

- Unsupervised learning

- PCA

Unsupervised learning: Sample points but no labels

↳ no classes, no Y-values, nothing to predict

↳ goal: discover structure of data

Examples of Unsupervised learning:

- **clustering:** partition data into groups of similar/nearby pts.
- **dimensionality reduction:** data often lies on an underlying lower dimension
 - matrices have low rank approximations
- clustering was about finding similarity, dim. reduction more about finding variation
- **density estimation:** fit continuous distribution to discrete data
 - e.g. MLE to fit gaussians
 - can do more complicated functions

Principal Component Analysis: given pts in \mathbb{R}^q , find k directions that capture most of the variance (dimensionality reduction)

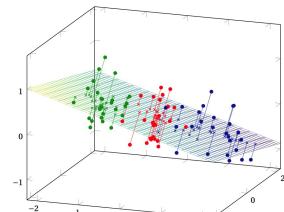


Figure 20.1: An example of PCA applied in a 3D dataset, with two principal components

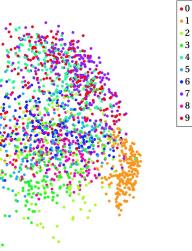
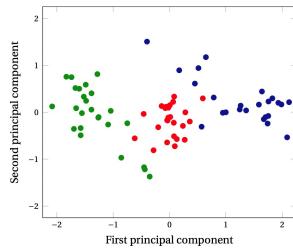


Figure 20.2: The MNIST dataset, projected into 2D (from 784D)

Why PCA?

- dim. reduction \rightarrow faster/cheaper compute
- removes irrelevant features/dimensions \rightarrow less likely to overfit
- PCA is like Subset Selection except the features aren't axis aligned
 - instead, they're lin. combos of input features
- PCA finds small basis for representing variations in complex things

↳ 2 dimensions not enough to separate data
(except 1 vs. 0)

} preprocessing

Computing the PCA:

$X \in \mathbb{R}^{n \times d}$ is the design matrix (no fictitious dim.)
assume centered X , i.e. $\mathbb{1}^T X = 0$

Picking just 1 Principal Component:

Consider unit vector \vec{w} : orthogonal projection of \vec{x} on \vec{w} : $\vec{\hat{x}} = (\vec{x}^T \vec{w}) \vec{w}$

if \vec{w} doesn't have unit length: $\vec{\hat{x}} = \frac{(\vec{x}^T \vec{w})}{\|\vec{w}\|^2} \vec{w}$

goal: pick best direction \vec{w} and project all the data onto \vec{w} to analyze data in 1-D
- can generalize to k-directions

The k vectors form our subspace/basis

if $\vec{v}_1, \dots, \vec{v}_k$ orthogonal, then $\vec{\hat{x}} = \sum_{i=1}^k (\vec{x}^T \vec{v}_i) \vec{v}_i$

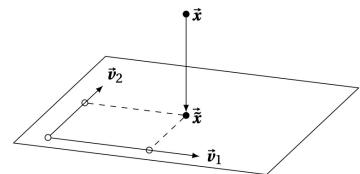


Figure 20.3: Example of projection onto a 2D subspace

Sometimes, we only want the principal coordinates $\vec{x}^T \vec{v}_i$ (don't actually project them)

Consider $X^T X$, which is square symmetric and PSD $\in \mathbb{R}^{d \times d}$

by spectral thm, its eigenvalues are real $\rightarrow 0 \leq \lambda_1 \leq \dots \leq \lambda_d$

Let $\vec{v}_1, \dots, \vec{v}_d$ be the corresponding orthogonal unit eigenvectors (principal components)

greater eigenvalue \rightarrow more important component/direction

Fitting Gaussians Via MLE:

Let's fit a Gaussian to data w/ MLE

↳ choose k Gaussian axes of greatest variance

Recall from MLE: estimate $\hat{\Sigma} = \frac{1}{n} X^T X$, X centered

→ PCA algo:

- Center X
- Normalize X (optional)
 - If units of measurement are different, then normalize (or else bad for Principal components)
 - Otherwise, usually don't normalize
- Want to maintain different variances (assuming same unit of measurement)
- Compute unit eigenvectors & eigenvalues of $X^T X$
- Optionally, choose k based on eigenvalue sizes
- Choose k eigenvectors corresponding to k largest eigenvalues
- Compute principal coordinates $\vec{x}_i^T \vec{v}_i$ of training or test data
 - 2 choices
 - 1) Un-center input training data before projecting
 - 2) translate test data by same vector we used to translate training data when we centered it

Maximizing Sample Variance

- 2nd derivation of PCA

goal: find \vec{w} that maximizes sample variance of projected data

Rayleigh Quotient
↓

$$\text{formally: } \vec{w}^* = \underset{\vec{w}}{\operatorname{argmax}} \left(\operatorname{Var}(\vec{x}_1, \dots, \vec{x}_n) \right) = \frac{1}{n} \sum_{i=1}^n (x_i \cdot \vec{w})^2 = \frac{1}{n} \frac{\|\vec{x}\vec{w}\|^2}{\|\vec{w}\|^2} = \frac{1}{n} \frac{\vec{w}^T \vec{x}^T \vec{x} \vec{w}}{\vec{w}^T \vec{w}}$$

$$\rightarrow \vec{w}^* = \underset{\vec{w}}{\operatorname{argmax}} \left(\frac{1}{n} \cdot \frac{\vec{w}^T \vec{x}^T \vec{x} \vec{w}}{\vec{w}^T \vec{w}} \right) = \frac{\|\vec{x}\vec{w}\|^2}{\|\vec{w}\|^2} = \frac{\|\vec{x}_i \vec{w}\|^2}{\|\vec{w}\|^2} = \lambda_i$$

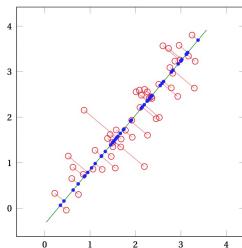


Figure 20.4: Points projected onto a line that maximizes the sample variance of the projected points

→ eigenvector of λ_{\max} maximizes variance $\frac{1}{n} \lambda_{\max}$

One can show that this \vec{v}_{\max} beats the other vectors since each candidate \vec{v} is a lin. combo of eigenvectors

↳ Rayleigh quotient is convex combo. of eigenvalues

k directions → pick k largest (λ, \vec{v}) pairs that are mutually orthogonal

minimizing mean Squared projection distance

3rd derivation: find \vec{w} that minimizes the mean Squared proj. distance

- analogous to LS regression but

- w/ one subtle difference: error isn't in fixed vertical direction
but orthogonal direction instead

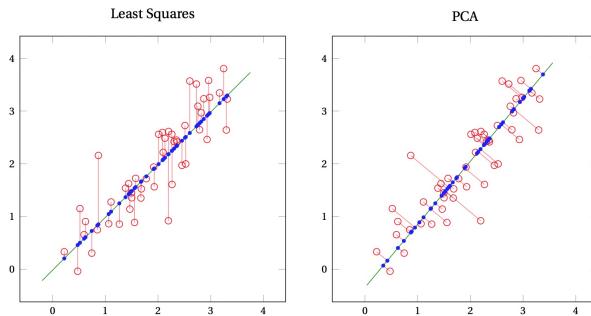


Figure 20.5: Comparison between least squares and PCA

$$\begin{aligned} \rightarrow \text{find } \vec{u} \text{ that minimizes } & \sum_{i=1}^n \|x_i - \hat{x}_i\|^2 = \sum_{i=1}^n \|x_i - \frac{x_i \cdot \vec{u}}{\|\vec{u}\|} \cdot \vec{u}\|^2 \\ &= \sum_{i=1}^n \|x_i\|^2 - \sum_{i=1}^n \left(x_i \cdot \frac{\vec{u}}{\|\vec{u}\|} \right)^2 \end{aligned}$$

Constant $\sum_{i=1}^n \|x_i\|^2$ (n · variance from 2nd derivation)

\Rightarrow minimizing Mean-Squared projection distance \equiv minimizing Variance

Eigenfaces: Use PCA for face recognition. $X \in \mathbb{R}^{n \times d}$, d pixels, n faces

200×200 image \rightarrow 40000-length vector

goal: find nearest neighbor in \mathbb{R}^d given image, compare w/ all training faces

issue: each query takes $\Theta(nd)$ time

\hookrightarrow use PCA to reduce faces to smaller dim. $\mathbb{R}^{d'}$

\rightarrow now $\Theta(nd') \leq \Theta(nd)$