

Lecture 5, Jan 26, 2024

Principal Component Analysis (PCA)

Definition

Dimensionality Reduction: Given a dataset $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ where $\mathbf{x}^{(i)} \in \mathbb{R}^D$, find a mapping $f: \mathbb{R}^D \mapsto \mathbb{R}^d$ where $d < D$ is a lower dimensional space.

- Dimensionality reduction is a type of unsupervised learning
 - PCA is a dimensionality reduction technique
 - Other techniques can include autoencoders, etc
- Dimensionality reduction can be used for a number of purposes:
 - Saving computational time/memory (helps with the curse of dimensionality)
 - Reduces overfitting
 - Visualize high-dimensional datasets
- We're essentially trying to create a summary of the data
- PCA is one of the only dimensionality reduction techniques with a closed-form solution
- PCA uses a linear model with the form $\mathbf{z} = \mathbf{U}^T(\mathbf{x} - \mathbf{b})$ where $\mathbf{U} \in \mathbb{R}^{D \times d}$ is an orthonormal matrix and $\mathbf{b} \in \mathbb{R}^D$
 - These orthonormal columns form a basis for a subspace \mathcal{S}
 - The projection of \mathbf{x} onto \mathcal{S} is the point $\tilde{\mathbf{x}} \in \mathcal{S}$ closes to \mathbf{x} (this is known as the *reproduction* of \mathbf{x})
 - \mathbf{z} is the *representation* or *code* of \mathbf{x}
- Choose $\mathbf{b} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}$
- Finding a general matrix \mathbf{U} is challenging, so we will start with a single column vector \mathbf{u}
 - We aim to minimize the reconstruction error: $\mathcal{L}(\mathbf{u}, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - (\mathbf{u}\mathbf{u}^T(\mathbf{x}^{(i)} - \mathbf{b}) + \mathbf{b})\|_2^2$
 - * $\hat{\mathbf{x}}^{(i)} = \mathbf{u}\mathbf{z} + \mathbf{b} = \mathbf{u}\mathbf{u}^T(\mathbf{x} - \mathbf{b}) + \mathbf{b}$
 - If the data is centered then $\mathbf{b} = \mathbf{0}$, so $\mathcal{L}(\mathbf{u}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)}\|_2^2$
- Expanding the reconstruction error:
 - $\mathcal{L}(\mathbf{u}) = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)})^T (\mathbf{x}^{(i)} - \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)})$
 - $= \frac{1}{N} \sum_{i=1}^n -2\mathbf{x}^{(i)T} \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)} + \mathbf{x}^{(i)T} \mathbf{u}\mathbf{u}^T \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)} + \text{const}$
 - $= \frac{1}{N} \sum_{i=1}^N -\mathbf{x}^{(i)T} \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)} + \text{const}$
 - So we can formulate the problem as minimizing $-\frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)T} \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)}$ subject to $\mathbf{u}^T \mathbf{u} = 1$
 - Equivalently, maximize $\frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)T} \mathbf{u}\mathbf{u}^T \mathbf{x}^{(i)} = \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}^{(i)}\|_2^2$ subject to $\mathbf{u}^T \mathbf{u} = 1$
- Note the mean of \mathbf{z} is zero since we centered \mathbf{x} so the objective function is equivalent to $\frac{1}{N} \sum_{i=1}^N \|\mathbf{z}^{(i)} - \bar{\mathbf{z}}\|_2^2$
 - Minimizing the reconstruction error is equivalent to maximizing the variance of the code vectors

- $$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \|\mathbf{z}^{(i)}\|_2^2 &= \frac{1}{N} \sum_{i=1}^N \mathbf{u}^T (\mathbf{x}^{(i)} - \boldsymbol{\mu})(\mathbf{x}^{(i)} - \boldsymbol{\mu})^T \mathbf{u} \\ &= \mathbf{u}^T \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \boldsymbol{\mu})(\mathbf{x}^{(i)} - \boldsymbol{\mu})^T \right] \mathbf{u} \\ &= \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} \\ &= \mathbf{u}^T \mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^T \mathbf{u} \\ &= \mathbf{a}^T \boldsymbol{\Lambda} \mathbf{a} \\ &= \sum_{j=1}^D \lambda_j \mathbf{a}_j^2 \end{aligned}$$

- We can decompose $\boldsymbol{\Sigma}$ since it is symmetric positive definite, as it is the empirical covariance matrix
 - $\mathbf{a} = \mathbf{Q}^T \mathbf{u}$ is a change of basis to the eigenbasis of $\boldsymbol{\Sigma}$
- Assuming all λ_i are sorted and distinct, we can choose $a_1 = \pm 1$ and $a_j = 0$ (since the first eigenvalue is the largest eigenvalue) in order to maximize the objective
 - Therefore $\mathbf{u} = \mathbf{Q} \mathbf{a} = \mathbf{q}_1$ which is just the top eigenvector
 - More generally, we can show that the k th principal component is given by the k th eigenvector of $\boldsymbol{\Sigma}$ (Courant-Fischer Theorem)
- Alternative derivation: we want to maximize a
 - The Lagrangian is $\mathcal{L}(\mathbf{u}, \gamma) = \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} + \gamma(1 - \mathbf{u}^T \mathbf{u})$
 - $\vec{\nabla}_{\mathbf{u}} \mathcal{L} = (\boldsymbol{\Sigma} + \boldsymbol{\Sigma}^T) \mathbf{u} - 2\gamma \mathbf{1} \mathbf{u} = 0 \implies 2\boldsymbol{\Sigma} \mathbf{u} = 2\gamma \mathbf{u} \implies \boldsymbol{\Sigma} \mathbf{u} = \gamma \mathbf{u}$
- We can also perform PCA with SVD:
 - If \mathbf{X} is a data matrix written in centered form, then the covariance matrix is $\boldsymbol{\Sigma} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$
 - Using an SVD, we can write $\boldsymbol{\Sigma} = \mathbf{V} \mathbf{S}_1 \mathbf{U}_1^T \mathbf{U}_1 \mathbf{S} \mathbf{V}^T = \frac{1}{N} \mathbf{V} \mathbf{S}_1^2 \mathbf{V}^T$
 - Since this is equal to $\mathbf{Q} \boldsymbol{\Lambda} \mathbf{Q}^T$ and spectral decompositions are unique, we must have that the columns of \mathbf{V} are the principal components and $\frac{\mathbf{S}_1^2}{N} = \boldsymbol{\Lambda}$
 - So to construct the PCA we can just take the first d columns
 - Using SVD is faster and more stable
- Note the code vectors given by PCA are de-correlated (i.e. their covariance matrix is diagonal)