

Библиотечные инструменты языка программирования Python

1. Поиск характеристик и визуализация теоретических законов распределений

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import scipy.special as sc import matplotlib.pyplot as plt %matplotlib inline</pre>
Вероятностные распределения в модуле scipy.stats	
	<p>Пусть R - обозначение вида закона распределения, $params$ - параметры распределения.</p> <p>Общий вид обращения к распределению: $sts.R(params)$ - закон распределения R с параметрами $params$.</p> <p>Некоторые виды распределений: $sts.uniform(a1, a2)$ - равномерный закон распределения; если $X \sim R(a, b)$, то $a1$ совпадает с a, $a2$ совпадает с $b - a$.</p> <p>$sts.norm(a, b)$ - нормальный закон распределения с параметрами a (математическое ожидание) и b (среднее квадратичное отклонение).</p> <p>$sts.expon(0, m)$ - экспоненциальный закон распределения; m совпадает с математическим ожиданием, т.е. если $X \sim Exp(\lambda)$, то m совпадает с $1/\lambda$.</p> <p>$sts.chi2(n)$ - закон распределения хи-квадрат с n степенями свободы.</p> <p>$sts.t(n)$ - закон распределения Стьюдента с n степенями свободы.</p> <p>$sts.f(k1, k2)$ - закон распределения Фишера со степенями свободы с $k1$ и $k2$.</p> <p>$sts.binom(n, p)$ - биномиальный закон распределения с параметрами n (общее число испытаний) и p (вероятность успеха в одном испытании).</p> <p>$sts.poisson(lm)$ - закон распределения Пуассона с параметром lm, параметр совпадает с традиционным параметром распределения λ (равен математическому ожиданию).</p>
Методы поиска функциональных характеристик вероятностных распределений в модуле scipy.stats	
	<p>Пусть R - обозначение типа закона распределения, $params$ - параметры распределения. Тогда:</p> <p>$R(params).cdf(x)$ - значение функции закона распределения R с параметрами $params$ в точке x.</p> <p>$R(params).pdf(x)$ - для непрерывной случайной величины значение плотности распределения в точке x.</p> <p>$R(params).pmf(x)$ - для дискретной случайной величины вероятность принять значение x.</p>
Методы поиска числовых характеристик случайной величины в модуле scipy.stats	
	<p>Пусть R - обозначение типа закона распределения, $params$ - параметры распределения. Тогда:</p> <p>$R(params).ppf(q)$ - квантиль порядка q.</p> <p>$R(params).mean()$ - математическое ожидание.</p> <p>$R(params).var(x)$ - дисперсия.</p> <p>$R(params).std(x)$ - среднее квадратичное (стандартное) отклонение.</p> <p>$R(params).median()$ - медиана.</p>

	<p><code>R(params).moment(k)</code> - начальный момент порядка k.</p> <p><code>R(params).stats('mvsk')</code> - математическое ожидание, дисперсия, коэффициент асимметрии и коэффициент эксцесса.</p>
Некоторые математические функции в модуле <code>scipy.special</code>	
	<p><code>sc.factorial(n)</code> - $n!$</p> <p><code>sc.comb(n, k)</code>, <code>sc.binom(n, k)</code> - C_n^k, число сочетаний из n по k.</p> <p><code>sc.perm(n, k)</code> - A_n^k, число размещений из n по k.</p> <p><code>sc.gamma(x)</code> - $\Gamma(x) = \int_0^{+\infty} t^{x-1} e^{-t} dt$, гамма-функция.</p>

2. Компьютерное моделирование выборок дискретных случайных величин, первичная обработка выборки

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import scipy.special as sc import matplotlib.pyplot as plt import random %matplotlib inline</pre>
Генераторы выборок дискретных распределений в библиотеке <code>numpy</code>	
	<p>Функция <code>np.random.choice(a, size, replace=True, p)</code> Возвращает выборку заданного объема <code>size</code> (если <code>size</code> - число) из массива <code>a</code>. По умолчанию повторную, если <code>replace=False</code>, то неповторную. Выбор элемента из массива <code>a</code> осуществляется с соответствующей вероятностью массива <code>p</code>. Если <code>size</code> - кортеж, то генерируется массив заданной формы. Другие варианты: если <code>a</code> - целое число, то генерируется случайное число из массива <code>arange(a)</code>; если параметр <code>p</code> не задан, то элементы из <code>a</code> равновозможны.</p>
	<p>Функция <code>np.random.binomial(n, p, size)</code> Возвращает выборку заданного объема <code>size</code> (если <code>size</code> - число) биномиального распределения с параметрами <code>n</code> и вероятностью <code>p</code>. Если <code>size</code> - кортеж, то генерируется массив заданной формы.</p>
	<p>Функция <code>np.random.poisson(lm, size)</code> Формирует выборку заданного объема <code>size</code> (если <code>size</code> - число) распределения Пуассона с параметром <code>lm</code>. Если <code>size</code> - кортеж, то генерируется массив заданной формы.</p>
	<p>В модуле <code>numpy.random</code> также имеются генераторы выборок следующих распределений: геометрического, гипергеометрического, отрицательного биномиального, распределения <code>logser</code> и др.</p>
Средства формирования вариационного и статистического рядов в библиотеке <code>numpy</code>	
	<p>Функция <code>np.bincount(x)</code> Вычисляет количество появлений в массиве каждого целого числа от 0 до <code>max(x)</code>.</p>
	<p>Функция <code>np.sort(list(set(x)))</code> Формирует набор уникальных элементов (вариант) выборки.</p>
Случайное перемешивание	
	<p>Функция <code>np.random.shuffle(x)</code> случайным образом перемешивает массив <code>x</code>.</p>

3. Компьютерное моделирование выборок непрерывных случайных величин, первичная обработка выборки

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts from statsmodels.distributions.empirical_distribution import ECDF import matplotlib.pyplot as plt %matplotlib inline</pre>
Генераторы выборок некоторых непрерывных распределений в библиотеке numpy	
	<p>Функция <code>np.random.normal(loc, scale, size)</code> Возвращает выборку заданного объема <code>size</code> (если <code>size</code> - число) нормального распределения $N(m, \sigma)$. Параметры: <code>loc=m</code>; <code>scale=σ</code>. Если <code>size</code> - кортеж, то генерируется массив заданной формы.</p>
	<p>Функция <code>np.random.uniform(low=0.0, high=1.0, size=None)</code> Возвращает выборку заданного объема <code>size</code> (если <code>size</code> - число) равномерного распределения $R(a, b)$. Параметры: <code>low=a</code>; <code>high=b</code>. Если <code>size</code> - кортеж, то генерируется массив заданной формы.</p>
	<p>Функция <code>np.random.exponential(scale=1.0, size=None)</code> Возвращает выборку заданного объема <code>size</code> (если <code>size</code> - число) экспоненциального распределения $Ex(\lambda)$. Параметр: <code>scale = 1/λ</code> (равен математическому ожиданию). Если <code>size</code> - кортеж, то генерируется массив заданной формы.</p>
	<p>Функция <code>np.random.chisquare(df, size=None)</code> Возвращает выборку заданного объема <code>size</code> (если <code>size</code> - число) распределения хи-квадрат с <code>df</code> степенями свободы. Если <code>size</code> - кортеж, то генерируется массив заданной формы.</p>
	<p>Функция <code>np.random.f(dfnum, dfden, size)</code> Возвращает выборку заданного объема <code>size</code> (если <code>size</code> - число) распределения Фишера со степенями свободы <code>dfnum</code>, <code>dfden</code>. Если <code>size</code> - кортеж, то генерируется массив заданной формы.</p>
	<p>В модуле <code>numpy.random</code> также имеются генераторы выборок следующих распределений: бета, гамма, Гумбеля, Лапласа, логистического, логнормального, степенного, Рэлея, треугольного, Ломакса (Парето II вида), фон Мизеса, Уайльда, Вейбулла и др.</p>
Построение эмпирической функции распределения в модуле <code>statsmodels.distributions.empirical_distribution</code>	
	<p>Функция <code>statsmodels.distributions.empirical_distribution.ECDF(x, side=right)</code> Возвращает эмпирическую функцию распределения. Параметры: <code>x</code> - массив (выборка); <code>side</code> - задает форму интервалов, по которым строятся ступени эмпирической функции: <code>right</code> (по умолчанию) - интервалы вида [...), открытые справа, <code>left</code> - интервалы вида (...], открытые слева.</p>
Средства визуализации: построение гистограммы и эмпирической функции распределения в модуле <code>matplotlib.pyplot</code>	

	<p>Функция <code>plt.hist(x, bins=None, density=None, weights=None, cumulative=False, histtype='bar', align='mid', orientation='vertical', log=False, color=None)</code></p> <p>Строит гистограмму и возвращает два массива: высот столбцов гистограммы и центров интервалов группировки.</p> <p>Параметры: <code>x</code> - массив (выборка); <code>bins</code> - число интервалов группировки или последовательность, задающая границы интервалов (все интервалы, кроме последнего, полуоткрытые вида [...)), или строка из списка, который приводится после перечня параметров; <code>density</code> - если <code>True</code>, то строится гистограмма относительных частот (суммарная площадь прямоугольников равна 1); <code>weights</code> - массив весов той же формы, что и <code>x</code>; <code>cumulative</code> - если <code>True</code>, то в сочетании с признаком <code>density=True</code> строит эмпирическую функцию распределения; <code>histtype</code> - кроме типа <code>'bar'</code> можно указать <code>'barstacked'</code> и <code>'step'</code>; <code>align</code> - задает расположение центров прямоугольников; <code>orientation</code> - установив значение <code>'horizontal'</code>, можно повернуть график на 90°; <code>log</code> - если <code>True</code>, для осей используется логарифмическая шкала; <code>color</code> - признак, устанавливающий цвет.</p> <p>По умолчанию число интервалов группировки равно 10.</p> <p>Список правил для выбора числа интервалов:</p> <p><code>bins='auto'</code> - максимальное из значений, получаемых по правилу Стерджесса и Фридмана - Диакониса;</p> <p><code>bins='fd'</code> - правило Фридмана - Диакониса;</p> <p><code>bins='sturges'</code> - правило Стерджесса;</p> <p><code>bins='doane'</code> - правило Дозна;</p> <p><code>bins='scott'</code> - правило Скотта;</p> <p><code>bins='stone'</code> - обобщение правила Скотта;</p> <p><code>bins='rice'</code> - правило Райса;</p> <p><code>bins='sqrt'</code> - правило квадратного корня.</p>
--	---

4. Точечное оценивание параметров распределения по выборке

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import matplotlib.pyplot as plt import seaborn %matplotlib inline</pre>
Точечные оценки параметров распределения в пакете numpy	
	<p>Функция <code>np.mean(a, axis)</code></p> <p>Возвращает выборочное среднее.</p> <p>Параметры: <code>a</code> - массив; в случае многомерного массива <code>a</code> можно указать ось (<code>axis</code>), вдоль которой вычисляется среднее.</p> <p>Функция <code>np.nanmean(a, axis)</code> при вычислении игнорирует пропущенные данные <code>nan</code> (важно при обработке реальных данных).</p> <p>Для вычисления выборочных начальных моментов порядка k можно использовать функцию <code>mean</code> применительно к k-й степени массива <code>a</code>.</p>
	<p>Функция <code>np.var(a, axis, ddof)</code></p> <p>Возвращает оценку дисперсии по выборке <code>a</code>.</p>

	<p>Параметры: a - массив; в случае многомерного массива a можно указать ось ($axis$), вдоль которой вычисляется дисперсия; $ddof$ по умолчанию равен 0 (вычисляется выборочная дисперсия), если задать $ddof=1$, то функция возвращает исправленную выборочную дисперсию.</p> <p>Функция <code>np.nanvar(a, axis, ddof)</code> при вычислении игнорирует пропущенные данные <code>nan</code>.</p>
	<p>Функция <code>np.std(a, axis, ddof)</code></p> <p>Возвращает корень из выборочной (или исправленной выборочной) дисперсии.</p> <p>Параметры: a - массив; в случае многомерного массива a можно указать ось, вдоль которой вычисляется дисперсия; $ddof$ по умолчанию равен 0 (вычисляется выборочная дисперсия), если задать $ddof=1$, то функция возвращает исправленную выборочную дисперсию.</p>
	<p>Функция <code>np.median(a, axis=None, out=None)</code></p> <p>Возвращает выборочную медиану (вычисляется как центральный элемент $a_{\frac{n-1}{2}}$ отсортированного по неубыванию массива a при нечетном n и как среднее арифметическое двух центральных значений при четном n).</p> <p>Параметры: a - массив; в случае многомерного массива a можно указать ось ($axis$), вдоль которой вычисляется среднее; out - массив, если он указан, в него помещаются вычисленные значения медиан.</p>
	<p>Функция <code>np.quantile(a, q, axis=None, out=None, interpolation='linear')</code></p> <p>Возвращает квантиль порядка q (указывается число из интервала $(0,1)$).</p> <p>Параметры: a – массив; в случае многомерного массива a параметр $axis$ - ось (кортеж осей), вдоль которой производятся вычисления; out - массив, если он указан, в него помещаются вычисленные значения квантилей; $interpolation$ - признак, определяющий метод интерполяции в ситуации, когда квантиль расположена между двумя значениями массива ('linear' по умолчанию, есть другие варианты).</p> <p>Первый квартиль Q_1 выборки X вычисляется с помощью функции <code>np.quantile(X, 0.25)</code>.</p> <p>Третий квартиль Q_3 выборки X вычисляется с помощью функции <code>np.quantile(X, 0.75)</code>.</p>
Точечные оценки параметров распределения в модуле <code>scipy.stats</code>	
	<p>Функция <code>sts.moment(x, moment=k, axis=0, nan_policy='propagate')</code></p> <p>Возвращает выборочный центральный момент порядка k.</p> <p>Параметры: x – выборка; $axis$ - ось, вдоль которой вычисляется оценка; nan_policy - определяет способ обработки пропущенных значений ('propagate' - возвращает <code>nan</code>, 'raise' - генерирует ошибку, 'omit' - игнорирует пропущенные данные).</p>
	<p>Функция <code>sts.skew(x, axis=0, bias=True, nan_policy='propagate')</code></p> <p>Возвращает выборочный коэффициент асимметрии.</p> <p>Параметры: x – выборка; $axis$ - ось, вдоль которой вычисляется оценка; $bias$ - признак (если <code>False</code> - применяется коррекция для устранения смещенности); nan_policy - определяет способ обработки пропущенных значений.</p>
	<p>Функция <code>sts.kurtosis(x, axis=0, fisher=True, bias=True, nan_policy='propagate')</code></p> <p>Возвращает выборочный коэффициент асимметрии.</p>

	<p>Параметры: <code>x</code> – выборка; <code>axis</code> - ось, вдоль которой вычисляется оценка; <code>fisher</code> - признак, если равен <code>True</code> (по умолчанию), то в формуле для эксцесса из отношения моментов вычисляется число 3; <code>bias</code> - признак (если <code>False</code> - применяется коррекция для устранения смещенности); <code>nan_policy</code> - определяет способ обработки пропущенных значений.</p>
	<p>Функция <code>sts.igr(x)</code> Вычисляет межквартильный размах - разность между третьим и первым квартилями.</p>
	<p>Функция <code>sts.describe(a, axis, ddof, bians, nan_policy)</code> Возвращает набор оценок основных параметров случайной величины: <code>nobs</code> - объем выборки; <code>minmax</code> - кортеж, содержащий максимальное и минимальное значение выборки; <code>mean</code> - выборочное среднее; <code>variance</code> - исправленная выборочная дисперсия s^2 (в случае задания <code>ddof=1</code> или по умолчанию) либо выборочная дисперсия (в случае задания <code>ddof=0</code>); <code>skewness</code> - коэффициент асимметрии; <code>kurtosis</code> - коэффициент эксцесса (в случае задания <code>bias=False</code>, коэффициенты асимметрии и эксцесса корректируются на величину смещения). Параметры: <code>a</code> - выборка; <code>axis</code> - задание оси (для многомерной выборки); <code>ddof</code> - признак смещенности (только для дисперсии); <code>bians</code> - признак коррекции (только для асимметрии и эксцесса); <code>nan_policy</code> - задает способ обработки пропущенных данных.</p>
Средства визуализации: построение гистограммы и боксплота в пакете <code>seaborn</code>	
	<p>Функция <code>boxplot(x=None, y=None, hue=None, data=None, order=None, hue_order=None, orient=None, color=None)</code> Строит боксплот. Параметры: <code>x, y, hue</code> - наименование признаков в наборе <code>data</code>; <code>data</code> - датафрейм, или массив <code>numpy</code>, или список; <code>order</code> и <code>hue_order</code> - строки, с помощью которых можно изменить порядок вывода признаков на график; <code>orient</code> - вертикальная или горизонтальная ориентация («v» или «u»); <code>color</code> - задание цвета. Построение боксплота в пакете удобно сочетать с построением гистограммы с помощью функции <code>histplot</code>.</p>

5. Интервальное оценивание параметров распределения по выборке

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import scipy.special as sc import matplotlib.pyplot as plt import statsmodels.api as sm import statsmodels.stats.weightstats import statsmodels.stats.proportion %matplotlib inline</pre>
Построение доверительных интервалов в предположении нормального распределения генеральной совокупности	
	<p>Функция <code>_zconfint_generic (mean, std_mean, alpha, alternative)</code> модуля <code>statsmodels.stats.weightstats</code> Возвращает границы доверительного интервала (1) (см. приложение 1).</p>

	<p>Параметры: mean - выборочное среднее, $\text{std_mean} = \frac{\sigma}{\sqrt{n}}$; α - уровень значимости; alternative - вид доверительного интервала («two-sided» - двусторонний, по умолчанию, «smaller» - левосторонний, «larger» - правосторонний).</p>
	<p>Функция <code>zconfint(x, alpha=0.05, alternative=«two-sided»)</code> модуля <code>statsmodels.api.stats</code></p> <p>Возвращает границы доверительного интервала (2).</p> <p>Параметры: x - выборка; α - уровень значимости; alternative - вид доверительного интервала («two-sided» - двусторонний (по умолчанию), «smaller» - левосторонний, «larger» - правосторонний).</p>
Построение доверительных интервалов в предположении биномиального распределения генеральной совокупности	
	<p>Функция <code>proportion_confint(count=m, nobs=n, alpha, method='normal')</code> модуля <code>statsmodels.stats.proportion</code></p> <p>Возвращает границы доверительного интервала (3) или (4).</p> <p>Параметры: count - число успехов; nobs - число испытаний; α - уровень значимости; method - вид доверительного интервала (если 'normal', строится интервал (6), если 'wilson', строится интервал (3)).</p>
	<p>Функция <code>samplesize_confint_proportion(proportion, half_length, alpha, method='normal')</code> модуля <code>statsmodels.stats.proportion</code></p> <p>Возвращает минимальный объем выборки, необходимый для достижения желаемой точности при интервальном оценивании вероятности события.</p> <p>Параметры: proportion - вероятность успеха; half_length - половина ширины требуемого интервала; α - уровень значимости (по умолчанию 0,05).</p>

6. Проверка гипотез о значениях параметров распределения

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import matplotlib.pyplot as plt %matplotlib inline</pre>
Проверка гипотезы о значении математического ожидания нормально распределенной генеральной совокупности при неизвестной дисперсии при двусторонней альтернативе	
	<p>Функция <code>ttest_1samp(a, popmean, axis = 0, nan_policy = 'propagate')</code> модуля <code>scipy.stats</code></p> <p>Возвращает: выборочное значение статистики $W = \frac{\bar{X} - m_0}{S / \sqrt{n}}$; достигаемый уровень значимости - p - значение (см. приложение 2).</p> <p>Параметры: a – выборка; popmean - гипотетическое значение математического ожидания; nan_policy - задает способ обработки пропущенных значений.</p>

7. Проверка гипотез о законе распределения, проблема нормализации выборки

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import matplotlib.pyplot as plt import pandas as pd %matplotlib inline</pre>
Экспорт данных из csv или Excell- файла в объект DataFrame библиотеки Pandas	
	<p>Функция <code>pd.read_csv('Data.xlsx', sep=',', header = 'infer', index_col=None)</code> Создает объект <code>DataFrame</code>.</p> <p>Параметры: <code>Data.xlsx</code> - строка с указанием пути к файлу; <code>sep</code> - разделитель (по умолчанию <code>,</code>); <code>header</code> - строка, содержащая имена столбцов (по умолчанию <code>'infer'</code> - в качестве имен используется первая строка данных); <code>index_col</code> - указывает, какой столбец в файле использовать в качестве индекса, если установить <code>index_col= False</code>, первый столбец данных в качестве индекса использоваться не будет.</p> <p>Функция <code>pd.read_excel()</code> имеет аналогичные параметры.</p>
Средства группировки выборки библиотеки numpy	
	<p>Функция <code>np.histogram(a, bins=10, range=None, weights=None, density= None)</code> Возвращает два массива: <code>hist</code> - массив высот столбцов гистограммы; <code>bin_edges</code> - массив границ интервалов.</p> <p>Параметры: <code>a</code> - одномерный массив (выборка); <code>bins</code> - число интервалов группировки (по умолчанию - 10) или последовательность, задающая границы интервалов; если <code>bins='auto'</code> - число интервалов выбирается как максимальное из величин, получаемых по правилу Стерджесса и Фридмана-Диакониса; <code>range</code> - начальная и конечная границы интервалов (если параметр не определен, то в качестве границ берутся минимальный и максимальный элементы выборки), элементы выборки вне области <code>range</code> игнорируются; <code>density</code> - если <code>True</code> - строится гистограмма относительных частот (суммарная площадь прямоугольников равна 1); <code>weights</code> - массив весов той же формы, что и <code>a</code>.</p>
Реализация критерия хи-квадрат проверки гипотезы о законе распределения в модуле scipy.stats	
	<p>Функция <code>sts.chisquare (f_obs, f_exp = None, ddof = 0, axis = 0)</code> Возвращает наблюдаемое значение статистики и p - значение (т.е. максимальное значение уровня значимости, при котором основная гипотеза принимается).</p> <p>Параметры: <code>f_obs</code> - наблюдаемые частоты (n_i); <code>f_exp</code> - частоты гипотетического (согласно основной гипотезе) распределения (np_i) (по умолчанию равные между собой); <code>ddof</code> - число параметров гипотетического распределения, оцениваемых по выборке. На вход функции можно подавать многомерный массив. Критерий будет применяться к каждому столбцу массива (если <code>axis = 1</code>, то к строке).</p> <p>Условие использования: все наблюдаемые и гипотетические частоты должны быть не менее 5 ($n_i \geq 5, np_i \geq 5$).</p>
Реализация критерия Шапиро-Уилка в модуле scipy.stats	
	<p>Функция <code>sts.shapiro (x, a=None, reta=False)</code> Возвращает наблюдаемое значение статистики p - значение; массив параметров (присутствует, если <code>reta=True</code>).</p>

	<p>Параметры: x - одномерный массив (выборка); a - массив внутренних параметров (если не заданы, вычисляется самой функцией); $reta$ - признак, нужно ли возвращать вычисленные параметры).</p>
Преобразование Бокса-Кокса модуля <code>scipy.stats</code>	
	<p>Функция <code>sts.boxcox(x, lambda=None, alpha=None)</code> Возвращает: (1) <code>boxcox</code> - массив, результат преобразования Бокса-Кокса; (2) если параметр <code>lambda=None</code>, то второй возвращаемый параметр <code>maxlog</code> - значение <code>lambda</code>, максимизирующее логарифм функции правдоподобия; (3) если <code>lambda=None</code> и <code>alpha</code> не <code>None</code>, возвращается кортеж, содержащий границы доверительного интервала. Параметры: x - входной одномерный массив положительных чисел (выборка); <code>lmbd</code> - если не <code>None</code>, преобразование выполняется для этого значения; <code>alpha</code> - если не <code>None</code>, то функция возвращает в качестве третьего аргумента $100(1-\alpha)\%$-й доверительный интервал для параметра <code>lambda</code>.</p>

8. Корреляционный анализ

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import matplotlib.pyplot as plt import pandas as pd import seaborn %matplotlib inline</pre>
Генерация многомерного нормального распределения в библиотеке <code>numpy</code>	
	<p>Функция <code>np.random.multivariate_normal(mean, cov, n)</code> Возвращает выборку объема n для многомерного нормального распределения с заданным вектором математических ожиданий <code>mean</code> и ковариационной матрицей <code>cov</code>.</p>
Средства визуализации диаграммы рассеивания	
	<p>Функция <code>plt.scatter(x, y)</code> модуля <code>matplotlib.pyplot</code> Строит диаграмму рассеивания признаков x и y. Параметры: x, y - два массива одинаковой длины.</p>
	<p>Функция <code>seaborn.pairplot(data, vars=None, kind='scatter', diag_kind='hist', height=4)</code> пакета <code>seaborn</code> Строит диаграммы рассеивания пар признаков из <code>vars</code>, а также визуализирует распределение отдельных признаков. Параметры: - датафрейм; <code>vars</code> - список имен переменных из <code>vars</code>, которые будут использованы для вывода диаграммы (если не задан, используются все числовые колонки <code>data</code>); <code>kind</code> - тип диаграммы рассеяния (обычная <code>'scatter'</code> или с линией регрессии <code>'reg'</code>); <code>diag_kind</code> - тип диагональных графиков (<code>'auto'</code>, <code>'hist'</code>, <code>'kde'</code>); <code>height</code> - высота каждой факеты (в дюймах).</p>
Расчет выборочных характеристик двумерной выборки	
	<p>Функция <code>np.cov(x, y=None, rowvar=True, bias=False, ddof=None)</code> библиотеки <code>numpy</code> Вычисляет выборочную ковариационную матрицу. Параметры: x - одномерный или двумерный массив. Если одномерный - вычисляется ковариация между x и y. Если двумерный - при значении <code>rowvar=True</code> (по</p>

	<p>умолчанию) вычисляется ковариация между строками массива x, при значении <code>rowvar=False</code> - между столбцами массива x; <code>bias</code> - признак, определяющий способ нормализации, по умолчанию (<code>False</code>) - производится деление на $n-1$, иначе на n; <code>ddof</code> - выполняет функцию, аналогичную признаку <code>bias</code>: при значении <code>ddof=1</code> производится деление на $n-1$, при значении <code>ddof=0</code> - деление на n.</p>
	<p>Функция <code>np.corrcoef(x, y=None, rowvar=True)</code> библиотеки <code>numpy</code> Вычисляет выборочную корреляционную матрицу. Параметры: x - одномерный или двумерный массив. Если одномерный - вычисляется коэффициент корреляции между x и y. Если двумерный: при значении <code>rowvar=True</code> (по умолчанию) вычисляется коэффициент корреляции между строками массива x, при значении <code>rowvar=False</code> - между столбцами массива x.</p>
	<p>Метод <code>data.corr(method='pearson')</code> библиотеки <code>pandas</code> Параметры: <code>data</code> - объект <code>DataFrame</code>, <code>method</code> - задает вид коэффициента корреляции <code>'pearson'</code>, <code>'spearman'</code>, <code>'kendall'</code> (по умолчанию <code>'pearson'</code>).</p>
Средства визуализация корреляционной матрицы	
	<p>Функция <code>seaborn.heatmap(data, annot=None, fmt='.2g', linewidth=0, linecolor='white', cbar=True, cbar_kws=None, cbar_ax=None)</code> пакета <code>seaborn</code> Принимает на вход прямоугольный массив данных и отображает данные с помощью цвета. Цветовая панель показывает соответствие цвета числовым значениям переменной. Параметры: <code>data</code> - объект <code>DataFrame</code>, <code>annot</code> - признак: если <code>True</code>, в каждую ячейку карты выводится значение признака; <code>fmt</code> - строка: задает формат для случая <code>annot=True</code>, <code>linewidth</code>, <code>linecolor</code> - размер и цвет линий, разделяющих ячейки; <code>cbar</code>, <code>cbar_kws</code>, <code>cbar_ax</code> - информация о необходимости вывода, цвете и расположении цветовой панели.</p>
Критерии значимости коэффициента корреляции Спирмена модуля <code>scipy.stats</code>	
	<p>Функция <code>sts.spearmanr(a,b=None, axis=0, nan_policy='propagate')</code> Проверяет гипотезу об отсутствии значимой монотонной связи. Возвращает r - выборочный коэффициент корреляции, p-value - достигаемый уровень значимости. Параметры: a, b - два одномерных или двумерных массива одинакового размера, <code>nan_policy</code> (<code>'propagate'</code>, <code>'raise'</code>, <code>'omit'</code>) - задает способ обработки пропущенных (NaN) значений.</p>

9. Регрессионный анализ: парная линейная регрессия

Загрузка основных модулей	
	<pre>import numpy as np import scipy.stats as sts import matplotlib.pyplot as plt import pandas as pd import seaborn from sklearn.linear_model import LinearRegression %matplotlib inline</pre>
Средства получения оценок линейной регрессии	
	<p>Вначале надо создать экземпляр класса <code>LinearRegression</code>, который будет представлять модель регрессии: <code>linreg = LinearRegression()</code>.</p> <p>Функция <code>linreg.fit(x, y)</code> Вычисляет оценки коэффициентов регрессии b_0, b_1. Ее параметры: x - двумерный массив, размера $n \times 1$ (независимая переменная); y - одномерный массив (вектор-строка) длины n (зависимая переменная).</p> <p>Коэффициент b_0 можно получить, вызвав <code>linreg.intercept_</code>.</p> <p>Коэффициент b_1 можно получить, вызвав <code>linreg.coef_</code>.</p> <p>Функция <code>linreg.score(x, y)</code> вычисляет коэффициент детерминации.</p>

Приложение 1

Параметр	Условия	Доверительный интервал	Номер формулы
m	$X \sim N(m, \sigma)$ σ известно	$\bar{X} - u_{\frac{1+\beta}{2}} \cdot \frac{\sigma}{\sqrt{n}} < m < \bar{X} + u_{\frac{1+\beta}{2}} \cdot \frac{\sigma}{\sqrt{n}}$	1
m	$X \sim N(m, \sigma)$ σ неизвестно	$\bar{X} - t_{\frac{1+\beta}{2}}(n-1) \cdot \frac{S}{\sqrt{n}} < m < \bar{X} + t_{\frac{1+\beta}{2}}(n-1) \cdot \frac{S}{\sqrt{n}}$	2
p	$X \sim B(n, p)$ $\sqrt{npq} \gg 1$	$\frac{P^* + \frac{1}{2} \frac{u_{\frac{1+\beta}{2}}^2}{n} - u_{\frac{1+\beta}{2}} \cdot \sqrt{\frac{P^*(1-P^*)}{n} + \frac{1}{4} \frac{u_{\frac{1+\beta}{2}}^2}{n^2}}}{1 + \frac{u_{\frac{1+\beta}{2}}^2}{n}} < p < \frac{P^* + \frac{1}{2} \frac{u_{\frac{1+\beta}{2}}^2}{n} + u_{\frac{1+\beta}{2}} \cdot \sqrt{\frac{P^*(1-P^*)}{n} + \frac{1}{4} \frac{u_{\frac{1+\beta}{2}}^2}{n^2}}}{1 + \frac{u_{\frac{1+\beta}{2}}^2}{n}}$	3
p	$X \sim B(n, p)$ $\sqrt{npq} \gg 1$	$P^* - u_{\frac{1+\beta}{2}} \cdot \sqrt{\frac{P^*(1-P^*)}{n}} < p < P^* + u_{\frac{1+\beta}{2}} \cdot \sqrt{\frac{P^*(1-P^*)}{n}}$	4

Модификация схемы проверки статистических гипотез с использованием p -значения.

Метод, основанный на использовании так называемого p -значения критерия, позволяет решить для всех уровней значимости одновременно, принять или отклонить основную гипотезу.

Определение. p -значением $p(x_1, x_2, \dots, x_n)$ нулевой гипотезы, проверяемой по выборке с помощью статистики критерия Z и критической области G_α , называется наименьший уровень значимости, при котором основная гипотеза при имеющейся выборке отклоняется:

$$p(x_1, x_2, \dots, x_n) = \min \{ \alpha \mid z_{\text{выб}} \in G_\alpha \}.$$

Здесь $z_{\text{выб}} = Z(x_1, x_2, \dots, x_n)$ – выборочное значение статистики.

Для всех значений уровня значимости, таких, что $\alpha \leq p(x_1, x_2, \dots, x_n)$, основная гипотеза принимается, при всех $\alpha > p(x_1, x_2, \dots, x_n)$ – отклоняется. Чем меньше p -значение, тем больше оснований отклонить нулевую гипотезу.

Вид формул, по которым вычисляются p -значения, зависит от вида критической области.

Справедливо следующее утверждение:

1. Если критическая область правосторонняя, т.е. имеет вид $(z_{1-\alpha}; +\infty)$, где $z_{1-\alpha}$ – квантиль порядка $1-\alpha$, то p -значение находится по формуле $p(x_1, x_2, \dots, x_n) = P\{Z > z_{\text{выб}} | H_0\}$.
2. Если критическая область левосторонняя, т.е. имеет вид $(-\infty; z_\alpha)$, где z_α – квантиль порядка α , то p -значение находится по формуле $p(x_1, x_2, \dots, x_n) = P\{Z < z_{\text{выб}} | H_0\}$.
3. Если критическая область двусторонняя, т.е. имеет вид $(-\infty; z_{\alpha/2}) \cup (z_{1-\alpha/2}; +\infty)$, то p -значение находится по формуле $p(x_1, x_2, \dots, x_n) = 2 \cdot \min\{p, 1-p\}$, где $p = P\{Z < z_{\text{выб}} | H_0\}$.