

Доказательство: оптимальные направления PCA совпадают с собственными векторами матрицы ковариаций

1. Постановка задачи

Цель метода главных компонент (PCA) — найти такие ортонормированные направления $\mathbf{v}_1, \dots, \mathbf{v}_k$, на которые проекция данных $\mathbf{X} \in \mathbb{R}^{n \times m}$ (данные центрированы, то есть среднее по каждому признаку равно нулю) даёт максимальную дисперсию. Это означает, что проекция данных на эти направления сохраняет наибольшее количество информации.

2. Дисперсия проекции на направление

Пусть $\mathbf{v} \in \mathbb{R}^m$ — вектор направления, такой что $\|\mathbf{v}\| = 1$. Тогда проекция наблюдения \mathbf{x}_i на \mathbf{v} — это скаляр $\mathbf{v}^T \mathbf{x}_i$. Дисперсия таких проекций по всем наблюдениям равна:

$$\text{Var}(\mathbf{v}^T \mathbf{x}) = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{v}^T \mathbf{x}_i)^2 = \mathbf{v}^T \Sigma \mathbf{v},$$

где $\Sigma = \frac{1}{n-1} \mathbf{X}^T \mathbf{X}$ — ковариационная матрица.

3. Задача максимизации

Цель — найти такой вектор \mathbf{v} , при котором выражение $\mathbf{v}^T \Sigma \mathbf{v}$ достигает максимального значения при условии $\mathbf{v}^T \mathbf{v} = 1$.

Это задача на экстремум функции при ограничении. Решу её с помощью метода множителей Лагранжа.

4. Метод Лагранжа

Рассмотрю вспомогательную функцию:

$$\mathcal{L}(\mathbf{v}, \lambda) = \mathbf{v}^T \Sigma \mathbf{v} - \lambda(\mathbf{v}^T \mathbf{v} - 1),$$

где λ — множитель Лагранжа. Чтобы найти стационарные точки, приравняю производные этой функции по компонентам \mathbf{v} к нулю:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = 2\Sigma \mathbf{v} - 2\lambda \mathbf{v} = 0.$$

Получаю условие:

$$\Sigma \mathbf{v} = \lambda \mathbf{v}.$$

Это уравнение означает, что вектор \mathbf{v} является собственным вектором матрицы Σ , а λ — соответствующим собственным значением.

5. Свойства матрицы ковариаций

Матрица Σ симметрична ($\Sigma = \Sigma^T$) и положительно полуопределённая, то есть $\mathbf{v}^T \Sigma \mathbf{v} \geq 0$ для любого вектора \mathbf{v} . Следовательно:

- Все собственные значения λ_i вещественные и неотрицательные.
- Собственные векторы, соответствующие различным собственным значениям, ортогональны.

6. Максимизация дисперсии

Так как выражение $\mathbf{v}^T \Sigma \mathbf{v}$ достигает максимума при собственном векторе, соответствующем наибольшему собственному значению λ_1 , этот вектор и есть направление первой главной компоненты.

Вторую компоненту выбираем аналогично, среди оставшихся направлений, при этом она должна быть ортогональна первой. Таким образом, каждая следующая компонента — это собственный вектор, соответствующий следующему по убыванию собственному значению, и ортогональный предыдущим.

7. Базис из собственных векторов

Собственные векторы матрицы Σ можно выбрать ортонормированными, так как она симметрична. Это означает, что они образуют ортонормированный базис, что соответствует требованиям PCA.

Полная сумма дисперсий всех направлений равна следу матрицы Σ :

$$\text{Tr}(\Sigma) = \sum_{i=1}^m \lambda_i.$$

Если выбрать только первые k компонент, то сохраняется доля дисперсии:

$$\gamma = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^m \lambda_i}.$$

8. Заключение

Таким образом, оптимальные направления в методе PCA — это собственные векторы матрицы ковариаций Σ , отсортированные по убыванию соответствующих собственных значений. Они:

- максимизируют дисперсию проекций,
- ортогональны между собой,
- образуют базис в пространстве признаков.