

Exploratory data analysis: PCA, clustering

- you have already collected data from various methods we have discussed.
 - 100 material \rightarrow [tensile strength, melting temp, bulk modulus, ...]
 - given material i , $x_i = [T, T_m, B, \dots]$

$$X = N \text{ rows} \times d \text{ columns}$$

↑
"dimension"

- distypically high thus we need a way to visualize our summary statistics.
- mean array (vector) of our samples: $1 \times d$ vector
 covariance: $d \times d$

Dimensionality reduction:

- reduce $X_{N \times d} \rightarrow Y_{N \times q}$ such that q is lower than d



- Principal component analysis (PCA)

$$f(\lambda) = \mu + V\lambda$$

$$= (1 \times d) + (d \times q) \times (q \times 1)$$

linear model parametrized by μ and V

eg: $q=2$ $[\lambda_0, \lambda_1] = \lambda$
 \downarrow
 $f(\lambda) = (1 \times d)$ vector

$$\text{error} = \sum_{i=1}^N \|x_i - f(\lambda_i)\|^2$$

$$= \sum_{i=1}^N \|x_i - \mu - V\lambda_i\|^2$$

one solution that minimizes the above

$$\hat{\mu} = \bar{x}_i$$

$$\hat{\lambda}_i = V^T(x_i - \mu)$$

$$= \sum_{i=1}^N \| \underbrace{x_i - \hat{\mu}} - V \underbrace{V^T(x_i - \hat{\mu})} \|^2$$

$$= \sum_{i=1}^N \| \tilde{x} - K \tilde{x} \|^2$$

we are minimizing reconstruction error when projecting our mean centered samples by a matrix 'K'

- the projection by K maximizes the variance along each dimensions.

$$X_{N \times d} \xrightarrow{K} Y_{N \times q}$$

$$C_\lambda = \frac{1}{N} \Lambda \Lambda^T \begin{cases} \Lambda = V^T(x - \mu) \\ \uparrow \\ \lambda_i = V^T(x_i - \mu) \end{cases}$$

$$= \frac{1}{N} V^T \tilde{x} \tilde{x}^T V$$

$$= \frac{1}{N} V^T C_{\tilde{x}} V$$

- spectral theorem $C_{\tilde{x}} = E^T D E$ \rightarrow diag(eigenvals)

$$C_\lambda = V^T E^T D E V$$

$$= (E V)^T D (E V)$$

our goal was to have 0 offdiagonals for C_λ

$$V = E^{-1} \Rightarrow C_\lambda = D \text{ diagonal matrix}$$

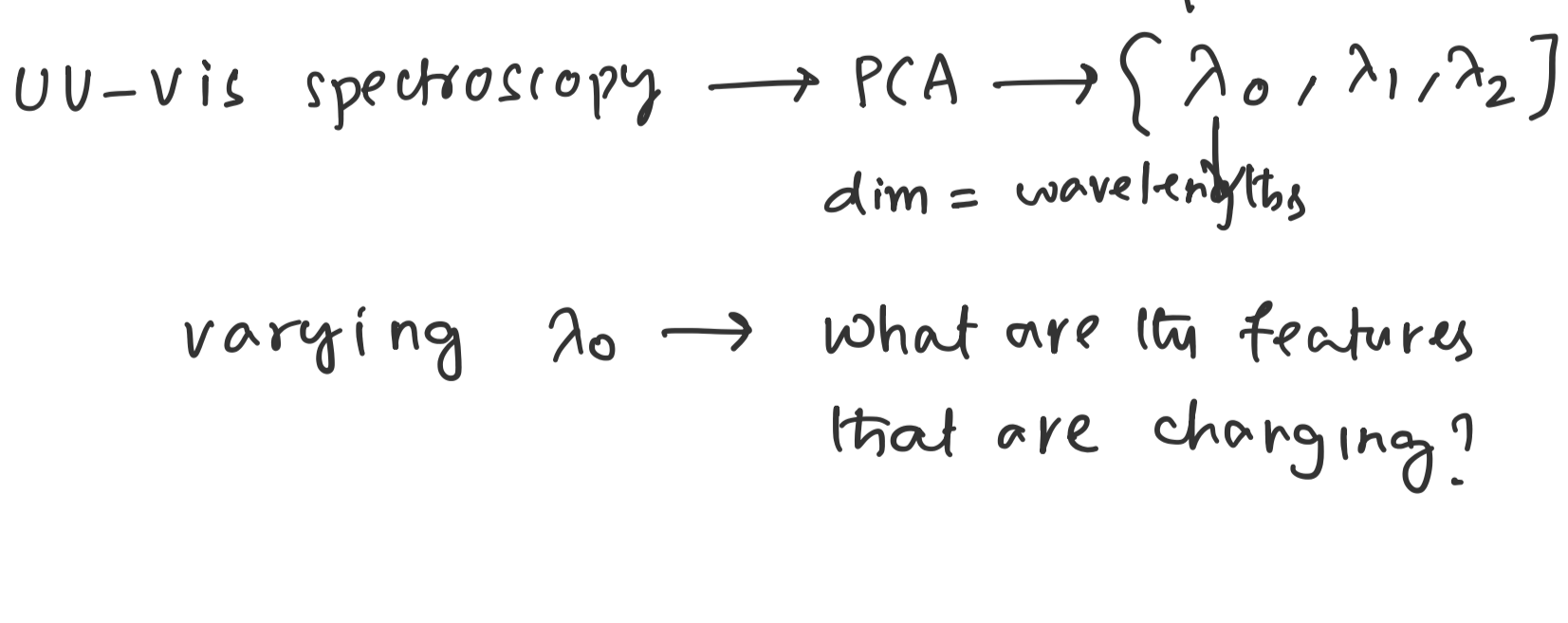
$$\hat{\mu} = \tilde{x}$$

$$\lambda_i = E(x_i - \hat{\mu})$$

(E is an orthogonal matrix, $E^{-1} = E^T$, $V^T = E$)

- each dimension in the projected space has a variance given by eigenvalues in D which can be used to quantify explained variance
 \downarrow
 find right set of dimension given a variance threshold

- generative model: (Asg 2)



Clustering:

- group samples based on similarities (distance metric, similarity function)
- within group similarities to be maximum
 across group " minimum
- k-mean clustering: 1) assigning: grouping based on template
 $\left(\begin{array}{l} \text{convergent} \\ \downarrow \\ \text{2) update: mean value of grouped points as template} \end{array} \right)$
 (when your templates does not change you stop iterating return the group labels)
- dim red: isomap, diffusion maps
 clustering: spectral clustering, density based clustering.