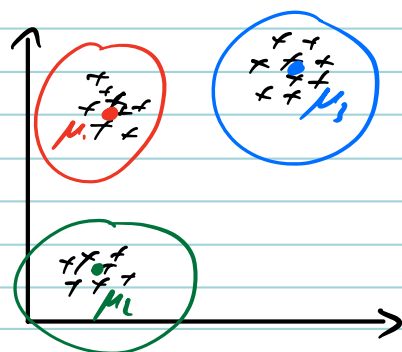


Dimensionality Reduction

Vector Quantisation:

If your data is clustered, you can approximate each input by its cluster assignment. E.g. GMMs give you a probability γ_i that \vec{x}_i is in cluster l .



$$\vec{x}_i \rightarrow \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ik} \end{pmatrix} \leftarrow \text{New } k\text{-dimensional representation.}$$

Covariances:

Random Variables $X^A, X^B \sim P(X^A, X^B)$ with $\mu^A = E[X^A] = 0$ $\mu^B = E[X^B] = 0$

Variance: $\text{Var}(X^A) = E[(X^A - \mu^A)^2] = E[X^A^2]$

Covariance: $\text{Cov}(X^A, X^B) = E[(X^A - \mu^A)(X^B - \mu^B)] = E[X^A X^B]$ $\text{COV}(X^A, X^A) = \text{VAR}(X^A)$

$$E[X^A X^B] = \begin{cases} > 0 & \text{positively correlated: } i) X^A \text{ is } > 0, X^B \text{ is } > 0 \text{ (and vice versa)} \\ = 0 & \text{uncorrelated} \\ < 0 & \text{negatively correlated: } i) X^A \text{ is } > 0, X^B \text{ is } < 0 \text{ (and vice versa)} \end{cases}$$

Covariance Matrix: If $\vec{x} \sim P$ is a vector $\vec{x} = [x_1, x_2, x_3, \dots, x_d]^T$

Assume data $D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\} \in \mathbb{R}^d$

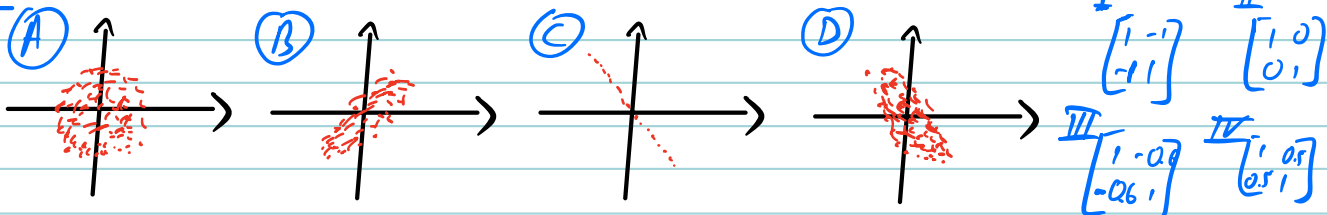
$$\vec{\mu} = E[\vec{x}] \approx \frac{1}{n} \sum_{i=1}^n \vec{x}_i \leftarrow \text{Weak law of large numbers.}$$

$$C = \text{Cov}(\vec{x}) = E[(\vec{x} - \vec{\mu})(\vec{x} - \vec{\mu})^T] = E[\vec{x} \vec{x}^T] \approx \frac{1}{n} \sum_{i=1}^n \vec{x}_i \vec{x}_i^T \leftarrow \text{both are indices } i.$$

\uparrow Covariance Matrix of all r.v. in X , i.e. x_1, x_2, \dots, x_d

$$C_{\alpha\beta} = \text{COV}(x_\alpha, x_\beta) \quad C_{\alpha\alpha} = \text{VAR}(x_\alpha)$$

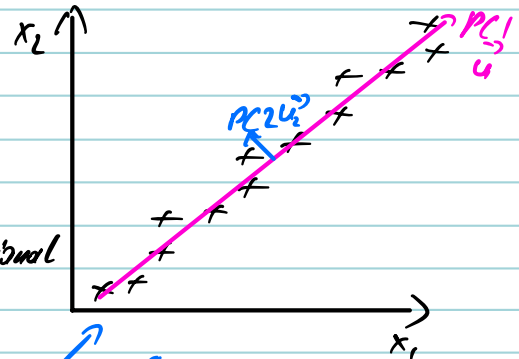
Match:



Principal Component Analysis: (Pearson 1901)

Data $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ but are truly from a lower-dimensional subspace recd.

Idea: Find basis vectors for this subspace and project data onto it. \Rightarrow Leads to r -dimensional representation



Everything "interesting" is along PC1
PC2 is likely only noise.

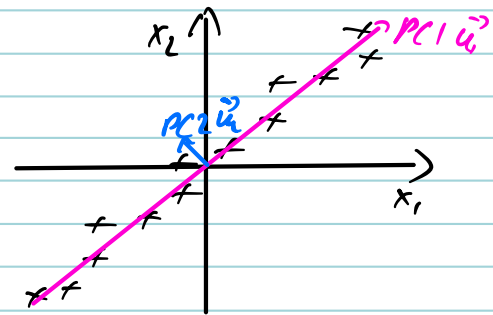
Step #1 of PCA: Center data

$$\vec{\mu} = \frac{1}{n} \sum_i \vec{x}_i \quad \text{subtract mean: } \underline{\underline{\vec{x}_i \leftarrow \vec{x}_i - \vec{\mu}}}}$$

Step #2: Find first principal component

PCA finds the subspace that contains maximum variance.

PC #1: Find \vec{v} s.t. after projection, $\vec{v}^T \vec{x}_i$, variance is maximized.



$$\max_{\vec{u}, \vec{u}^T \vec{u} = 1} \frac{1}{n} \sum_{i=1}^n (\vec{x}_i^T \vec{u})^2 = \max_{\vec{u}, \vec{u}^T \vec{u} = 1} \sum_i (\vec{x}_i^T \vec{u}) (\vec{x}_i^T \vec{u}) = \max_{\vec{u}} \sum_{i=1}^n \vec{u}^T \vec{x}_i \vec{x}_i^T \vec{u} = \max_{\vec{u}, \vec{u}^T \vec{u} = 1} \vec{u}^T \left(\sum_{i=1}^n \vec{x}_i \vec{x}_i^T \right) \vec{u} = \max_{\vec{u}, \vec{u}^T \vec{u} = 1} \vec{u}^T C \vec{u}$$

↑ We only care about the direction.
↑ enforces $\vec{u}^T \vec{u} = 1$ (if not $\vec{u}^T \vec{u} = 1$ will set $\lambda \rightarrow \infty$)
↑ Covariance matrix

Lagrangian:

$$\max_{\vec{u}} \min_{\lambda} \underbrace{\vec{u}^T C \vec{u} - \lambda (\vec{u}^T \vec{u} - 1)}_{\mathcal{L}(\vec{u}, \lambda)} \quad \frac{\partial \mathcal{L}}{\partial \vec{u}} = 0 \quad \leftarrow \text{must hold at optimum.}$$

$$\Rightarrow 2C\vec{u} - 2\lambda\vec{u} = 0$$

$$\therefore \underline{\underline{C\vec{u} = \lambda\vec{u}}}$$

\vec{u} is an eigenvector of C

C has d eigenvectors: u_1, u_2, \dots, u_d s.t. $Cu_i = \lambda_i u_i$

$u_i^T u_j = 1$
 $u_i^T u_j = 0$ if $i \neq j$ \Leftarrow orthogonal unit vectors.

Sort eigenvectors such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

u_1 is the first (aka leading) principal component. $\Rightarrow U = [u_1, \dots, u_d]$
 u_2 is the second, ...

$$\underline{\underline{PCA:}} \quad \vec{x}_i \in \mathbb{R}^d \rightarrow U \underbrace{(\vec{x}_i - \vec{\mu})}_{\in \mathbb{R}^d}$$

Reconstruction:

PCA dimensionality reduction: $\vec{z}_i = U^T(\vec{x}_i - \vec{\mu})$

PCA reconstruction: $\hat{x}_i = U\vec{z}_i + \vec{\mu}$

Quiz: proof that if $r=d$ the reconstruction is perfect (i.e. $\hat{x}_i = x_i$).

PCA de-correlates dimensions:

Correlation matrix of z_1, \dots, z_n :

$$C_z = \frac{1}{n} \sum_{i=1}^n \vec{z}_i \vec{z}_i^T = \frac{1}{n} \sum_i U^T \vec{x}_i (U^T \vec{x}_i)^T = \frac{1}{n} \sum_{i=1}^n U^T \vec{x}_i \vec{x}_i^T U = \frac{1}{n} U^T \left(\sum_{i=1}^n \vec{x}_i \vec{x}_i^T \right) U$$
$$= \frac{1}{n} U^T C U \Rightarrow [C_z]_{ij} = u_i^T C u_j \begin{cases} \lambda_i & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

↑ ↑
eigenvectors

$$C_z = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} \Rightarrow C_z = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$$

How to pick r ? λ_i is the variance within r^{th} PC.

If you project onto r dimensions you

lose $\left(\frac{\sum_{i=1}^r \lambda_i}{\sum_{i=1}^n \lambda_i} \right)$ fraction of the total variance.

Denoising: Pick smallest r such that $\frac{\sum_{i=1}^r \lambda_i}{\sum_{i=1}^n \lambda_i} \geq 0.95$

Project out
5% variance
as noise.

Singular Value Decomposition:

(for centered data)

$$X = U S V^T$$

↑ ↑ ←
Principal components eigenvalues of C $S V^T = U^T X$ ← projected data