

Preguntas Conceptuales.

1º ¿Que pasa si no estandarizamos características con escalas muy diferentes?

Si no se estandarizan, las componentes principales estarán controladas por variables con mayor varianza a nivel de escala, por lo que varias variables relevantes en escala pequeña pueden quedar ignoradas.

2º ¿cómo determinar el número óptimo de componentes?

Elegir el mínimo número de componentes que expliquen el umbral objetivo, también conservar con un valor > 1 .

3º si PCA reduce dimensionalidad conservando la mayor parte de la varianza, se puede usar para compresión de reconstrucción aproximada.

Con mis: $n = 150$ muestras, $P = 4$ características

Original: almacena $(150 \cdot 4) = 600$ números

con 2 componentes sería $n \cdot K = 2 \cdot 150 = 300$ números.

con 2 componentes principales sería $P \cdot K = 2 \cdot 4 = 8$

Media $P = 4$

total sería $= 300 + 8 + 4 = 312$ números.

razón de compresión $= \frac{600}{312} = 1.9231$

Se reduce un 1.9231 del tamaño almacenado.

4. Relación entre PCA y SVD.

Sea X la matriz de datos centrada ($n \times p$). La SVD de X es

$$X = U \Sigma V^T$$

$U (n \times r)$, $\Sigma (r \times r)$ diagonal con valores singulares σ_i y

$V (p \times r)$ vectores singulares pequeños.

Calcular los SVD de la matriz X y usar V y los σ_i para obtener direcciones y varianzas.

