

Recitation Week 9

Ashwin Bhola

CDS, NYU

Oct 30th, 2019

1. Suppose $x_1, \dots, x_n \in \mathbb{R}^d$ are datapoints you want to represent in $k < d$ dimensions.
 1. Explain how to do this using PCA
 2. How can you implement PCA using SVD?
 3. How to determine an optimal value for k ?

1. PCA steps:

- i. Center your data: If $X = \begin{bmatrix} \cdots & x_1 & \cdots \\ & \vdots & \\ \cdots & x_n & \cdots \end{bmatrix} \in \mathbb{R}^{n \times d}$ is your data matrix, then $A[i, j] = X[i, j] - \frac{1}{n} \sum_p X[p, j]$
- ii. Use the new matrix A of centered data instances to construct the covariance matrix $S = A^T A \in \mathbb{R}^{d \times d}$
- iii. Spectral decomposition of S: $S = V D V^T$
- iv. Choose the top k eigenvalues $(\lambda_1, \dots, \lambda_k)$ and the eigenvectors (v_1, \dots, v_k) from V corresponding to those eigenvalues

- v. Construct your new data matrix as $X_{new} = AV_k = \begin{bmatrix} \cdots & a_1 & \cdots \\ & \vdots & \\ \cdots & a_n & \cdots \end{bmatrix} \begin{bmatrix} \vdots & \vdots & \vdots \\ v_1 & \cdots & v_k \\ \vdots & \vdots & \vdots \end{bmatrix} \in \mathbb{R}^{n \times k}$

2. PCA using SVD: Let's do SVD on our centered data matrix A : $A = M\Sigma N^T$
- $$\Rightarrow A^T A = N\Sigma^T \Sigma N^T = N D N^T$$

We see that $D = \Sigma^T \Sigma \Rightarrow \lambda_i = \sigma_i^2$

Thus we can use the first k right singular vectors to perform the projection

3. Optimal value of k can be chosen using
1. Scree plot
 2. Fraction of variance explained by top k eigenvalues

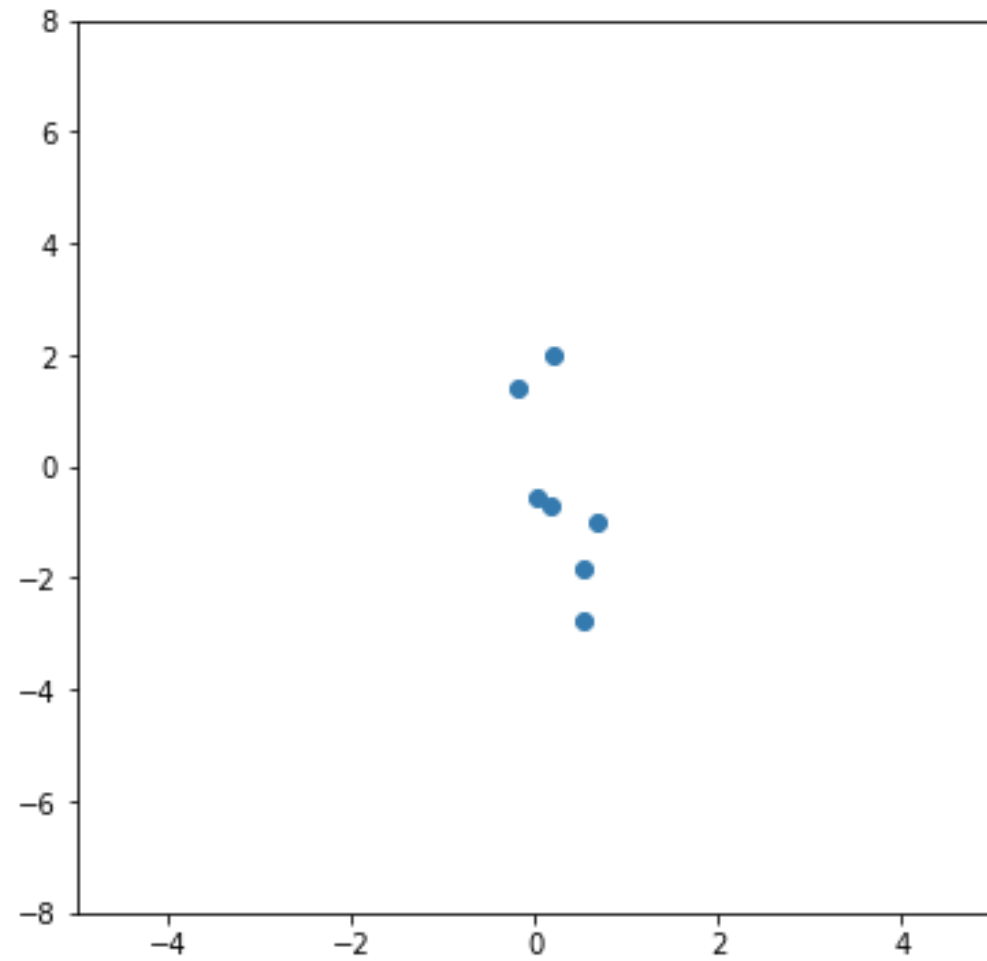
1. Let $X \in \mathbb{R}^{n \times d}$ be your matrix of data points. Suppose you are implementing PCA. Someone suggests that you should standardize your data before calculating the eigenvalues. How do you standardize the data? Is it really required?

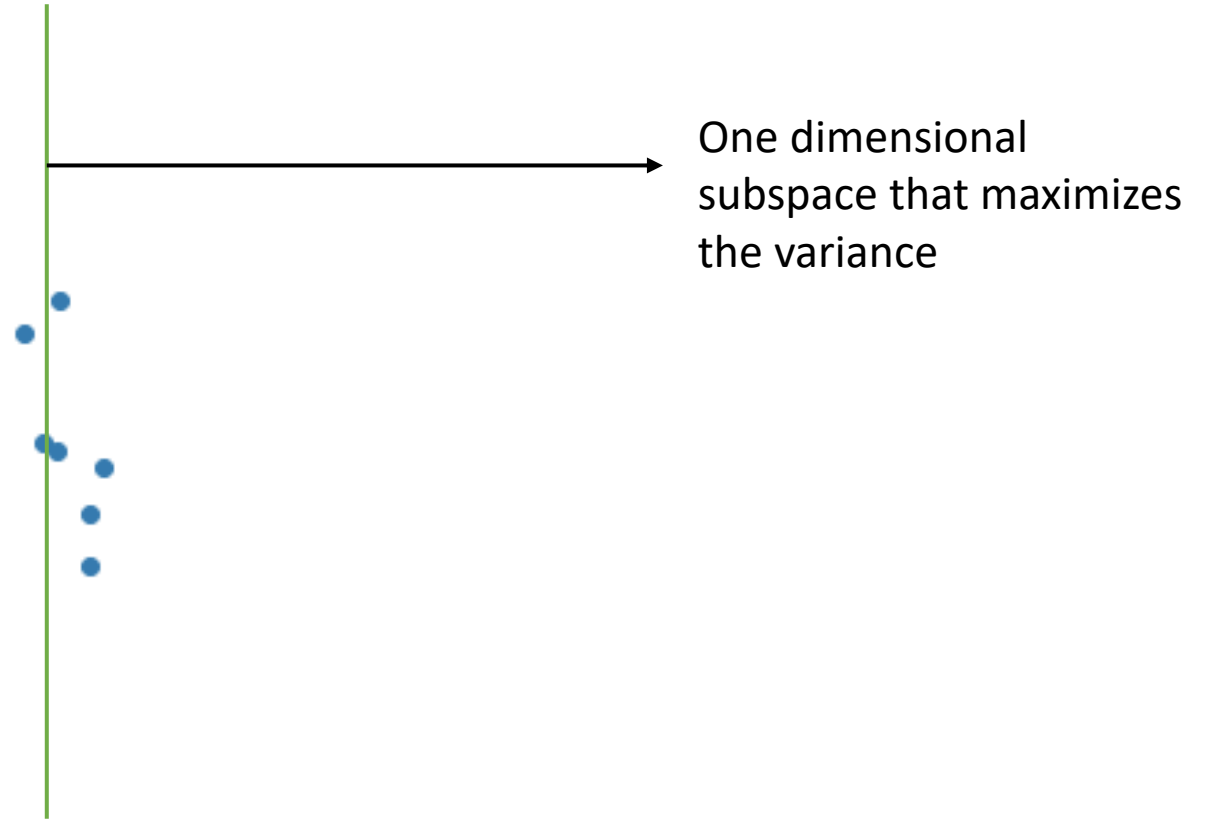
Solution: If $X = \begin{bmatrix} \cdots & x_1 & \cdots \\ & \vdots & \\ \cdots & x_n & \cdots \end{bmatrix} \in \mathbb{R}^{n \times d}$ is your data matrix, then $A[i, j] = \frac{X[i, j] - \mu_j}{\sigma_j}$, where

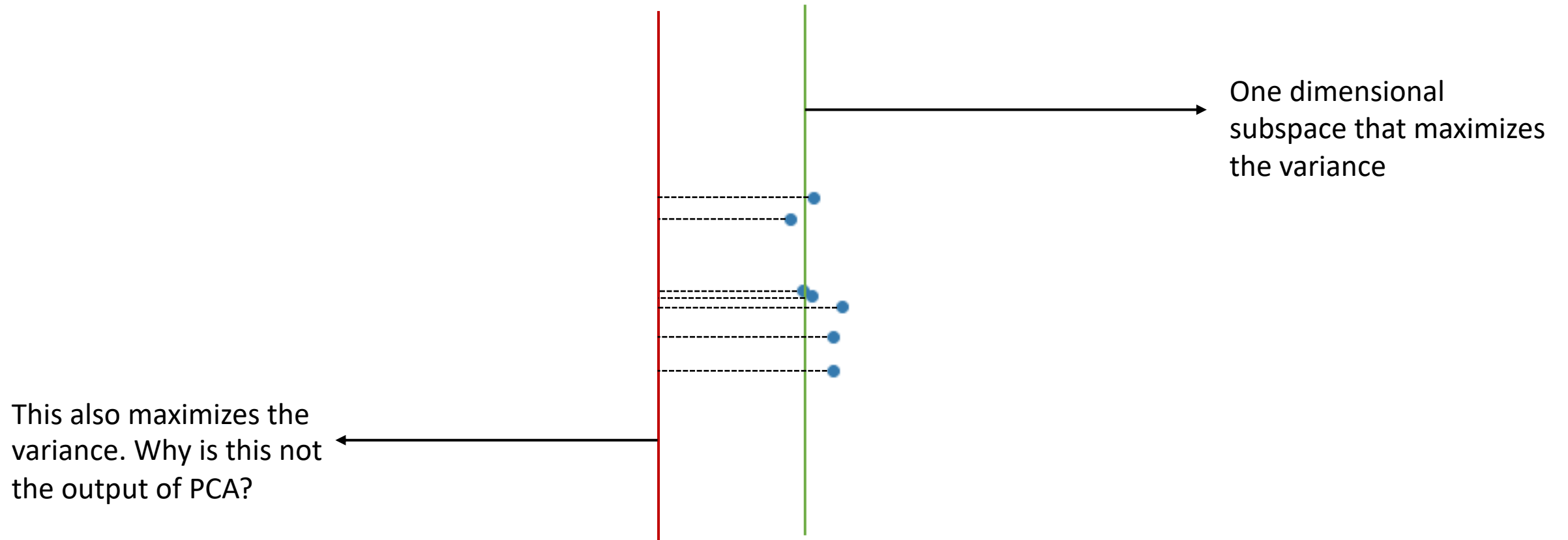
$$\mu_j = \frac{1}{n} \sum_p X[p, j] \text{ and}$$

$$\sigma_j^2 = \sum_p (X[p, j] - \mu_j)^2$$

Centering the data instances is necessary because otherwise, the principal components get influenced by the location of the data in the feature space. This is equivalent to saying that if we don't center the data, then we are trying to approximate the data by a subspace and not an affine space which is less efficient. You may not want to standardize the data if you think the differences in scale are informative. In general, features with larger scale influences the principal components more which may not be desired in some scenarios

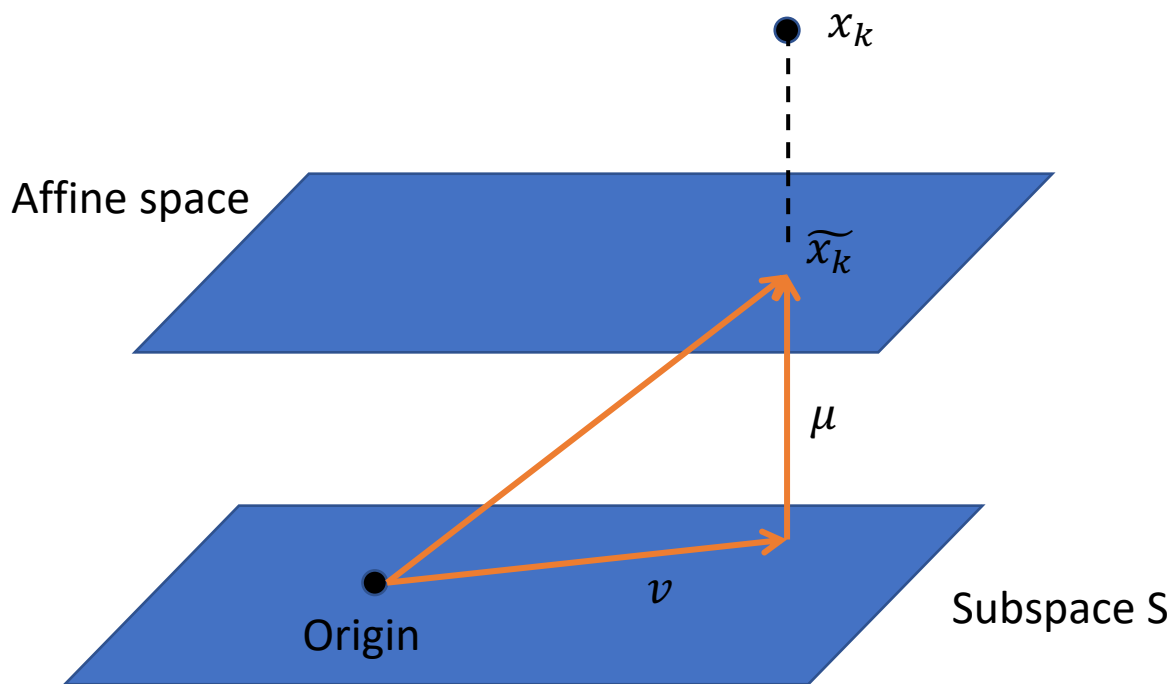






Because we force our principal components to pass through origin which works perfectly only if the data is centered

1. PCA can also be thought of as finding a d -dimensional affine space such that the sum of l_2 distance between the data points and their projection on the affine space is minimum. How to formulate this optimization mathematically?



Let v_1, \dots, v_d be an orthonormal basis of subspace S parallel to the affine space we are trying to find. Assume that μ is the distance between this subspace and the affine space we are trying to find.

$$\text{Objective: } \min \sum_{k=1}^n \|x_k - \tilde{x}_k\|^2$$

Now, $\tilde{x}_k = \mu + v$, where $v \in S$

Let $v = \sum_{i=1}^d (\beta_k)_i v_i = V \beta_k$ (V is a matrix with v_i as the i^{th} column)

Thus, the objective function is: $\min \sum_{k=1}^n \|x_k - \mu - V \beta_k\|^2$

And this minimization is carried to find the affine space. Thus, we minimize the function over $\mu \in \mathbb{R}^p, V \in \mathbb{R}^{p \times d}, \beta'_k \in \mathbb{R}^d$

