

PCA

1. Cara kerja algoritma PCA

Principal Component Analysis (PCA) adalah teknik statistik yang digunakan untuk mengurangi dimensi data dengan tetap mempertahankan sebanyak mungkin variabilitas yang ada di data asli. PCA melakukan transformasi data ke dalam set baru dari variabel yang disebut komponen utama (principal components), yang merupakan kombinasi linear dari variabel asli. Komponen utama diurutkan berdasarkan jumlah variabilitas yang mereka jelaskan dalam data. Langkah-langkah utama dalam algoritma PCA adalah sebagai berikut.

Standarisasi Data:

Langkah pertama dalam PCA adalah menstandarisasi data sehingga semua fitur memiliki rata-rata (mean) 0 dan variansi (variance) 1. Hal ini penting karena PCA sangat bergantung pada skala data, dan tanpa standarisasi, fitur dengan variabilitas lebih besar akan mendominasi komponen utama.

Membentuk Matriks Kovariansi:

Setelah data distandarisasi, langkah berikutnya adalah menghitung matriks kovariansi dari data. Matriks kovariansi ini menunjukkan seberapa banyak variabel berhubungan satu sama lain. Matriks ini membantu dalam menentukan arah di mana data memiliki varians terbesar.

Dekomposisi Eigen:

Dengan menggunakan matriks kovariansi, PCA selanjutnya melakukan dekomposisi eigen untuk menemukan eigenvalues (nilai eigen) dan eigenvectors (vektor eigen). Eigenvectors adalah arah (komponen utama) yang menunjukkan arah terbesar dari varians data, sementara eigenvalues menunjukkan besarnya varians dalam arah tersebut.

Urutkan Komponen Utama:

Komponen utama diurutkan berdasarkan nilai eigen dari terbesar ke terkecil. Komponen utama yang pertama (PC1) akan menjadi arah dengan varians terbesar, diikuti oleh komponen utama kedua (PC2), dan seterusnya. Jumlah komponen utama yang dipilih biasanya ditentukan berdasarkan seberapa banyak varians yang ingin dijelaskan oleh model.

Transformasi Data:

Terakhir, data asli diproyeksikan ke ruang baru yang dibentuk oleh komponen utama yang telah dipilih. Hal ini akan menghasilkan data baru yang berdimensi lebih rendah tetapi tetap mempertahankan sebagian besar informasi dari data asli.

4. Hasil Evaluasi Model

From Scratch	From Scikit-Learn
<pre> Data setelah PCA: [[-5.91600974 0.09629415] [2.13804043 -1.86573936] [3.8931435 2.38558965] ... [3.17743598 -0.55307688] [2.16683244 0.19072479] [5.26946903 -0.57551978]] Komponen utama: [[0.02154493 -0.0010277] [-0.00291384 -0.00205451] [0.00797286 0.00346774] [0.02498181 0.02011766] [-0.01022103 -0.00339059] [0.99331706 0.02196462] [0.06550577 -0.08504673] [0.04377031 0.02897547] [-0.01793922 -0.01093757] [0.01963611 -0.99230546] [-0.0065307 0.00476126] [0.00901843 0.01008831] [0.02204657 -0.00142238] [0.00355568 -0.03925532] [0.03072379 0.06385944] [-0.04791266 0.01984478] [0.03577614 0.00908157]] Persentase explained variance dari tiap komponen: [0.58025585 0.10370574]</pre>	<pre> Data setelah PCA: [[-5.91600974 -0.09629415] [2.13804043 1.86573936] [3.8931435 -2.38558965] ... [3.17743598 0.55307688] [2.16683244 -0.19072479] [5.26946903 0.57551978]] Komponen utama (Scikit-learn PCA): [[0.02154493 -0.00291384 0.00797286 0.02498181 -0.01022103 0.99331706 0.06550577 0.04377031 -0.01793922 0.01963611 -0.0065307 0.00901843 0.02204657 0.00355568 0.03072379 -0.04791266 0.03577614] [0.0010277 0.00205451 -0.00346774 -0.02011766 0.00339059 -0.02196462 0.08504673 -0.02897547 0.01093757 0.99230546 -0.00476126 -0.01008831 0.00142238 0.03925532 -0.06385944 -0.01984478 -0.00908157]] Persentase explained variance dari tiap komponen (Scikit-learn PCA): [0.58025585 0.10370574]</pre>

Hasil data setelah PCA dari implementasi scratch dan scikit-learn serupa dalam struktur, tetapi mungkin ada perbedaan kecil dalam nilai numerik karena perbedaan dalam cara implementasi dan pembulatan. Komponen utama (eigenvectors) dari PCA Scratch dan SKPCA umumnya serupa, tetapi ada kemungkinan perbedaan kecil dalam tanda dan skala komponen. Hal ini dapat disebabkan oleh perbedaan dalam urutan eigenvectors (komponen utama) dan cara normalisasi yang diterapkan. Sementara itu, persentase explained variance yang diberikan oleh PCA Scratch dan SKPCA cocok, menunjukkan bahwa kedua metode menghasilkan hasil yang konsisten dalam hal informasi yang dibawa oleh komponen utama.

Perbedaan kecil dalam hasil dapat disebabkan oleh perbedaan dalam precision numerik dan stabilitas algoritma, terutama jika menggunakan metode eigenvalue yang berbeda atau pemrosesan data. Implementasi manual mungkin memiliki perbedaan dalam cara penanganan data, normalisasi, dan cara menghitung matriks kovarians dan eigenvalues/eigenvectors dibandingkan dengan implementasi yang dioptimalkan seperti di scikit-learn. Sorting dari eigenvalues dan komponen utama serta penandaan (sign) dari komponen utama mungkin berbeda, yang dapat menyebabkan perbedaan dalam nilai komponen utama dan urutan.

5. *Improvement*

Untuk meningkatkan performa model dapat dilakukan optimasi hyperparameter dengan menggunakan grid search atau random search untuk menemukan hyperparameter terbaik, termasuk jumlah komponen PCA jika digunakan dalam model pembelajaran mesin. Kombinasi dengan teknik lain juga perlu dipertimbangkan. Jika data tidak linier, pertimbangkan untuk menggunakan kernel PCA yang dapat menangkap struktur non-linier dalam data. Optimasi juga dapat dilakukan pada implementasi PCA dengan menggunakan Singular Value Decomposition (SVD) yang seringkali lebih stabil secara numerik dibandingkan metode eigenvalue. Sementara untuk stabilitas numerik dapat menggunakan algoritma yang lebih stabil secara numerik untuk menghitung matriks kovarians dan eigenvalues. Terakhir, dapat juga menguji model dengan berbagai dataset untuk memastikan bahwa model yang dihasilkan dapat bekerja dengan baik pada berbagai data.