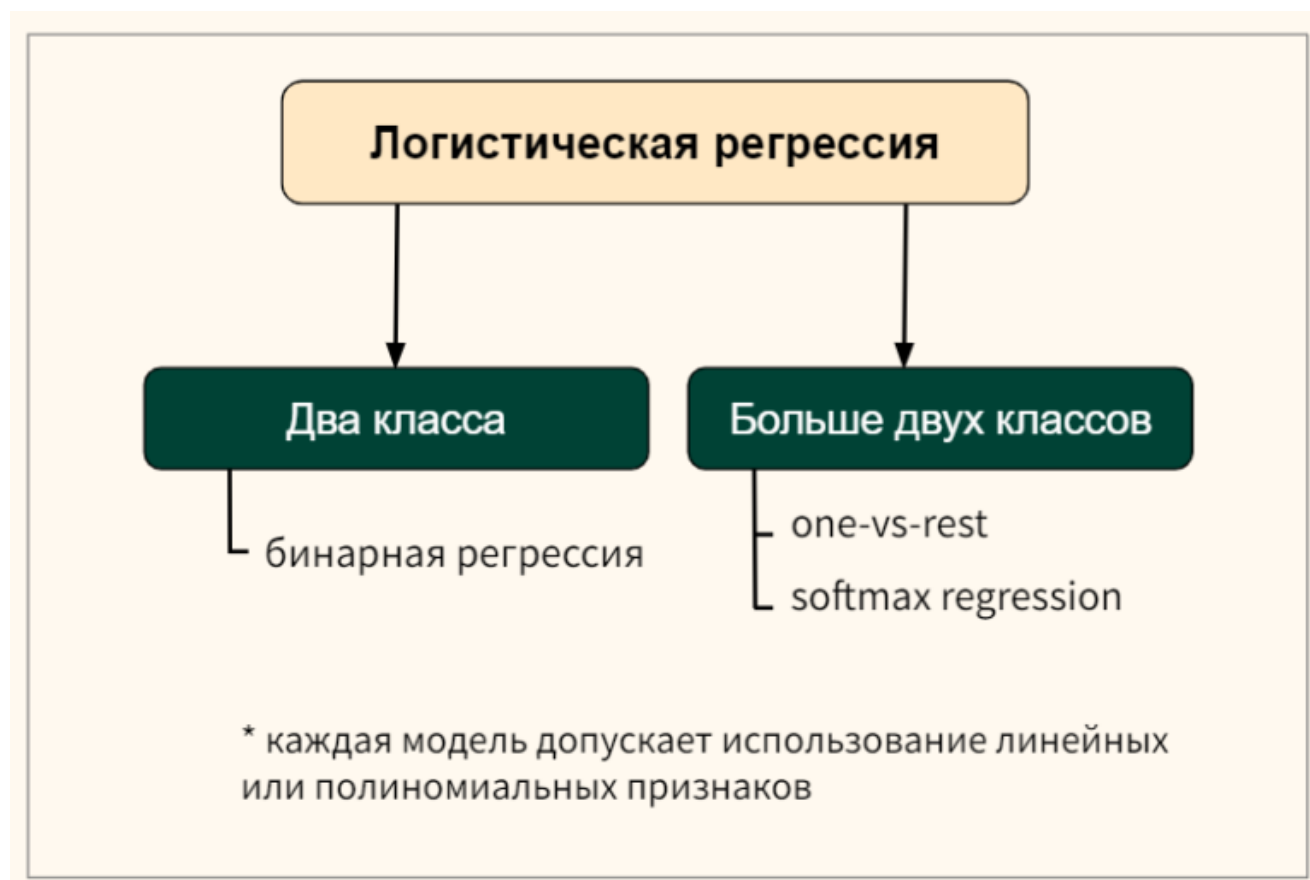
```
1 (array([ 0.09046097,  1.12593099, -1.21639196]),
2  array([[ 1.86357908,  1.89698292],
3         [ 0.86696131, -1.43973164],
4         [-2.73054039, -0.45725129]]))
```

```
1 y_pred = model.predict(X)
2
3 accuracy_score(y, y_pred)
```

```
1 0.898876404494382
```

Подведем итог

Сегодня мы разобрали множество разновидностей и подходов к использованию логистической регрессии. Давайте систематизируем изученный материал с помощью следующей схемы.



Рассмотрим обучение нейронных сетей.

