

Министерство науки и высшего образования Российской Федерации

ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ
ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ
ИТМО

Лабораторная работа №2
по дисциплине
"Линейная алгебра и анализ данных"

Семестр II

Выполнили:
Бобров Михаил гр. J3113
ИСУ: 465234
Чунихина Валерия гр. J3114
ИСУ: 468026

Отчет сдан:
24.04.2025
05:55

Санкт-Петербург, 2025

Содержание

1	Введение	3
2	Теория	5
2.1	Доказать	5
2.2	Доказательство	5
3	Практика	6
3.1	Метод Гаусса	6
3.2	Центрирование	8
3.3	Вычисление матрицы ковариации	9
3.4	Поиск собственных значений матрицы	9
3.5	Вычисление собственных векторов	11
3.6	Вычисление доли объясненной дисперсии	13
3.7	Алгоритм PCA	13
3.8	Вычисляем MSE	15
3.9	Выбор числа главных компонент	16
3.10	Обработка пропущенных значений в данных	16
3.11	Исследуем влияние шума на PCA	17
3.12	Применяем наш PCA к реальному датасету	17
4	Выводы	19

1 Введение

Цель работы

Изучить метод главных компонент (РСА) как инструмент снижения размерности данных, реализовать его основные этапы с нуля на Python, и исследовать его свойства на практике.

Задачи:

1. Изучить математические основы РСА
2. Доказать, что главные компоненты — это собственные векторы матрицы ковариаций, максимизирующие дисперсию данных.
2. Реализовать центрирование данных (вычитание среднего).
3. Вычислить матрицу ковариаций для центрированных данных.
4. Решить СЛАУ методом Гаусса для нахождения собственных векторов.
5. Найти собственные значения матрицы ковариаций методом бисекции.
6. Вычислить собственные векторы для полученных значений.
7. Определить долю объяснённой дисперсии для выбора числа компонент.

Hard Level:

8. Реализовать полный алгоритм РСА (от центрирования до проекции на компоненты).
9. Визуализировать данные в пространстве главных компонент (2D/3D).
10. Оценить ошибку восстановления данных после снижения размерности.

Expert Level:

11. Автоматически выбирать число компонент на основе порога дисперсии.

12. Обрабатывать пропущенные значения в данных перед применением PCA.

13. Исследовать устойчивость PCA к шуму и сравнить результаты до/после добавления шума.

14. Применить PCA к реальному датасету и оценить влияние на качество модели.

15. Проанализировать полученные результаты и сделать выводы.

2 Теория

2.1 Доказать

Направления главных компонент совпадают с собственными векторами матрицы ковариаций C , упорядоченными по убыванию соответствующих собственных значений.

2.2 Доказательство

Для симметричной матрицы C существует ортогональное разложение:

$$C = V\Lambda V^T \quad (1)$$

где:

- $V = [v_1 | v_2 | \dots | v_m]$ - матрица собственных векторов
- $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$ - собственные значения, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$

Дисперсия проекции на направление w :

$$\text{Var}(w) = w^T C w = w^T V \Lambda V^T w = u^T \Lambda u = \sum_{i=1}^m \lambda_i u_i^2 \quad (2)$$

где $u = V^T w$, причем $\|u\| = \|w\| = 1$.

Максимум $\sum_{i=1}^m \lambda_i u_i^2$ при $\sum u_i^2 = 1$ достигается при:

$$u = (1, 0, \dots, 0)^T \Rightarrow w = v_1 \quad (3)$$

что дает максимальную дисперсию λ_1 .

При дополнительном условии ортогональности $w \perp v_1$:

$$u = (0, 1, 0, \dots, 0)^T \Rightarrow w = v_2 \quad (4)$$

с дисперсией λ_2 , и т.д.

Следовательно

- Главные компоненты – собственные векторы v_i матрицы C
- Порядок определяется убыванием λ_i
- Ортогональность следует из симметричности C



3 Практика

3.1 Метод Гаусса

Реализовали стандартный метод Гауса:

* Сначала создаём расширенную матрицу

* Потом приводим матрицу к ступенчатому виду (находим ведущий элемент -> нормируем строку -> переставляем (если надо) -> вычитаем строку)

* Определяем ранг матрицы

* Для единственного решения - извлекаем решение из последнего столбца

* Для случая бесконечного множества решений:

**** Находим базисные и свободные переменные

**** Строим частное решение

**** Строим базисные решения для свободных переменных

```

1 def gauss_solver(A: 'Matrix', b: 'Matrix') -> List['Matrix']:
2     if A.shape[0] != A.shape[1]:
3         raise ValueError('Матрица не квадратная')
4     if A.shape[0] != b.shape[0] or b.shape[1] != 1:
5         raise ValueError('Матрица и вектор несовместимы')
6
7     n = A.shape[0]
8     #расширенная матрица [A|b]
9     augmented = Matrix(n, n + 1)
10    for i in range(1, n + 1):
11        for j in range(1, n + 1):
12            augmented[i, j] = A[i, j]
13        augmented[i, n + 1] = b[i, 1]
14
15    #прямой ход метода Гаусса
16    rank = 0
17    for col in range(1, n + 1):
18        #ищем ненулевого элемента в текущем столбце
19        pivot_row = None
20        for row in range(rank + 1, n + 1):
21            if abs(augmented[row, col]) > 1e-5:
22                pivot_row = row
23                break
24
25        if pivot_row is None:
26            continue #пропускаем нулевой столбец
27
28        #перестановка строк
29        rank += 1
30        if pivot_row != rank:
31            for j in range(col, n + 2):
32                augmented[rank, j], augmented[pivot_row, j] = augmented[
pivot_row, j], augmented[rank, j]
33
34        #нормировка текущей строки
35        norm_val = augmented[rank, col]
36        for j in range(col, n + 2):
37            augmented[rank, j] /= norm_val
38
39        #обнуление оставшихся элементов
40        for row in range(1, n + 1):
41            if row != rank and abs(augmented[row, col]) > 1e-10:
42                factor = augmented[row, col]
43                for j in range(col, n + 2):
44                    augmented[row, j] -= factor * augmented[rank, j]
45
46    for row in range(rank + 1, n + 1):
47        if abs(augmented[row, n + 1]) > 1e-10:
48            raise ValueError('У системы нет решений')
49
50    #единственное решение
51    if rank == n:
52        solution = Matrix(n, 1)
53        for i in range(1, n + 1):
54            solution[i, 1] = augmented[i, n + 1]
55        return [solution]
56
57    #бесконечное множество решений
58    #тогда определяю базисные и свободные переменные
59    basic_vars = []

```

```

60     for row in range(1, rank + 1):
61         for col in range(1, n + 1):
62             if abs(augmented[row, col] - 1) < 1e-10:
63                 basic_vars.append(col)
64                 break
65
66     free_vars = [col for col in range(1, n + 1) if col not in basic_vars]
67     solutions = []
68
69     #частное решение
70     particular = Matrix(n, 1)
71     for row, col in enumerate(basic_vars, 1):
72         particular[col, 1] = augmented[row, n + 1]
73     if any(abs(particular[i, 1]) > 1e-10 for i in range(1, n + 1)):
74         solutions.append(particular)
75
76     #базисные решения
77     for free_col in free_vars:
78         vec = Matrix(n, 1)
79         vec[free_col, 1] = 1
80         for row, basic_col in enumerate(basic_vars, 1):
81             sum_ = 0
82             for j in range(basic_col + 1, n + 1):
83                 sum_ += augmented[row, j] * vec[j, 1]
84             vec[basic_col, 1] = -sum_
85
86         if any(abs(vec[i, 1]) > 1e-10 for i in range(1, n + 1)):
87             solutions.append(vec)
88
89     return solutions

```

Листинг 1: Метод Гауса

3.2 Центрирование

```

1 def center_data(X: 'Matrix') -> 'Matrix':
2     """
3     Вход: матрица данных X (*nm)
4     Выход: центрированная матрица X_centered (*nm)
5     """
6     n_rows, m_cols = X.shape
7     colum_means = []
8
9     for col in range(1, m_cols + 1):
10         colum_sum = 0
11         for row in range(1, n_rows + 1):
12             colum_sum += X[row, col]
13
14         mean_val = colum_sum / n_rows
15         colum_means.append(mean_val)
16
17     X_centered = Matrix(n_rows, m_cols)
18     for row in range(1, n_rows+1):
19         for col in range(1, m_cols+1):
20             tek_val = X[row, col]
21             mean_val_col = colum_means[col-1]
22             X_centered[row, col] = tek_val - mean_val_col

```



```

23
24     return X_centered

```

Листинг 2: Центрирование матрицы

3.3 Вычисление матрицы ковариации

```

1 def covariance_matrix(X_centered: 'Matrix') -> 'Matrix':
2     """
3     Вход: центрированная матрица X_centered (*nm)
4     Выход: матрица ковариаций C (*mm)
5     """
6     n_rows, m_cols = X_centered.shape
7
8     if n_rows == 1:
9         raise ValueError('n равно 1 не( можем избежать деления на ноль)')
10
11     C = Matrix(m_cols, m_cols)
12
13     scale = 1/(n_rows-1)
14
15     for i in range(1, m_cols + 1):
16         for j in range(1, m_cols + 1):
17             res = 0
18             for k in range(1, n_rows + 1):
19                 res += X_centered[k, i]*X_centered[k, j]
20             C[i, j] = res*scale
21
22     return C

```

Листинг 3: Вычисление матрицы ковариации

3.4 Поиск собственных значений матрицы

Собственные значения находятся методом бисекции, который мы применяем к вычислению значения характеристического значения (должен быть = 0).

Шаг, с которым мы идём при бисекции, по умолчанию мы взяли 5000, для больших матриц $m*100$

```

1 from matrix import Matrix, calculate_determinant
2 from typing import List
3 from pca.root_finder import *
4 import math
5
6 def find_eigenvalues(C: 'Matrix', tol: float = 1e-6) -> List[float]:
7     """Вход:
8     C: матрица ковариаций (*mm)
9     tol: допустимая погрешность

```

```

10 Выход: список вещественных собственных значений
11 """
12 m = C.shape[0]
13
14 min_val = float('inf')
15 max_val = -float('inf')
16
17 for i in range(1, m+1):
18     r = sum(abs(C[i, j]) for j in range(1, m+1) if i != j)
19     center = C[i, i]
20     min_val = min(min_val, center - r)
21     max_val = max(max_val, center + r)
22
23 #задаём функцию для вычисления значения характеристического многочлена
24 def charact_polynom(lambd):
25     C_minus_lambd_E = Matrix(m, m)
26     for i in range(1, m+1):
27         for j in range(1, m+1):
28             if i == j:
29                 C_minus_lambd_E[i, j] = C[i, j] - lambd
30             else:
31                 C_minus_lambd_E[i, j] = C[i, j]
32     det = calculate_determinant(C_minus_lambd_E)
33     # print(C_minus_lambd_E)
34     return det
35
36 root_finder = RootFinder(charact_polynom)
37
38 n_intervals = max(5000, 100*m) #количество интервалов для поиска
39 step = (max_val - min_val) / n_intervals
40
41 intervals = []
42 #проверочка на границы
43 x = min_val - step / 2
44 x_next = x + step
45 prev_val = charact_polynom(x)
46
47 eigenvalues = []
48
49 #цикл с шагом step/2, чтобы не пропустить значения на границах отрезка с
    обычным шагом step
50 for _ in range(2 * n_intervals + 1):
51     curr_val = charact_polynom(x_next)
52
53     if prev_val*curr_val <= 0 or abs(prev_val) < tol:
54         intervals.append((x, x_next))
55     x = x + step / 2
56     x_next = x_next + step / 2
57     prev_val = curr_val
58
59 #бисекция для интервалов
60 for a, b in intervals:
61     root, _ = root_finder.bisection(a, b, tol)
62     if root is not None:
63         if not any(abs(root - x) < tol for x in eigenvalues):
64             eigenvalues.append(root)
65
66 if len(eigenvalues) < m:
67     for i in range(1, m + 1):
68         diag_val = C[i, i]

```

```

69         if not any(abs(diag_val - x) < tol for x in eigenvalues):
70             eigenvalues.append(diag_val)
71     eigenvalues.sort(reverse=True)
72
73     return eigenvalues

```

Листинг 4: Поиск собственных значений матрицы методом бисекции

3.5 Вычисление собственных векторов

Собственные векторы мы находим в соответствии с полученными собственными значениями.

Вывод в функции у нас реализован как список пар (вектор, соответствующее ему собственное значение). С таким форматом нам будет дальше удобно работать.

Для РСА нам нужны ортогональные векторы, поэтому мы придумали использовать метод Грамма-Шмидта, он позволяет, сохраняя собственные значения векторов, сделать систему ортогональной. Это происходит за счёт вычисления скалярного произведения, а затем вычитаем проекции. Под конец мы нормализуем векторы.

$$(C - \lambda I)v = 0 \quad (5)$$

где:

- C - исходная матрица
- λ - собственное значение матрицы C
- I - единичная матрица того же размера, что и C
- v - искомый собственный вектор (столбец)
- 0 - нулевой вектор

```

1 from matrix import Matrix
2 from typing import List, Tuple
3 from pca.gauss_solver import gauss_solver
4 import math
5
6 def gram_schmidt(eigen_pairs: List[Tuple[Matrix, float]]) -> List[Tuple[
    Matrix, float]]:

```

```

7      """
8      Применяет процесс Грам-Шмидта- к списку пар вектор(, собственное значение)
9      Возвращает новый список ортогональных векторов с сохранением собственных значений
10     """
11     ortho_pairs = []
12
13     for v, lambda_ in eigen_pairs:
14         #создаем копию вектора для ортогонализации
15         w = Matrix(v.shape[0], 1)
16         for i in range(1, v.shape[0] + 1):
17             w[i, 1] = v[i, 1]
18
19         #вычитаем проекции на уже имеющиеся ортогональные векторы
20         for u, _ in ortho_pairs:
21             #вычисляем скалярное произведение
22             dot_product = sum(w[i, 1] * u[i, 1] for i in range(1, w.shape
[0] + 1))
23             #вычитаем проекцию
24             for i in range(1, w.shape[0] + 1):
25                 w[i, 1] -= dot_product * u[i, 1]
26
27         #нормализуем вектор
28         norm = math.sqrt(sum(w[i, 1] ** 2 for i in range(1, w.shape[0] +
1)))
29         if norm > 1e-10:
30             for i in range(1, w.shape[0] + 1):
31                 w[i, 1] /= norm
32             ortho_pairs.append((w, lambda_))
33
34     return ortho_pairs
35
36 def find_eigenvectors(C: 'Matrix', eigenvalues: List[float]) -> List[Tuple
[Matrix, float]]:
37     """
38     Находит собственные векторы матрицы C для заданных собственных значений.
39     Вход:
40     C: матрица ковариаций (*mm)
41     eigenvalues: список собственных значений
42     Выход: список собственных векторов каждый( вектор - объект Matrix размера *m1)
43     """
44     m = C.shape[0]
45     eigen_pairs = []
46
47     for lambda_ in eigenvalues:
48         #создаем матрицу (C - lambda*I)
49         A = Matrix(m, m)
50         b = Matrix(m, 1) #нулевой вектор правой части
51
52         for i in range(1, m + 1):
53             for j in range(1, m + 1):
54                 if i == j:
55                     A[i, j] = C[i, j] - lambda_
56                 else:
57                     A[i, j] = C[i, j]
58
59         #решаем систему (C - lambda*I)v = 0
60         solutions = gauss_solver(A, b)
61
62         if solutions:
63             for sol in solutions:

```

```

64         eigen_pairs.append((sol, lambda_))
65
66     ortho_pairs = gram_schmidt(eigen_pairs)
67     return ortho_pairs

```

Листинг 5: Поиск собственных значений матрицы методом бисекции

3.6 Вычисление доли объясненной дисперсии

Вычисляем долю, где числитель представляет сумму дисперсий, объясняемых первыми k главными компонентами. Знаменатель - общая сумма всех собственных значений (полная дисперсия). Результат γ показывает, какая доля общей дисперсии объясняется выбранными компонентами.

$$\gamma = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^m \lambda_i} \quad (6)$$

где:

- λ_i - собственные значения матрицы ковариации
- k - количество учитываемых главных компонент
- m - общее количество собственных значений

```

1 def explained_variance_ratio(eigenvalues: List[float], k: int) -> float:
2     """
3     Вход:
4     eigenvalues: список собственных значений
5     k: число компонент
6     Выход: доля объяснённой дисперсии
7     """
8     eigenvalues.sort(key=int, reverse=True)
9     return sum(eigenvalues[:k]) / sum(eigenvalues)

```

Листинг 6: Вычисление доли объединённой дисперсии

3.7 Алгоритм PCA

1. Центрирование данных.

2. Вычисление матрицы выборочных ковариаций.
3. Нахождение собственных значений и векторов.
4. Проекция данных на главные компоненты.

Для удобства обращения мы создали класс для хранения всех промежуточных объектов

```

1 class PCA_DATA:
2     mean_columns = []
3     centered_data = None
4     covariance_matrix = None
5     X_proj = None
6     eigenvalues = []
7     eigenvectors_matrix = None
8     variance = 0
9
10    def __init__(self, **kwargs) -> None:
11        for key, value in kwargs.items():
12            setattr(self, key, value)

```

Листинг 7: Класс для хранения объектов

```

1 def pca(X: Matrix, k: int, show_logs=True) -> PCA_DATA:
2     """
3     Вход:
4     X: матрица данных (*nm)
5     k: число главных компонент
6     Выход:
7     X_proj: проекция данных (*nk)
8     : доля объяснённой дисперсии
9     """
10    n, m = X.shape
11    centered_data, mean_columns = center_data(X)
12    if show_logs:
13        print("Отцентрированы данные...")
14    cov_matrix = covariance_matrix(centered_data)
15    if show_logs:
16        print("Найдена матрица ковариаций...")
17    eigenvalues = find_eigenvalues(cov_matrix)
18    if show_logs:
19        print("Найдены собственные значения")
20    eigenvectors = find_eigenvectors(cov_matrix, eigenvalues)
21    if show_logs:
22        print("Найдены собственные векторы")
23
24    eigenvectors = list(sorted(eigenvectors, key=lambda x: x[1], reverse=True))
25
26    eigenvalues = [round(i[1], 6) for i in eigenvectors]
27    eigenvectors = [i[0] for i in eigenvectors]
28
29    eigenvectors_matrix = Matrix(m, k)
30    for i in range(k):
31        for j in range(eigenvectors[i].shape[0]):
32            eigenvectors_matrix[j + 1, i + 1] = eigenvectors[i][j + 1, 1]

```

```

33
34 X_proj = centered_data * eigenvectors_matrix
35 if show_logs:
36     print("Найдена проекция...")
37 variance = explained_variance_ratio(eigenvalues, k)
38 if show_logs:
39     print("найден дисперсия...\n")
40
41 pca_data = PCA_DATA(
42     mean_columns=mean_columns,
43     eigenvalues=eigenvalues,
44     eigenvectors_matrix=eigenvectors_matrix,
45     covariance_matrix=cov_matrix,
46     variance=variance,
47     X_proj=X_proj,
48     centered_data=centered_data,
49 )
50
51 return pca_data
52
53 def reconstruct_X(pca_data: PCA_DATA) -> Matrix:
54     scaled_data = pca_data.X_proj * pca_data.eigenvectors_matrix.transpose()
55     for row in range(1, scaled_data.shape[0] + 1):
56         for col in range(1, scaled_data.shape[1] + 1):
57             scaled_data[row, col] += pca_data.mean_columns[col - 1]
58     return scaled_data

```

Листинг 8: PCA и функция восстановления матрицы

3.8 Вычисляем MSE

Вычисляем среднеквадратическую ошибку восстановления данных

$$\text{MSE} = \frac{1}{n \cdot m} \sum_{i=1}^n \sum_{j=1}^m (X_{\text{orig}}^{(i,j)} - X_{\text{recon}}^{(i,j)})^2 \quad (7)$$

где:

- n - количество строк в матрице данных
- m - количество столбцов в матрице данных
- $X_{\text{orig}}^{(i,j)}$ - элемент исходной матрицы в строке i , столбце j
- $X_{\text{recon}}^{(i,j)}$ - элемент восстановленной матрицы в строке i , столбце j

```

1 def reconstruction_error(X_orig: Matrix, X_recon: Matrix) -> float:
2     """
3     Возвращает среднеквадратическую ошибку восстановления
4     """
5     MSE = 1 / (X_orig.shape[0] * X_orig.shape[1])
6     sum_error = 0
7     for row in range(1, X_orig.shape[0] + 1):
8         for col in range(1, X_orig.shape[1] + 1):
9             sum_error += (X_orig[row, col] - X_recon[row, col]) ** 2
10    MSE *= sum_error
11    return MSE

```

Листинг 9: Функция для расчёта MSE

3.9 Выбор числа главных компонент

Мы реализовали автоматический выбор числа главных компонент на основе порога объяснённой дисперсии (далее интегрируем в полный алгоритм pca).

```

1 def auto_select_k(eigenvalues: List[float], threshold: float = 0.95) ->
2     int:
3     """
4     Вход:
5     eigenvalues: список собственных значений
6     threshold: порог объяснённой дисперсии
7     Выход: оптимальное число главных компонент k
8     """
9     for i in range(1, len(eigenvalues) + 1):
10         if explained_variance_ratio(eigenvalues, i) >= threshold:
11             return i
12     return -1

```

Листинг 10: Функция выбора числа главных компонент

3.10 Обработка пропущенных значений в данных

```

1 def handle_missing_values(X: 'Matrix') -> 'Matrix':
2     """
3     Вход: матрица данных X (*nm) с возможными NaN
4     Выход: матрица данных X_filled (*nm) без NaN
5     """
6     mean_values = [0 for _ in range(X.shape[1])]
7
8     for col in range(1, X.shape[1] + 1):
9         colum_sum = 0
10        colum_not_null_count = 0
11        for row in range(1, X.shape[0] + 1):
12            if X[row, col] is not None:
13                colum_sum += X[row, col]
14                colum_not_null_count += 1

```


Ошибка без шума: 0.10040878429494081
Ошибка с шумом: 0.10560372373406768

Рис. 1: Сравнение метрики Ассигасу с шумом и без

```
15     if colum_sum == 0:
16         for row in range(1, X.shape[0] + 1):
17             X[row, col] = 0
18     else:
19         for row in range(1, X.shape[0] + 1):
20             if X[row, col] is not None:
21                 continue
22             X[row, col] = colum_sum / colum_not_null_count
23
24     return X
```

Листинг 11: Функция выбора числа главных компонент

3.11 Исследуем влияние шума на PCA

Добавляем Гауссовский шум к данным. Сравниваем результаты PCA с помощью метрики Ассигасу.

```
1 def add_noise_and_compare(X: 'Matrix', noise_level: float = 0.1):
2     """
3     Вход:
4     X: матрица данных (*nm)
5     noise_level: уровень шума доля( от стандартного отклонения)
6     Выход: результаты PCA до и после добавления шума.
7     В этом задании можете проявить творческие способности, поэтому выходные данные не
8     →
9     типизированы.
10    """
11    for row in range(1, X.shape[0] + 1):
12        for col in range(1, X.shape[1] + 1):
13            X[row, col] += random.gauss(0, noise_level)
14    return X
```

Листинг 12: Создаём шум (гауссовский)

3.12 Применяем наш PCA к реальному датасету

Датасет мы взяли с Kaggle

<https://www.kaggle.com/dongedong/seed-from-uci>

Наблюдения из графика:

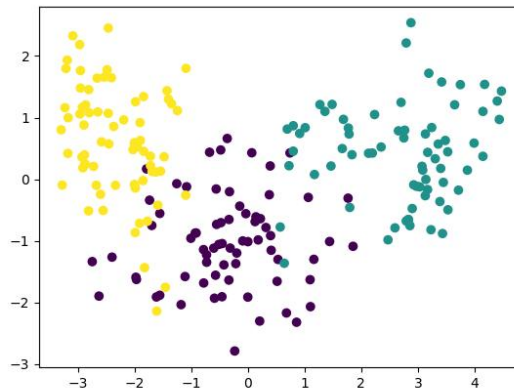


Рис. 2: Проекция данных на первые две компоненты

```
----- Подсчет метрик -----
Accuracy scores:
Без PCA: 0.9428571428571428
Мой PCA: 0.9476190476190476
Sklearn PCA: 0.9142857142857143
```

Рис. 3: Сравнений метрик Ассигасу (наша больше на 0.03!!)

Данные распределены вдоль двух осей: PC1 (горизонтальная) и PC2 (вертикальная)

Основная вариация данных идёт по горизонтали (PC1), так как разброс точек по этой оси значительно больше

Точки образуют вытянутое "облако" с наклоном примерно в 45 градусов

Наибольшая плотность данных сосредоточена в центре (около точки 0,0)

Есть несколько выбросов в верхнем правом и нижнем левом углах

4 Выводы

В этой работе мы разобрались, как работает метод главных компонент (РСА) — мощный инструмент для уменьшения размерности данных. Вот что у нас получилось:

Что сделано:

Написали функции для всех этапов РСА: центрирование данных, расчёт ковариаций, поиск собственных значений и векторов.

Реализовали проекцию данных на главные компоненты и оценку качества через долю объяснённой дисперсии и ошибку восстановления.

Ко всем функциям написали тесты.

Добавил визуализацию и обработку пропущенных значений.

Выбрали метрику и проанализировали результаты с РСА.

Много подебажили.

Что узнали:

РСА эффективно снижает размерность данных, сохраняя основную информацию.

Главные компоненты — это направления максимальной дисперсии, заданные собственными векторами.

Автоматический выбор k через порог дисперсии упрощает анализ.

Шум ухудшает качество РСА, но метод остается устойчивым при умеренных уровнях.

Применение к реальным данным (датасету seeds) подтвердило практическую пользу.

Итог: Работа помогла нам понять математику РСА и научиться применять его на практике.

Ссылка на исходный код:

<https://github.com/bobrnebobr/MatrixLinalgLab.git>