

	VIETTEL AI RACE	Public 118
	DIMENSIONALITY REDUCTION & PCA	Lần ban hành: 1

1. Giới thiệu

Dimensionality Reduction (giảm chiều dữ liệu) là một kỹ thuật quan trọng trong Machine Learning. Dữ liệu thực tế có thể có số chiều rất lớn (hàng nghìn). Việc giảm chiều giúp tiết kiệm lưu trữ, tăng tốc tính toán và có thể coi như nén dữ liệu. Một phương pháp tuyến tính cơ bản là Principal Component Analysis (PCA).

2. Một chút toán

2.1 Norm 2 của ma trận

$$\|A\|_2 = \max_x \|Ax\|_2 / \|x\|_2 \quad (1)$$

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 \quad (2)$$

Giải bằng nhân tử Lagrange cho thấy norm 2 của ma trận chính là singular value lớn nhất của A.

Vector tương ứng là right-singular vector của A.

2.2 Biểu diễn vector trong các hệ cơ sở khác nhau

$$x = Uy, \quad y = U^{-1}x \quad (7)$$

Nếu U trực giao: $U^{-1} = U^T$, do đó $y = U^T x$.

2.3 Trace

Một số tính chất:

$$\text{trace}(A) = \text{trace}(A^T)$$

$$\text{trace}(kA) = k \text{trace}(A)$$

$$\text{trace}(AB) = \text{trace}(BA)$$

$$\|A\|_F^2 = \text{trace}(A^T A)$$

$$\text{trace}(A) = \text{tổng các trị riêng của } A$$

2.4 Kỳ vọng và ma trận hiệp phương sai

$$\text{Một chiều: } \bar{x} = (1/N) \sum x_n, \quad \sigma^2 = (1/N) \sum (x_n - \bar{x})^2$$

$$\text{Đa chiều: } \bar{x} = (1/N) \sum x_n, \quad S = (1/N) (X - \bar{x}1^T)(X - \bar{x}1^T)^T$$

3. Principal Component Analysis (PCA)

Mục tiêu: Tìm hệ cơ sở trực chuẩn sao cho phương sai dữ liệu tập trung ở K thành phần đầu.

$$\text{Dữ liệu chuẩn hoá: } \dot{X} = X - \bar{x}1^T$$

$$\text{Ma trận hiệp phương sai: } S = (1/N) \dot{X} \dot{X}^T$$

	VIETTEL AI RACE	Public 118
	DIMENSIONALITY REDUCTION & PCA	Lần ban hành: 1

Hàm mất mát: $J = \sum_{i=K+1}^D u_i^T S u_i$

Tối ưu tương đương chọn K vector riêng ứng với K trị riêng lớn nhất của S .

4. Các bước PCA

- Tính kỳ vọng \bar{x}
- ▣ Chuẩn hoá dữ liệu: $\tilde{X} = X - \bar{x}1^T$
- Tính ma trận hiệp phương sai S
- Tính trị riêng & vector riêng, sắp xếp λ giảm dần
- Chọn K vector riêng lớn nhất $\rightarrow U_K$
- Tính tọa độ mới: $Z = U_K^T \tilde{X}$
- Xấp xỉ khôi phục: $x \approx U_K z + \bar{x}$