

Линейная регрессия и основные библиотеки Python для анализа данных и научных вычислений

Это задание посвящено линейной регрессии. На примере прогнозирования роста человека по его весу Вы увидите, какая математика за этим стоит, а заодно познакомитесь с основными библиотеками Python, необходимыми для дальнейшего прохождения курса.

Материалы

- Лекции данного курса по линейным моделям и градиентному спуску
- [Документация](#) по библиотекам NumPy и SciPy
- [Документация](#) по библиотеке Matplotlib
- [Документация](#) по библиотеке Pandas
- [Pandas Cheat Sheet](#)
- [Документация](#) по библиотеке Seaborn

Задание 1. Первичный анализ данных с Pandas

В этом задании мы будем использовать данные [SOCR](#) по росту и весу 25 тысяч подростков.

[1]. Если у Вас не установлена библиотека Seaborn - выполните в терминале команду `conda install seaborn`. (Seaborn не входит в сборку Anaconda, но эта библиотека предоставляет удобную высокоуровневую функциональность для визуализации данных).

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
```

Считаем данные по росту и весу (`weights_heights.csv`, приложенный в задании) в объект Pandas DataFrame:

```
In [2]: data = pd.read_csv('weights_heights.csv', index_col='Index')
```

Чаще всего первое, что надо сделать после считывания данных - это посмотреть на первые несколько записей. Так можно отловить ошибки чтения данных (например, если вместо 10 столбцов получился один, в названии которого 9 точек с запятой). Также это позволяет познакомиться с данными, как минимум, посмотреть на признаки и их природу (количественный, категориальный и т.д.).

Один из эффективных методов первичного анализа данных - отображение попарных зависимостей признаков. Создается $m \times m$ графиков (m - число признаков), где по диагонали рисуются гистограммы распределения признаков, а вне диагонали - scatter plots зависимости двух признаков. Это можно делать с помощью метода *scatter_matrix* Pandas Data Frame или *pairplot* библиотеки Seaborn.

Чтобы проиллюстрировать этот метод, интересней добавить третий признак. Создадим признак *Индекс массы тела* ([BMI](#)). Для этого воспользуемся удобной связкой метода *apply* Pandas DataFrame и lambda-функций Python.

```
In [6]: def make_bmi(height_inch, weight_pound):
        METER_TO_INCH, KILO_TO_POUND = 39.37, 2.20462
        return (weight_pound / KILO_TO_POUND) / \
            (height_inch / METER_TO_INCH) ** 2
```

```
In [7]: data['BMI'] = data.apply(lambda row: make_bmi(row['Height'],
            row['Weight']), axis=1)
```

[3]. Постройте картинку, на которой будут отображены попарные зависимости признаков , 'Height', 'Weight' и 'BMI' друг от друга. Используйте метод *pairplot* библиотеки Seaborn.

```
In [8]: # Ваш код здесь
```

Часто при первичном анализе данных надо исследовать зависимость какого-то количественного признака от категориального (скажем, зарплаты от пола сотрудника). В этом помогут "ящики с усами" - boxplots библиотеки Seaborn. Box plot - это компактный способ показать статистики вещественного признака (среднее и квантили) по разным значениям категориального признака. Также помогает отслеживать "выбросы" - наблюдения, в которых значение данного вещественного признака сильно отличается от других.

[4]. Создайте в DataFrame *data* новый признак *weight_category*, который будет иметь 3 значения: 1 – если вес меньше 120 фунтов. (~ 54 кг.), 3 - если вес больше или равен 150 фунтов (~68 кг.), 2 – в остальных случаях. Постройте «ящик с усами» (boxplot), демонстрирующий зависимость роста от весовой категории. Используйте метод *boxplot* библиотеки Seaborn и метод *apply* Pandas DataFrame. Подпишите ось *y* меткой «Рост», ось *x* – меткой «Весовая категория».

```
In [9]: def weight_category(weight):
        pass
        # Ваш код здесь

data['weight_cat'] = data['Weight'].apply(weight_category)
# Ваш код здесь
```

[5]. Постройте scatter plot зависимости роста от веса, используя метод *plot* для Pandas DataFrame с аргументом *kind='scatter'*. Подпишите картинку.

In [10]: `# Ваш код здесь`

Задание 2. Минимизация квадратичной ошибки

В простейшей постановке задача прогноза значения вещественного признака по прочим признакам (задача восстановления регрессии) решается минимизацией квадратичной функции ошибки.

[6]. Напишите функцию, которая по двум параметрам w_0 и w_1 вычисляет квадратичную ошибку приближения зависимости роста y от веса x прямой линией $y = w_0 + w_1 * x$:

$$error(w_0, w_1) = \sum_{i=1}^n (y_i - (w_0 + w_1 * x_i))^2$$

Здесь n – число наблюдений в наборе данных, y_i и x_i – рост и вес i -ого человека в наборе данных.

In [11]: `# Ваш код здесь`

Итак, мы решаем задачу: как через облако точек, соответствующих наблюдениям в нашем наборе данных, в пространстве признаков "Рост" и "Вес" провести прямую линию так, чтобы минимизировать функционал из п. 6. Для начала давайте отобразим хоть какие-то прямые и убедимся, что они плохо передают зависимость роста от веса.

[7]. Проведите на графике из п. 5 Задания 1 две прямые, соответствующие значениям параметров $(w_0, w_1) = (60, 0.05)$ и $(w_0, w_1) = (50, 0.16)$. Используйте метод *plot* из *matplotlib.pyplot*, а также метод *linspace* библиотеки NumPy. Подпишите оси и график.

In [12]: `# Ваш код здесь`

Минимизация квадратичной функции ошибки - относительно простая задача, поскольку функция выпуклая. Для такой задачи существует много методов оптимизации. Посмотрим, как функция ошибки зависит от одного параметра (наклон прямой), если второй параметр (свободный член) зафиксировать.

[8]. Постройте график зависимости функции ошибки, посчитанной в п. 6, от параметра w_1 при $w_0 = 50$. Подпишите оси и график.

In [13]: `# Ваш код здесь`

Теперь методом оптимизации найдем "оптимальный" наклон прямой, приближающей зависимость роста от веса, при фиксированном коэффициенте $w_0 = 50$.

[9]. С помощью метода `minimize_scalar` из `scipy.optimize` найдите минимум функции, определенной в п. 6, для значений параметра w_1 в диапазоне $[-5, 5]$. Проведите на графике из п. 5 Задания 1 прямую, соответствующую значениям параметров $(w_0, w_1) = (50, w_{1_opt})$, где w_{1_opt} – найденное в п. 8 оптимальное значение параметра w_1 .

```
In [14]: # Ваш код здесь
```

```
In [15]: # Ваш код здесь
```

При анализе многомерных данных человек часто хочет получить интуитивное представление о природе данных с помощью визуализации. Увы, при числе признаков больше 3 такие картинки нарисовать невозможно. На практике для визуализации данных в 2D и 3D в данных выделяют 2 или, соответственно, 3 главные компоненты (как именно это делается - мы увидим далее в курсе) и отображают данные на плоскости или в объеме.

Посмотрим, как в Python рисовать 3D картинки, на примере отображения функции $z(x, y) = \sin(\sqrt{x^2 + y^2})$ для значений x и y из интервала $[-5, 5]$ с шагом 0.25.

```
In [16]: from mpl_toolkits.mplot3d import Axes3D
```

Создаем объекты типа `matplotlib.figure.Figure` (рисунок) и `matplotlib.axes._subplots.Axes3DSubplot` (ось).

```
In [17]: fig = plt.figure()
ax = fig.gca(projection='3d') # get current axis

# Создаем массивы NumPy с координатами точек по осям X и Y.
# Используем метод meshgrid, при котором по векторам координат
# создается матрица координат. Задаем нужную функцию Z(x, y).
X = np.arange(-5, 5, 0.25)
Y = np.arange(-5, 5, 0.25)
X, Y = np.meshgrid(X, Y)
Z = np.sin(np.sqrt(X**2 + Y**2))

# Наконец, используем метод *plot_surface* объекта
# типа Axes3DSubplot. Также подписываем оси.
surf = ax.plot_surface(X, Y, Z)
ax.set_xlabel('X')
ax.set_ylabel('Y')
ax.set_zlabel('Z')
plt.show()
```

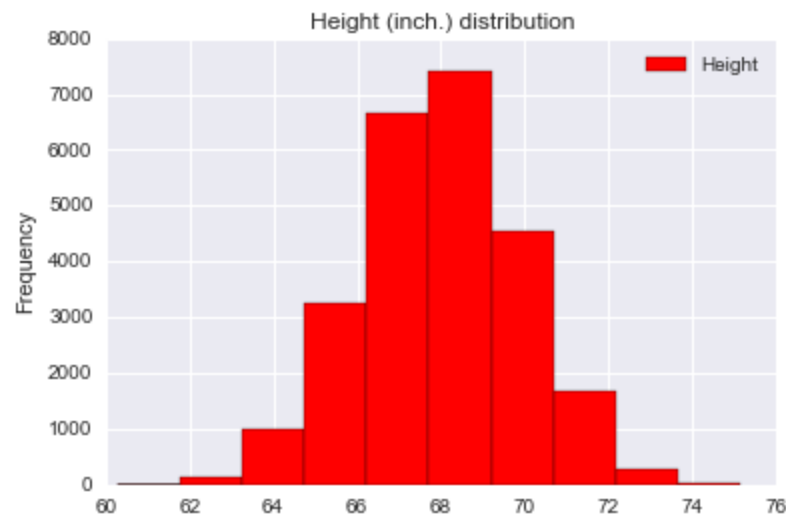
```
//anaconda/lib/python2.7/site-packages/matplotlib/collections.py:590: FutureWarning: elementwise comparison failed; returning scalar instead, but in the future will perform elementwise comparison
    if self._edgecolors == str('face'):
```

После этого стоит построить гистограммы распределения признаков - это опять-таки позволяет понять природу признака (степенное у него распределение, или нормальное, или какое-то еще). Также благодаря гистограмме можно найти какие-то значения, сильно не похожие на другие - "выбросы" в данных. Гистограммы удобно строить методом *plot* Pandas DataFrame с аргументом *kind='hist'*.

Пример. Построим гистограмму распределения роста подростков из выборки *data*. Используем метод *plot* для DataFrame *data* с аргументами *y='Height'* (это тот признак, распределение которого мы строим)

```
In [3]: data.plot(y='Height', kind='hist',  
                 color='red', title='Height (inch.) distribution')
```

```
Out[3]: <matplotlib.axes._subplots.AxesSubplot at 0x10e2d30d0>
```



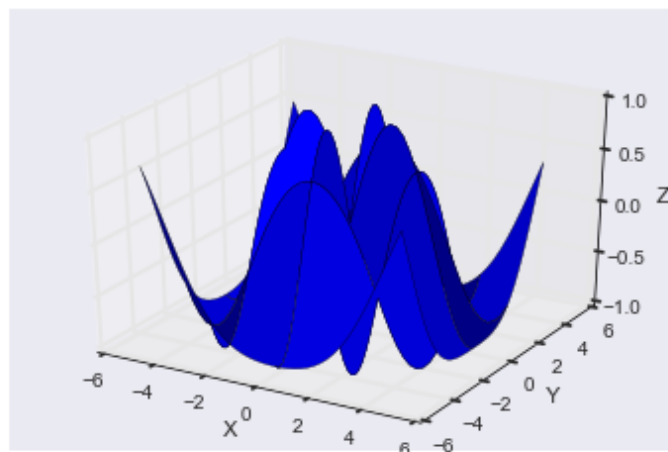
Аргументы:

- *y='Height'* - тот признак, распределение которого мы строим
- *kind='hist'* - означает, что строится гистограмма
- *color='red'* - цвет

[2]. Посмотрите на первые 5 записей с помощью метода *head* Pandas DataFrame. Нарисуйте гистограмму распределения веса с помощью метода *plot* Pandas DataFrame. Сделайте гистограмму зеленой, подпишите картинку.

```
In [4]: # Ваш код здесь
```

```
In [5]: # Ваш код здесь
```



[10]. Постройте 3D-график зависимости функции ошибки, посчитанной в п.6 от параметров w_0 и w_1 . Подпишите ось x меткой «Intercept», ось y – меткой «Slope», а ось z – меткой «Error».

In [18]: # Ваш код здесь

In [19]: # Ваш код здесь

In [20]: # Ваш код здесь

[11]. С помощью метода *minimize* из `scipy.optimize` найдите минимум функции, определенной в п. 6, для значений параметра w_0 в диапазоне $[-100, 100]$ и w_1 - в диапазоне $[-5, 5]$. Начальная точка – $(w_0, w_1) = (0, 0)$. Используйте метод оптимизации L-BFGS-B (аргумент `method` метода *minimize*). Проведите на графике из п. 5 Задания 1 прямую, соответствующую найденным оптимальным значениям параметров w_0 и w_1 . Подпишите оси и график.

In [21]: # Ваш код здесь

In [22]: # Ваш код здесь

Критерии оценки работы

- Выполняется ли тетрадка IPython без ошибок? (15 баллов)
- Верно ли отображена гистограмма распределения роста из п. 2? (3 балла). Правильно ли оформлены подписи? (1 балл)
- Верно ли отображены попарные зависимости признаков из п. 3? (3 балла). Правильно ли оформлены подписи? (1 балл)
- Верно ли отображена зависимость роста от весовой категории из п. 4? (3 балла). Правильно ли оформлены подписи? (1 балл)
- Верно ли отображен scatter plot роста от веса из п. 5? (3 балла). Правильно ли оформлены подписи? (1 балл)
- Правильно ли реализована функция подсчета квадратичной ошибки из п. 6? (10 баллов)
- Правильно ли нарисован график из п. 7? (3 балла) Правильно ли оформлены подписи? (1 балл)
- Правильно ли нарисован график из п. 8? (3 балла) Правильно ли оформлены подписи? (1 балл)
- Правильно ли используется метод `minimize_scalar` из `scipy.optimize`? (6 баллов). Правильно ли нарисован график из п. 9? (3 балла)
Правильно ли оформлены подписи? (1 балл)
- Правильно ли нарисован 3D-график из п. 10? (6 баллов) Правильно ли оформлены подписи? (1 балл)
- Правильно ли используется метод `minimize` из `scipy.optimize`? (6 баллов). Правильно ли нарисован график из п. 11? (3 балла). Правильно ли оформлены подписи? (1 балл)