ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ «НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ ИТМО»

ФАКУЛЬТЕТ ТЕХНОЛОГИЙ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА



Отчёт по лабораторной работе № 1 по дисциплине "Линейная алгебра и обработка данных"

Выполнили студенты:

Воробьев Андрей Павлович

гр. J3111 ИСУ 465440

Мавров Артём Николаевич

гр. J3113 ИСУ 466574

Отчет сдан: 23.04.2025

Санкт-Петербург 2025

Содержание

Вве	дение	3
Теоретическая часть 4		4
.1	Математическое обоснование	4
Практическая часть 6		
.1	Метод Гаусса для решения СЛАУ	6
.2	Центрирование данных	8
.3	Вычисление матрицы ковариаций	8
.4	Поиск собственных значений	9
.5	Поиск собственных векторов	11
.6	Вычисление доли объяснённой дисперсии	12
.7	Полный алгоритм РСА	12
.8	Визуализация проекций данных на первые две главные	
	компоненты	13
.9	Вычисление среднеквадротичной ошибки восстановления	
	данных	13
.10	Автоматический выбор числа главных компонент на основе порога объяснённой дисперсии	14
.11	Обработка пропущенных значений в данных	14
.12	Иследование влияния шума на РСА	15
.13	Применение РСА к реальному датасету	16
.14	Тестирование	17
Заключение		18

Введение

В данной лабороторной работе предстоит математически обосновать, что оптимальные направления PCA совпадают с собственными векторами матрицы ковариаций, реалезовать полный цикл PCA без использования сторонних библиотек, дополнить функционал PCA, исследовать влияние шума и применить PCA к реальным данным.

Данная лабороторная работа была выполнена на языке Python из-за лёгкого синтаксиса, простой возможностью расширения, а также в случае необходисости ускорения работы программы отдельные сложные части когда могут быть перенесены на C++, C.

Реализация представлена в виде двух классов. Первый класс Matrix необходим для хранения и выполнения базовых операций с разряженными матрицами. Во втором классе PCA непосредственно представлен полный цикл PCA и дополнительные функции. Также были сделаны тесты.

Полный листинг кода можно найти по ссылке.

Теоретическая часть

Математическое обоснование

Постановка задачи

Пусть $X \in \mathbb{M}at_{n \times m}(R)$ — матрица данных, центрированная по столбцам. Нас интересует вектор направления $w \in \mathbb{R}^m$, вдоль которого проекции наблюдений обладают наибольшей дисперсией.

Для отдельного наблюдения x_i проекция на направление w равна

$$z_i = x_i^{\mathsf{T}} w,$$

а вектор всех таких проекций записывается как

$$z = Xw$$
.

Поскольку строки X уже центрированы, дисперсия проекций может быть выражена через евклидову норму:

$$Var(z) = \frac{1}{n} \sum_{i=1}^{n} z_i^2 = \frac{1}{n} ||Xw||^2 = \frac{1}{n} w^{\top} X^{\top} Xw.$$

Обозначив ковариационную матрицу как $\Sigma = \frac{1}{n} X^{\top} X$, получим:

$$\operatorname{Var}(z) = w^{\top} \Sigma w.$$

Теперь задача сводится к максимизации квадратичной формы при условии нормированности вектора:

$$\max_{\|w\|=1} w^{\top} \Sigma w.$$

Свойства матрицы Σ

Пусть $X \in \mathbb{R}^n$ — случайный вектор с математическим ожиданием $\mathbb{E}[X]$. Тогда матрица ковариаций определяется как

$$\Sigma = \operatorname{Cov}(X) = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T].$$

Рассмотрим квадратичную форму $x^T \Sigma x$ для произвольного вектора $x \in \mathbb{R}^n$:

$$x^T \Sigma x = x^T \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T]x.$$

Благодаря линейности математического ожидания можно записать:

$$x^T \Sigma x = \mathbb{E}[x^T (X - \mathbb{E}[X])(X - \mathbb{E}[X])^T x].$$

Заметим, что

$$x^{T}(X - \mathbb{E}[X])(X - \mathbb{E}[X])^{T}x = ((X - \mathbb{E}[X])^{T}x)^{2},$$

то есть:

$$x^T \Sigma x = \mathbb{E}[(x^T (X - \mathbb{E}[X]))^2] \ge 0.$$

Тогда ковариационная матрица Σ симметрична и неотрицательно определённа, а значит, диагонализируема в ортонормальном базисе собственных векторов v_1, \ldots, v_m с соответствующими собственными значениями $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_m \geq 0$.

Разложение вектора w по базису собственных векторов

Любой вектор w можно представить в виде линейной комбинации собственных векторов:

$$w = \sum_{i=1}^m \alpha_i v_i$$
, где $\sum_{i=1}^m \alpha_i^2 = 1$.

Подставляя это в выражение $w^{\top}\Sigma w$, получаем:

$$w^{\top} \Sigma w = \left(\sum_{i=1}^{m} \alpha_i v_i\right)^{\top} \Sigma \left(\sum_{j=1}^{m} \alpha_j v_j\right) = \sum_{i=1}^{m} \sum_{j=1}^{m} \alpha_i \alpha_j v_i^{\top} \Sigma v_j.$$

Так как v_j — собственный вектор: $\Sigma v_j = \lambda_j v_j$, и $v_i^\top v_j = \delta_{ij}$, то:

$$v_i^{\top} \Sigma v_i = \lambda_i \delta_{ii},$$

что даёт:

$$w^{\top} \Sigma w = \sum_{i=1}^{m} \alpha_i^2 \lambda_i.$$

Оптимальный выбор w

Так как $\sum \alpha_i^2 = 1$, выражение $w^\top \Sigma w$ — это выпуклая комбинация собственных значений. Она достигает максимума, когда весь вес сосредоточен на наибольшем собственном значении:

$$\alpha_1 = 1, \quad \alpha_2, \dots, \alpha_m = 0 \quad \Rightarrow \quad w = v_1.$$

Остальные компоненты главных направлений

Аналогично, чтобы найти второе направление с максимальной дисперсией, необходимо искать максимум того же функционала на подпространстве, ортогональном v_1 . Тогда оптимальным будет v_2 , и далее по порядку.

Вывод: Оптимальные направления PCA — это собственные векторы ковариационной матрицы Σ , упорядоченные по убыванию собственных значений.

Практическая часть

Метод Гаусса для решения СЛАУ

Метод Гаусса реалезован классически с оптимезацией перестановки строк. На каждом шаге і выбирается введущая строка с максимальным элементом в позиции [i, i], для уменьшения погрешности. При необходимости строки меняются местами, только вместо явной перестановки строк, меняются лишь два числа в списке order, за счёт чего этот процесс значительно ускоряется. При дальнейшем обращении к матрице по индексу для индекса строки необходимо сначала обращаться к order, если порядок важен. После выбора ведущей строки и перестоновки из каждой строки ниже і вычетается ведущая строка домноженая на factor. После преведения матрицы к верхнетреугольному виду в зависимости от количетсва нулевых строк даётся либо частное либо общее решение. В случае частного ршения каждая неизвестная считается сверху вниз. В случае общего сначала определяются свободные переменные и анологично с частным для каждого вектора с свободной переменной считается искомый вектор неизвестных.

На вход матрица A и вектор b из класса Matrix, где A - матрица коэффицентов, b - вектор неизвестных.

```
def gauss_solver(A: 'Matrix', b: 'Matrix') -> List['Matrix']:
           order = list(range(A.num_rows))
           for i in range(A.num_rows - 1):
               max_value = 0.0
               col = i - 1
6
               while round(abs(max_value), A.accuracy) == 0.0 and col !=
                   A.num_rows:
                   col += 1
                   for j in range(i, A.num_rows):
                       value = A[order[j] + 1, col + 1]
                       if abs(value) > abs(max_value):
11
                            max_value = value
                            pivot_row = j
13
               if round(max_value, A.accuracy) == 0.0:
                   continue
               if order[pivot_row] != i:
18
                   order[i], order[pivot_row] = order[pivot_row], order[
                      i]
20
               for j in range(i + 1, A.num_rows):
21
                   factor = A[order[j] + 1, col + 1] / max_value
                   if round(factor, A.accuracy) == 0.0:
                        continue
24
                   for k in range(col, A.num_rows):
26
                       A[order[j] + 1, k + 1] -= factor * A[order[i] +
                           1, k + 1
                   b[order[j] + 1, 1] -= factor * b[order[i] + 1, 1]
30
           A.accuracy = 5
31
```

```
amt_null_rows = 0
32
           mn_null_row = A.num_rows
33
           for i in range(A.num_rows):
34
                flag_null_row = True
               for j in range(i, A.num_columns):
                    if round(A[order[i] + 1, j + 1], A.accuracy) != 0.0:
                        flag_null_row = False
38
                        break
39
                if flag_null_row:
40
                    mn_null_row = min(i, mn_null_row)
41
                    amt_null_rows += 1
42
43
           if amt_null_rows == 0:
44
               Ans = Matrix(f"{A.num_rows} 1")
45
               Ans[A.num_rows, 1] = b[order[A.num_rows - 1] + 1, 1] / A[
46
                   order[A.num_rows - 1] + 1, A.num_columns]
47
               for i in range(A.num_rows - 2, -1, -1):
                    value = 0
49
50
                    for j in range(A.num_columns, i + 1, -1):
                        value += A[order[i] + 1, j] * Ans[j, 1]
                    Ans[i + 1, 1] = (b[order[i] + 1, 1] - value) / A[
54
                       order[i] + 1, i + 1]
                basis = []
56
                basis.append(Ans)
58
           else:
59
               null_on_diag = []
               last_ind = A.num_rows
61
               for i in range(mn_null_row - 1, -1, -1):
                    for j in range(i, last_ind):
63
                        if round(A[order[i] + 1, j + 1], A.accuracy) !=
                           0.0:
                            for k in range(j + 2, last_ind + 1):
65
                                 null_on_diag.append(k - 1)
                             last_ind = j
67
                            break
69
               basis = []
70
               for i in range(amt_null_rows):
71
                    basis.append(Matrix(f"{A.num_rows} 1"))
72
73
               j = 0
               for i in null_on_diag:
                    basis[j][i + 1, 1] = 1
76
                    j += 1
78
               for k in range(amt_null_rows):
                    for i in range(null_on_diag[k] - 1, -1, -1):
                        if i in null_on_diag:
81
```

```
82
                             continue
83
                        first_not_null = A.num_columns
84
                        for j in range(i, A.num_columns):
85
                             if round(A[order[i] + 1, j + 1]) != 0.0:
                                 first_not_null = j
                                 break
88
89
                        value = 0
90
                        for j in range(A.num_columns - 1, first_not_null,
                             -1):
                             value += A[order[i] + 1, j + 1] * basis[k][j
                                + 1, 1]
93
                        if value != 0:
94
                             basis[k][i + 1, 1] = (-value) / A[order[i] +
                                1, first_not_null + 1]
97
           return basis
```

Центрирование данных

Центрирвоание данных происходит по формуле $X_{centered} = X - mean(X)$. Реализация разбита на две функции: mean_vector которая вычесляет вектор средних значений для каждой колонки матрицы X и center_data, которая возвращает центрированную матрицу.

```
def mean_vector(X: 'Matrix') -> 'Matrix':
          mean_vector = Matrix(f"{X.num_columns} 1")
          for i in range(1, X.num_columns + 1):
3
               mean\_vector[i, 1] = sum(X[j, i] for j in range(1, X.
                  num_rows + 1)) / X.num_rows
          return mean_vector
  def center_data(X: 'Matrix') -> 'Matrix':
      X_vector = PCA.mean_vector(X)
9
      X_centered = Matrix(f"{X.num_rows} {X.num_columns}")
      for i in range(1, X.num_rows + 1):
           for j in range(1, X.num_columns + 1):
               X_{centered}[i, j] = X[i, j] - mean_X_vector[j, 1]
14
      return X_centered
```

Вычисление матрицы ковариаций

Вычисления происходит по формуле: $C = \frac{1}{n-1}X^TX$.

```
covariance_matrix /= X_centered.num_rows - 1

return covariance_matrix
```

Поиск собственных значений

Поиск собственных значений осуществляется посредству применения формулы:

$$\det(C - \lambda I) = 0$$

Алгоритм нахождения собственных значений основан на нескольких идеях. Так как по ТЗ требуется находить значения методом бисекции, необходимо сначала определять интервал, на котором будут исследоваться отрезки.

Границы интервала получаются из информации о матрице. Известно, что у матрицы ковариаций вещественные и неотрицательные корни. Из этого делается вывод, что нижняя граница -0, а верхняя — след матрицы (это берётся из теоремы Виета для характеристического многочлена).

Из-за сложностей с нахождением собственных значений было решено дополнительно разбивать интервал с помощью уже найденных собственных значений. Для этого в начале реализована проверка: значение 1 добавляется в список, если оно является корнем характеристического многочлена. Также отдельно проверяется 0, так как это частое значение, которое тяжело заметить с помощью обычной бисекции. Это особенно актуально, так как у ковариационных матриц часто наблюдаются нулевые собственные значения из-за линейной зависимости признаков.

Большие интервалы разбиваются на маленькие отрезки в количестве num_intervals. Далее для каждого такого отрезка считается значение характеристического многочлена в начале и в конце. Если значения имеют противоположные знаки, значит на этом отрезке гарантированно есть корень, и запускается алгоритм бисекции. Метод постепенно сужает отрезок, пока длина интервала не станет меньше tolerance. После этого значение округляется и добавляется в список найденных собственных значений.

Если после прохода найдено меньше собственных значений, чем число строк матрицы (т.е. возможное число независимых признаков), количество миниинтервалов увеличивается в 10 раз. Это делается для того, чтобы попытаться слишком близкие друг к другу собственные значения, которые могли быть пропущены при грубом разбиении.

Когда достигается максимальное количество миниинтервалов, алгоритм останавливается. Это нужно так как могут быть кратные собственные значения.

```
def characteristic_polynomial_value(covariance_matrix: 'Matrix',
     value: float) -> float:
           matrix = deepcopy(covariance_matrix)
           for row in range(1, matrix.num_rows + 1):
               matrix[row, row] -= value
           return matrix.get_determinant()
6
  def find_eigenvalues(covariance_matrix: 'Matrix',
8
                        tolerance: Optional[float] = 1e-10,
                        intial_num_intervals: Optional[int] = 10,
                        max_num_intervals: Optional[int] = 1000) -> List
11
                           [float]:
       eigenvalues = []
12
      matrix = deepcopy(covariance_matrix)
13
```

```
14
       zero_eigenvalue = \
           (PCA.characteristic_polynomial_value(matrix, -tolerance)
16
            * PCA.characteristic_polynomial_value(matrix, tolerance) <=
            or PCA.characteristic_polynomial_value(matrix, 0) < 1e-6)
       one_eigenvalue = abs(PCA.characteristic_polynomial_value(matrix,
19
          1)) < tolerance
       eigenvalues.append(1)
20
21
       num_intervals = intial_num_intervals
23
       while len(eigenvalues) - 1 != matrix.num_rows - zero_eigenvalue
24
          and num_intervals <= max_num_intervals:</pre>
           for interval_index in range(len(eigenvalues) + 1):
               if interval_index == 0:
26
                   lower_bound = tolerance
                   higher_bound = eigenvalues[0] - tolerance
               elif interval_index == len(eigenvalues):
29
                   lower_bound = eigenvalues[-1] + tolerance
30
                   higher_bound = matrix.get_tracer() + tolerance
               else:
                    lower_bound = eigenvalues[interval_index - 1] +
                       tolerance
                   higher_bound = eigenvalues[interval_index] -
34
                       tolerance
               if higher_bound - lower_bound < tolerance:</pre>
35
                    continue
37
               interval_length = (higher_bound - lower_bound) /
38
                  num intervals
               current_num_intervals = num_intervals
39
               if interval_length < tolerance:</pre>
40
                    current_num_intervals = int((higher_bound -
41
                       lower_bound) / tolerance) + 1
                    interval_length = (higher_bound - lower_bound) /
42
                       num_intervals
43
               for interval in range(current_num_intervals):
44
                    if len(eigenvalues) - 1 == matrix.num_rows -
45
                       zero_eigenvalue:
                        break
46
                    low = lower_bound + interval * interval_length
47
                   high = lower_bound + (interval + 1) * interval_length
48
                    characteristic_value_low = PCA.
49
                       characteristic_polynomial_value(matrix, low)
                    characteristic_value_high = PCA.
                       characteristic_polynomial_value(matrix, high)
                    if characteristic_value_low == 0:
                        eigenvalues.append(round(low, int(math.log10(1 /
                           tolerance))))
                        continue
```

```
if characteristic_value_high == 0:
                        eigenvalues.append(round(high, int(math.log10(1 /
56
                             tolerance))))
                        continue
                    if characteristic_value_low *
                       characteristic_value_high < 0:</pre>
                        if characteristic_value_low > 0:
                             low, high = high, low
                        eigenvalue = (low + high) / 2
                        while True:
                             characteristic_value = PCA.
                                characteristic_polynomial_value(matrix,
                                eigenvalue)
                             if abs(high - low) < tolerance / 10:</pre>
65
                                 eigenvalues.append(round(eigenvalue, int(
66
                                    math.log10(1 / tolerance))))
                                 eigenvalues.sort()
                                 break
68
                             if characteristic_value > 0:
69
                                 high = eigenvalue
                             else:
71
                                 low = eigenvalue
                             eigenvalue = (low + high) / 2
73
           num_intervals *= 10
74
       if zero_eigenvalue:
76
           eigenvalues.append(0)
       if not one_eigenvalue:
           eigenvalues = [i for i in eigenvalues if i != 1]
80
       return sorted(eigenvalues)
81
```

Поиск собственных векторов

Собственные вектора вычесляются по формуле: $(C - \lambda I)v = 0$, где λ - собственное число, C - матрица ковариаций, I - единичная матрица.

```
def find_eigenvectors(C: 'Matrix', eigenvalues: List[float]) -> Dict[
    float, List['Matrix']]:
        eigenvectors = {}
        b = Matrix(f"{C.num_rows} 1")
        I = Matrix.eye(C.num_rows)

for eigenvalue in eigenvalues:
        A = C - (float(eigenvalue) * I)
        eigenvectors[eigenvalue] = PCA.gauss_solver(A, b)

return eigenvectors
```

Вычисление доли объяснённой дисперсии

Доля объяснённой дисперсии считается по формуле:

$$\gamma = \frac{\sum_{i=0}^{k} \lambda_i}{\sum_{i=0}^{m} \lambda_i}$$

```
def explained_variance_ratio(eigenvalues: List[float], k: int) ->
    float:
    return sum(eigenvalues[-k:]) / sum(eigenvalues)
```

Полный алгоритм РСА

Полный алгоритм РСА состоит из:

- 1. Центрирования данных.
- 2. Вычисления матрицы выборочных ковариаций.
- 3. Нахождения собственных значений и векторов.
- 4. Проекции данных на главные компоненты.

```
def get_components(eigenvectors: Dict[float, List['Matrix']], k: int)
      -> 'Matrix':
           components = Matrix(f"{eigenvectors[next(iter(eigenvectors))
              ][0].num_rows} {k}")
           count = 0
           for eigenvalue in sorted(eigenvectors.keys(), reverse=True):
               if count == k:
                   break
               for vector in eigenvectors[eigenvalue]:
                   for element in range(1, vector.num_rows + 1):
                       components[element, count + 1] = vector[element,
                          1]
                   count += 1
                   if count == k:
                       break
12
13
           return components
14
  def RSA(X: 'Matrix', k: int) -> Tuple['Matrix', 'Matrix', float]:
           centered_X = PCA.center_data(X)
17
           covariance_matrix = PCA.covariance_matrix(centered_X)
18
           eigenvalues = PCA.find_eigenvalues(covariance_matrix)
19
           eigenvectors = PCA.find_eigenvectors(covariance_matrix,
20
              eigenvalues)
           for eigenvalue in eigenvectors.keys():
21
               eigenvectors[eigenvalue] = [vector.norm_vector() for
                  vector in eigenvectors[eigenvalue]]
           components = PCA.get_components(eigenvectors, k)
23
           projection = centered_X * components
           variance = PCA.explained_variance_ratio(eigenvalues, k)
26
           return projection, components, variance
27
```

Визуализация проекций данных на первые две главные компоненты

Для визуализации используется библиотка Matplotlib.

```
def plot_pca_projection(X_proj: 'Matrix') -> Figure:
          x = [X_proj[i, 1] for i in range(1, X_proj.num_rows + 1)]
2
          y = [X_proj[i, 2] for i in range(1, X_proj.num_rows + 1)]
          fig, ax = plt.subplots(figsize=(8, 6))
           ax.scatter(x, y, s=30, color='steelblue', alpha=0.7,
6
              edgecolors='k', linewidths=0.5)
           for i in range(len(x)):
               ax.text(x[i] + 5, y[i], str(i + 1), fontsize=9, color='
                  darkred')
          ax.set_title("Projection onto the first 2 principal
9
              components", fontsize=14)
          ax.set_xlabel("Main component 1", fontsize=12)
          ax.set_ylabel("Main component 2", fontsize=12)
          ax.grid(True)
          ax.set_aspect('equal', adjustable='box')
13
14
          return fig
```

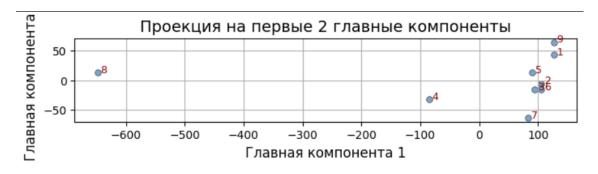


Рис. 1: * Пример визуализации

Вычисление среднеквадротичной ошибки восстановления данных

Среднеквадротичная ошибка вычисляется по формуле:

$$MSE = \frac{1}{nm} \sum_{i=0}^{n} \sum_{j=0}^{m} (X_{orig}^{ij} - X_{recon}^{ij})^{2}$$

Реализация разбита на две функции reconstruction восстанавливает исходную матрицу из проекции, матрицы компонентв, вектора средних значений. reconstruction_error считает среднеквадротичную ошибку по вышепреведённой формуле.

```
def reconstruction(projection: 'Matrix', components: 'Matrix',
     mean_vector: 'Matrix') -> 'Matrix':
          X_recon = projection * components.transpose()
          for i in range(1, X_recon.num_rows + 1):
               for j in range(1, X_recon.num_columns + 1):
                   X_recon[i, j] += mean_vector[j, 1]
          return X_recon
  def reconstruction_error(X_orig: 'Matrix', X_recon: 'Matrix') ->
     float:
      error = 0
      for i in range(1, X_orig.num_rows + 1):
          for j in range(1, X_orig.num_columns + 1):
               error += (X_orig[i, j] - X_recon[i, j]) ** 2
13
14
      return error / (X_orig.num_rows * X_orig.num_columns)
```

Автоматический выбор числа главных компонент на основе порога объяснённой дисперсии

Автоматический выбор числа главных компонент происходит по формуле:

$$k = \left\{ k : \frac{\sum_{i=1}^{k} \lambda_i}{\sum_{i=1}^{m} \lambda_i} \ge threshold \right\}$$

Обработка пропущенных значений в данных

Пропущенные значения обробатывались по следующему принципу:

$$X_{filled} = \left\{ egin{array}{ll} X^{ij} & ext{,если } X^{ij}
eq NaN \ mean(X) & ext{,иначе.} \end{array}
ight.$$

```
size_string = input_string[0].split()
           X.num_rows = int(size_string[0])
6
           X.num_columns = int(size_string[1])
           X.square = X.num_rows == X.num_columns
           X.values = []
           X.row_sizes = [0]
           X.column_indices = []
           X.accuracy = int(input_string[1])
13
14
           missing_values = []
           columns_sums = [0] * X.num_columns
16
           columns_sizes = [0] * X.num_columns
18
           for row in range(X.num_rows):
19
               row_string = input_string[row + 2].split()
20
               X.row_sizes.append(X.row_sizes[-1])
21
               for column in range(X.num_columns):
                    element = row_string[column]
23
                    if element == 'nan':
24
                        missing_values.append((row, column))
                    else:
                        element = round(float(element), X.accuracy)
                        if element != 0:
28
                            X.values.append(element)
29
                            X.column_indices.append(column)
30
                            X.row_sizes[-1] += 1
31
                            columns_sums[column] += element
                            columns_sizes[column] += 1
33
           columns_means = []
35
           for sum, size in zip(columns_sums, columns_sizes):
36
               if size != 0:
                    columns_means.append(sum / size)
               else:
                    columns_means.append(0)
40
           for row, column in missing_values:
41
               X[row + 1, column + 1] = columns_means[column]
42
43
           return X
```

Иследование влияния шума на РСА

Шум добовляется посредством генерации матрицы случайных чисел, которые зависят от среднего вектора и введённого параметра noise_level. Далее эта матрица прибавляется к матрице к которой мы хотим добовить шум. Применяется алгоритм РСА и вычесляется среднеквадратичная ошибка.

```
X.num_rows = int(size_string[0])
6
           X.num_columns = int(size_string[1])
           X.square = X.num_rows == X.num_columns
9
           X.values = []
           X.row_sizes = [0]
           X.column_indices = []
12
           X.accuracy = int(input_string[1])
13
14
           missing_values = []
           columns_sums = [0] * X.num_columns
16
           columns_sizes = [0] * X.num_columns
17
18
           for row in range(X.num_rows):
19
               row_string = input_string[row + 2].split()
20
               X.row_sizes.append(X.row_sizes[-1])
21
               for column in range(X.num_columns):
                    element = row_string[column]
                    if element == 'nan':
24
                        missing_values.append((row, column))
25
                    else:
26
                        element = round(float(element), X.accuracy)
27
                        if element != 0:
                            X.values.append(element)
29
                            X.column_indices.append(column)
30
                            X.row_sizes[-1] += 1
31
                             columns_sums[column] += element
32
                             columns_sizes[column] += 1
34
           columns_means = []
35
           for sum, size in zip(columns_sums, columns_sizes):
36
               if size != 0:
37
                    columns_means.append(sum / size)
               else:
                    columns_means.append(0)
           for row, column in missing_values:
41
               X[row + 1, column + 1] = columns_means[column]
42
43
           return X
44
```

При небольшом шуме (noise_level = 0.1) PCA сохраняет высокую точность вплоть до тысячных. Также замечена отрицательная корреляция между количством компонент и погрешностью: чем больше количество компонент, тем выше устойчивость у шуму. Это может говорить о том, что шум как-бы расползается по осям и большее их количество компенсирует его.

Применение РСА к реальному датасету

Реализованный алгоритм PCA был применён к датасету "pokemon.csv". Время выполнения значительно возрасло с большим размером матрицы.

```
ds = file.read()
file.close()

X = PCA.handle_missing_values(ds)

projection, components, variance = PCA.RSA(X, k)
mean_vector = PCA.mean_vector(X)

reconstructed_X = PCA.reconstruction(projection, components, mean_vector)
error = PCA.reconstruction_error(X, reconstructed_X)

return projection, error
```

Тестирование

Для тестирования работы PCA были сгенерированы матрицы. Эти матрицы спроецированы на компоненты и реконструированы обратно. Далее проверялось лежит ли средняя квадратичная ошибка в пределах погрещности.

Тесты можно найти в репозитроии в разделе tests.

Заключение

В ходе лабороторной работы был реалезован полный цикл РСА, математически обоснован выбор собственных векторов матрицы ковариаций, проведено тестирование программы, иследовано влияние шумов на работу РСА.

Из недочётов реализованного РСА, можно выделить:

- Низкая скорость работы на больших матрицах. Так как реализованный РСА использовал класс разряженных матриц многие опреации такие как вызов элемента по индексу, поиск определителя и т.д. более сложные по сравнению с плотными.
 - Для улучшения скорости можно: использовать плотные матрицы пожертвовав памятью, переписать сложные блоки кода на C++ или C.
- Погрешность при вычислениях. Хоть и погрешность не велика, она может негативно влиять на желаемый результат.
 - Для улучшения погрешности можно использовать библиотеку Python Decimal, или реализовать улучшение дробных чисел самостоятельно, но это негативно скажется на алгоритмической и пространственной сложности алгоритма.

В процессе выполнения лабороторной работы были получены и закрплены знания теории линейных операторов, получен опыт работы с применением теории линейной алгебры на практике.