

Subspace Learning

*Lecturer: Dr. Raman Arora**Scribe: Corbin Rosset*

1 What is Subspace Learning

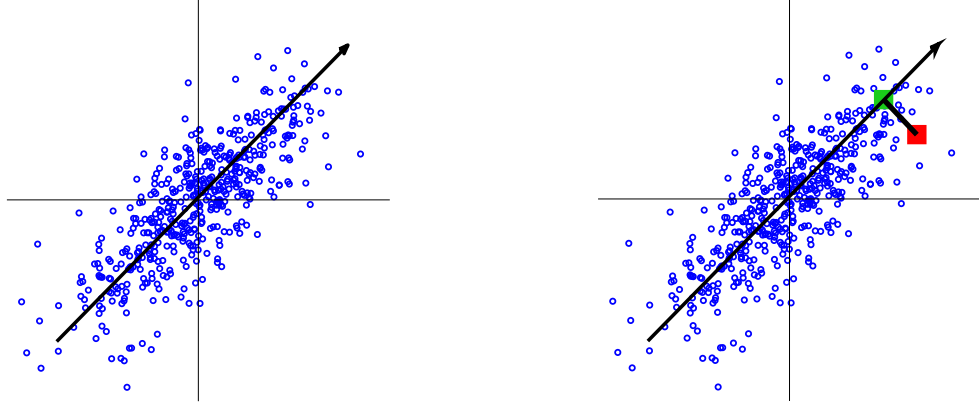
Many raw data sets in machine learning capture high dimensional representations of complex natural phenomena. Yet it is often the case that natural processes are largely controlled by finitely many simpler mechanisms that operate in lower dimensions. For example, a pixel-by-pixel image representation of a data set comprised of rotations of a single object rotated in 3-space under constant light conditions may appear to be high dimensional: shadows cast by surfaces on the face, as well as the different perspectives, make the images appear highly varied, even though the only real parameter in this scenario is the angle of rotation of the object. Instead of representing each image by an array of pixels, each image is merely a linear transformation of a source image. Such transformations of data are common and they only provide a few degrees of freedom with respect to the dimensionality of the dataset. Our task is to seek a low dimensional representation of the data that captures these latent variables.



Figure 1: Seemingly high dimensional variance caused by the rotation of a single object

2 Problem Specification

We assume that the data can be represented in a substantially lower dimensional space such that when the data is transformed back (“reconstructed”) into its original space, the loss of information is minimal. In this way, subspace learning is intuitively a version of lossy compression. Concretely, say the original space is \mathbb{R}^2 and the data therein is plotted in Figure 2:



(a) Lower dimensional subspace superimposed on the original data

(b) The learned subspace minimizes squared error of projection

Figure 2: The goal of subspace learning is to find a low dimensional linear subspace that captures the most information or variance in the data. It is not to be confused with linear regression.

Given n data points in \mathbb{R}^d , find a projection matrix $U \in \mathbb{R}^{d \times k}$ such that the projection $\hat{x} \in \mathbb{R}^k$ of a data vector $x \in \mathbb{R}^d$ given by $\hat{x} = U^T U x$ minimizes the empirical reconstruction error:

$$error = \frac{1}{n} \sum_{i=1}^n \|x_i - U U^T x_i\|_2^2 = \frac{1}{n} \sum_{i=1}^n \left(\|x_i\|_2^2 - \|U^T x_i\|_2^2 \right) \quad (1)$$

The above equality is derived:

$$\begin{aligned} \|x_i - U U^T x_i\|_2^2 &= (x_i - U U^T x_i)^T (x_i - U U^T x_i) \\ &= x_i^T x_i - x_i^T U U^T x_i - x_i^T U U^T x_i + x_i^T U^T U U U^T x_i \\ &= x_i^T x_i - x_i^T U U^T x_i && \text{recall } U^T U = I \\ &= \|x_i\|_2^2 - \|U^T x_i\|_2^2 \\ &= \|x_i\|_2^2 - \|U^T x_i\|_2^2 \end{aligned}$$

In Equation 1, the term $\frac{1}{n} \sum_{i=1}^n \|x_i\|_2^2$ is constant given any data set X . So in order to minimize the reconstruction, or projection, error, we must maximize $\frac{1}{n} \sum_{i=1}^n \|U^T x_i\|_2^2$. We will show that this is equivalent to maximizing the variance captured by the subspace, that is, maximizing the trace of the empirical covariance matrix of X :

$$\underbrace{\left[\frac{1}{n} \sum_{i=1}^n \|U^\top x_i\|_2^2 \right]}_{\text{Variance captured by } U} = \left[\frac{1}{n} \sum_{i=1}^n x_i^\top U U^\top x_i \right] = \text{trace} \left(\frac{1}{2} \sum_{i=1}^n U^\top x_i x_i^\top U \right) = \text{trace} \left(U^\top \underbrace{\left[\frac{1}{n} \sum_{i=1}^n x_i x_i^\top \right]}_{\text{Emp. Covariance}} U \right) \quad (2)$$

Since U is unique, minimizing the error is equal to maximizing the trace of the empirical covariance matrix of X ; the diagonal of any covariance matrix contains the variance of each data point.

3 Purpose and Applications

As we will see, PCA is a type of subspace learning. PCA is ubiquitous in machine learning, data analysis, information retrieval, and is expanding into data intensive fields like bioinformatics and even drug discovery. The purpose of PCA is two fold. Firstly, to reduce the dimensionality of the feature space: many machine learning models such as clustering have sample complexities that depend exponentially on the dimension of the data. High dimensionality can then lead to poor generalization. The second objective is to learn a unifying and governing structure behind noisy data, specifically, learning uncorrelated components of data. Lower dimensional data can also make the structure more interpretable, allowing causal relationships to be more easily identified. Moreover, independent components of influence behind data makes a system more robust to corruption and noise.

4 PCA

In principle component analysis, we attempt to find a subspace comprised of k normalized linearly independent and uncorrelated basis vectors $\{u_i\}$, $i = 1, \dots, k$ where the i 'th component y_i of x is $y_i = u_i^\top x$ and $\text{Var}(y_1) \geq \text{Var}(y_2) \geq \dots \geq \text{Var}(y_k)$. These principle components in fact yield from the top k eigenvectors of the empirical covariance matrix of X : $C_{XX} = E[xx^\top]$. If a covariance matrix is not at hand, then it suffices to use singular value decomposition of a data matrix normalized and centered for each attribute.

We inspect the mean and variance of y_i :

$$E[y_i] = E[u_i^\top x_i] = u_i^\top E[x_i] = 0 \quad (3)$$

since the data is centered, $E[x_i] = 0$. This simplifies the value of the variance:

$$\text{Var}[y_i] = E[(u_i^\top x_i)^2] - (E[u_i^\top x_i])^2 = u_i^\top E[x_i x_i^\top] u_i = E[(u_i^\top x_i)^2] - 0 \quad (4)$$

4.1 Finding the First Principle Component

The first principal component should capture a largest amount of variance contained in the data vector x . From the definition above we solve the optimization problem:

$$\begin{aligned} & \underset{u_1 \in \mathbb{R}^d}{\text{maximize}} && \text{Var}(u_1^\top x) \\ & \text{subject to} && u_1^\top u_1 = 1 \end{aligned}$$

where $\text{Var}(u_1^\top x) = E[(u_1^\top x)^2] = E[u_1^\top x x^\top u_1] = u_1^\top E[xx^\top] u_1 = u_1^\top C_{xx} u_1$. We solve the maximization problem by defining the Lagrangian: $\mathcal{L} = u_1^\top C_{xx} u_1 + \lambda_1 (1 - u_1^\top u_1)$, for $\lambda_1 \in \mathbb{R}$. Setting $\partial \mathcal{L} / \partial u_1 = 2C_{xx} u_1 - 2\lambda_1 u_1 = 0$ yields $C_{xx} u_1 = \lambda_1 u_1$. Notice u_1 meets the definition of being an eigenvector of the covariance matrix, and the corresponding eigenvalue $\lambda_1 = \lambda_1 u_1^\top u_1 = u_1^\top C_{xx} u_1 = \text{Var}(u_1^\top x)$ is the maximum eigenvalue of C_{xx} , which satisfies the optimization problem.

4.2 Subsequent Principal Components

The second principal component should capture the largest variance in the data set given it is linearly uncorrelated and orthogonal from the first. We would suspect the second principle component to be the second eigenvector of C_{xx} , and eigenvectors corresponding to unique eigenvalues of a symmetric matrix (such as the covariance matrix) are perpendicular.

For y_1 and y_2 to be uncorrelated, their covariance must be zero, i.e. $E[y_1 y_2] - E[y_1]E[y_2] = 0$. This amounts to $E[y_1 y_2] = 0$ by Equation 3. Expanding, $E[y_1 y_2] = E[(u_1^\top x)(u_2^\top x)] = u_1^\top C_{xx} u_2 = \lambda_1 u_1^\top u_2 = 0$, which gives $u_1^\top u_2 = 0$. Conveniently, the constraint necessary for principle components to be uncorrelated is merely that they are orthogonal (and unit normalized). We can now find the second principle component:

$$\begin{aligned} & \underset{u_2 \in \mathbb{R}^d}{\text{maximize}} && \text{Var}(u_2^\top x) \\ & \text{subject to} && u_2^\top u_2 = 1, u_1^\top u_2 = 0 \end{aligned}$$

We proceed as before, defining the Lagrangian as $\mathcal{L} = u_2^\top C_{xx} u_2 + \lambda_2(1 - u_2^\top u_2) + \lambda_{12} u_1^\top u_2$. Setting $\partial \mathcal{L} / \partial u_2 = 2C_{xx} u_2 - 2\lambda_2 u_2 + \lambda_{12} u_1 = 0$ yields $C_{xx} u_2 + \frac{\lambda_{12}}{2} u_1 = \lambda_2 u_2$. Multiplying on left by u_1^\top : $u_1^\top C_{xx} u_2 + \frac{\lambda_{12}}{2} u_1^\top u_1 = \lambda_2 u_1^\top u_2$ yields $\lambda_{12} = 0$. We get finally we get $C_{xx} u_2 = \lambda_2 u_2$, then $\text{Var}(u_2^\top x) = u_2^\top C_{xx} u_2 = \lambda_2$, and u_2 is the eigenvector corresponding to the second largest eigenvalue. Note in this proof we could have also taken the derivative of the Lagrangian with respect to u_1 .

By induction, the variance of the projected data is maximized by the top k eigenvectors of the covariance matrix of data.

4.3 Solving for the Projection Matrix

We saw that the principal components are given as top- k singular vectors of $X = [x_1, \dots, x_n]$, or given as top- k eigenvectors of XX^\top . We must find a projection UU^\top to minimize the reconstruction error. Then U will provide a geometric description of the dominant subspace structure for all data points. In this section we will solve for U .

Theorem 1 (Principal subspace). *Given a set of points $\{x_j\}_{j=1}^n$ in \mathbb{R}^d , the k -dimensional subspace $S \subseteq \mathbb{R}^d$ that best fits these points, i.e. minimizes the residual error*

$$\begin{aligned} & \underset{U \in \mathbb{R}^{d \times k}, y_j \in \mathbb{R}^k}{\text{minimize}} && \sum_{j=1}^n \|x_j - Uy_j\|_2^2 \\ & \text{subject to} && U^\top U = I_k \end{aligned}$$

is spanned by the top- k eigenvectors of XX^\top , where $X \in \mathbb{R}^{d \times n} = [x_1, \dots, x_n]$.

As alluded to in Figure 2, if U were known, this would be a least squares regression problem. Since it is not known, PCA is a joint optimization problem. We solve for U by again defining the lagrangian:

$$\mathcal{L} = \sum_{j=1}^n \|x_j - Uy_j\|_2^2 + \text{trace}\left((I - U^\top U) \Lambda\right)$$

for $\Lambda = \Lambda^\top$. Setting $\frac{\partial \mathcal{L}}{\partial y_j} = -2U^\top(x_j - Uy_j) = 0$ yields $y_j = U^\top x_j$. Upon substitution, notice that

$$\begin{aligned}
\sum_{j=1}^n \|x_j - Uy_j\|_2^2 &= \sum_{j=1}^n (x_j - UU^\top x_j)^\top (x_j - UU^\top x_j) \\
&= \sum_{j=1}^n (x_j^\top x_j - x_j^\top UU^\top x_j) \\
&= \text{trace}(X^\top X - X^\top UU^\top X)
\end{aligned}$$

Upon inspection, minimizing, $\text{trace}(X^\top X - X^\top UU^\top X)$ is equivalent to maximizing $\text{trace}(X^\top UU^\top X)$ over all $U \in \mathbb{R}^{d \times k}$ since X is given. We now have a new objective function: $\text{trace}(X^\top X - X^\top UU^\top X)$, so we rewrite

$$\begin{aligned}
&\underset{U \in \mathbb{R}^{d \times k}}{\text{maximize}} && \text{trace}(X^\top X - X^\top UU^\top X) \\
&\text{subject to} && U^\top U = I_k
\end{aligned} \tag{5}$$

The Lagrangian is given by:

$$\mathcal{L} = \text{trace}(X^\top UU^\top X) + \text{trace}\left(\left(I_k - U^\top U\right) \Lambda\right)$$

for a $\Lambda = \Lambda^\top$. Setting $\partial \mathcal{L} / \partial U = 2XX^\top U - 2U\Lambda = 0$ yields $XX^\top U = U\Lambda$.

There is one small problem here: Λ is not guaranteed to be diagonal. Write $\Lambda = U^\top XX^\top U$ a singular value decomposition for XX^\top . The objective for Equation 5 is then $\text{trace}(\Lambda)$ since $\text{trace}(U^\top XX^\top U) = \text{trace}(X^\top UU^\top X)$. If U is a solution to Equation 5 then so is $\tilde{U} = UV$ for any orthogonal V . Choose V to be eigenvectors of $\Lambda = V\tilde{\Lambda}V^\top$. Then $XX^\top \tilde{U} = (XX^\top U)V = (U^\top)^{-1}\Lambda V = UV\tilde{\Lambda}V^\top V = (UV)\tilde{\Lambda} = \tilde{U}\tilde{\Lambda}$ is the solution to Equation 5 with objective function $\text{trace}(\tilde{\Lambda})$. Notice, the columns of \tilde{U} are eigenvectors of XX^\top with corresponding eigenvalues in $\tilde{\Lambda}$.

Since objective is $\text{trace}(\tilde{\Lambda})$, \tilde{U} must be the top- k eigenvectors. Also, the top k eigenvectors of $XX^\top \equiv$ top k left singular vectors of X , meaning that SVD also gives an optimal solution to PCA. We can further write that $Y = U_k^\top X = [y_1, \dots, y_n] \in \mathbb{R}^{k \times n}$ is a more compact representation for X as typically $k \ll d$.

5 PCA Algorithms

5.1 Exact Solution: When the Covariance Matrix is computable

Below is the algorithm to compute the optimal projection matrix for PCA. The input is a data matrix $X \in \mathbb{R}^{d \times n}$ for n examples each with dimension d .

```

pca( $k : \mathbb{N}, X = [x_1, \dots, x_n] : \mathbb{R}^{d \times n}$ )
1    $X \leftarrow X - X\mathbf{1}_n\mathbf{1}_n^\top$ ;
2    $\hat{C}_{xx} \leftarrow \frac{1}{n-1}XX^\top$ ;
3    $U, \lambda \leftarrow \text{eig}(\hat{C}_{xx})$ ;
4   return  $U_{1:k}, \lambda_{1:k}$ ;

```

Algorithm 1: Principal component analysis.

The computational cost of computing the eigendecomposition above of a matrix of size $d \times d$ is $O(d^3)$. Furthermore, covariance matrix of $X \in \mathbb{R}^{d \times k}$ is $\frac{1}{n}XX^\top \in \mathbb{R}^{d \times d}$. Usually, if PCA is applied to a dataset, it means d is quite large, so the exact computation of U can be prohibitive. For this reason, efficient approximation algorithms for U are given below.

5.2 Approximation 1: Power Iteration Method

Given a real matrix $A \in \mathbb{R}^{d \times d}$, such that $A^\top = A$, such that

- A has d linearly independent eigenvectors, $\{v_1, \dots, v_d\}$
- The eigenvalues can be ordered in magnitude as $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_d| > 0$

input symmetric matrix $A \in \mathbb{R}^{d \times d}$; number of iterations m

1. Sample $v^{(0)} \in \mathbb{R}^d$ randomly
2. **for** $i = 1, 2, \dots, m$ **do**
3. $\tilde{v}^{(i)} = A v^{(i-1)}$
4. $v^{(i)} = \frac{\tilde{v}^{(i)}}{\|v^{(i)}\|_2}$
5. **end for**

output $\hat{v}_1 = v^{(m)}$, the approximate eigenvector

The algorithm returns an approximation for the first principal component. We prove this by first noting that any vector $v^{(0)} \in \mathbb{R}^d$ can be written as $v^{(0)} = c_1 v_1 + c_2 v_2 + \dots + c_d v_d$ where $\{v_1, \dots, v_d\}$ is the set of linearly independent eigenvectors of A . Then

$$\begin{aligned} A v^{(0)} &= c_1 \lambda_1 v_1 + \dots + c_d \lambda_d v_d \\ A^2 v^{(0)} &= c_1 \lambda_1^2 v_1 + \dots + c_d \lambda_d^2 v_d \\ &\dots \\ A^m v^{(0)} &= c_1 \lambda_1^m v_1 + \dots + c_d \lambda_d^m v_d \end{aligned}$$

Dividing by λ_1^m gives

$$\frac{1}{\lambda_1^m} A^m v^{(0)} = c_1 v_1 + c_2 \left(\frac{\lambda_2}{\lambda_1}\right)^m \dots + c_d \left(\frac{\lambda_d}{\lambda_1}\right)^m$$

As $m \rightarrow \infty$, $\frac{\lambda_j}{\lambda_1} \rightarrow 0$ for all $j > 1$ and only the term $c_1 v_1$ survives for large m : $\frac{1}{\lambda_1^m} A^m v^{(0)} = c_1 v_1$ where $\hat{v}_1 = A^m v^{(0)}$ is approximately an eigenvector corresponding to λ_1 . If we choose $v^{(0)} \not\propto v_1 \Rightarrow c_1 \neq 0$, then when convergence is reached, for any $y \not\propto v_1$, we can approximate λ_1 as the ratio of two terms:

$$\begin{aligned} \frac{1}{\lambda_1^m} \langle A^m v^{(0)}, y \rangle &= c_1 v_1^\top y \\ \frac{1}{\lambda_1^{m+1}} \langle A^{m+1} v^{(0)}, y \rangle &= c_1 v_1^\top y \end{aligned}$$

which is

$$\frac{y^\top A^{m+1} v^{(0)}}{y^\top A^m v^{(0)}} = \frac{\lambda_1^{m+1}}{\lambda_1^m} = \lambda_1$$

We stop iterating when relative change in estimated eigenvalue is small. To find subsequent eigenvectors, run the power iteration method on $A = A - \lambda_1 \hat{v}_1 \hat{v}_1^\top$; recall $A = U\Lambda U^\top = \sum_{i=1}^n \lambda_i v_i v_i^\top$

Note: the power iteration method converges with fewer iterations if there is a large eigengap, meaning $\lambda_i \gg \lambda_j$ for any $i > j$. Typically, the first few eigenvectors are discovered readily, while it takes exponentially more iterations to converge to less dominant eigenvalues/eigenvectors. The other consideration to note is that if $\lambda_i = \lambda_j$ for $i > j$, then the returned eigenvector approximations \hat{v}_1, \hat{v}_2 will span the space of the true eigenvectors v_1 and v_2 , that is, the algorithm may not recover the true eigenvectors, but rather a linear combination thereof.

5.3 Approximation 2: Orthogonal Iteration Method

Below is an algorithm to approximate the top- k eigenvectors simultaneously.

input symmetric matrix $A \in \mathbb{R}^{d \times d}$; number of iterations m

1. Sample $U_0 \in \mathbb{R}^{d \times k}$ such that $U_0^\top U_0 = I_k$
2. **for** $i = 1, 2, \dots, m$ **do**
3. $\tilde{U}^{(i)} = AU^{(i-1)}$.
4. $U^{(i)} = [\text{matrix of orthonormal basis vectors for range } (\tilde{U}^{(i)})]$
5. **end for**

output $U = U^{(m)}$

where U_0 can be constructed from any randomly sampled matrix by performing Gram-Schmidt.

6 Dimensionality Considerations

6.1 Not Enough Samples

What if the dimension d of any vector in the space spanned by the data is much larger than n the number of examples? [2] This might occur, say, with having only a few hundred very high resolution pictures of millions of pixels. Trivially, it makes no sense to find k principle components for $k > n - 1$ because at least $d - n + 1$ eigenvalues of the data matrix are zero.

As before, define the empirical covariance matrix to be $C = \frac{1}{n}XX^\top$, an eigenvector u_i of which satisfies $\frac{1}{n}XX^\top u_i = \lambda_i u_i$. If we premultiply by X to obtain:

$$\frac{1}{n}XX^\top(Xu_i) = \lambda_i(Xu_i)$$

and define $v_i = Xu_i$ an eigenvector for $XX^\top \in \mathbb{R}^{n \times n}$, then we can perform the eigendecomposition in $O(n^3)$ time rather than $O(d^3)$. To determine the eigenvectors, again premultiply the above by X^\top to give:

$$\frac{1}{n}(X^\top X)(X^\top v_i) = \lambda_i(X^\top v_i)$$

and notice that $X^\top v_i$ is an eigenvector for C with eigenvalue λ_i . The only drawback is that $X^\top v_i$ may not be unit norm, so normalize. It is remarkable that to calculate PCA in this case, only the inner product

was used. This sets the foundation for kernel PCA algorithms, which operate in specialized inner product spaces.

6.2 Prior Assumption of Sparsity

Compressed sensing is a form of dimensionality reduction with the assumption that the data is sparse in some well-chosen basis. This is often the case, in fact, JPEG-2000 compression uses the assumption that images are sparse in a wavelet-basis [4]

References

- [1] Bengio, Yoshua and Courville, Aaron and Vincent, Pascal, “Representation learning: A review and new perspectives”, *Cambridge University Press*, 2006.
- [2] Bishop, Christopher M. *Pattern Recognition and Machine Learning*. New York: Springer, 2006. Print.
- [3] Mohri, Mehryar, Afshin Rostamizadeh, and Ameet Talwalkar. ”Foundations of Machine Learning.” MIT Press, 2012
- [4] Shalev-Shwartz, Shai, and Shai Ben-David. *Understanding Machine Learning*. New York, NY: Cambridge UP, 2014. Print.