

Notes on PCA and applications to k -means

Chris Schwiegelshohn

Consider a set of n samples x_1, \dots, x_n , with $x_i \in \mathbb{R}$. The *empirical mean* of these samples is defined as

$$\mathbf{E}[x_i] \triangleq \frac{1}{n} \sum_{i=1}^n x_i$$

and the *empirical variance* is defined as

$$\mathbf{Var}[x_i] \triangleq \mathbf{E}[(x_i - \mathbf{E}[x_i])^2] = \mathbf{E}[x_i^2] - \mathbf{E}[x_i]^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \sum_{i=1}^n x_i \right)^2.$$

As a matter of fact, we often normalize the data such that the mean is 0 and then the variance reduces to $\frac{1}{n} \sum_{i=1}^n x_i^2$. In this case, we may interpret the variance as the (scaled) squared Euclidean norm of the vector containing the samples. In general, for a vector $\mathbf{x} \in \mathbb{R}^n$, with $\mathbf{x} = (x_1, \dots, x_n)^T$, $\|\mathbf{x}\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}$, hence for normalized data sets, $\mathbf{Var}[x_i] = \frac{1}{n} \|\mathbf{x}\|_2^2$, where \mathbf{x} is the vector containing all the samples: $\mathbf{x} = (x_1, \dots, x_n)$. We also note that $\mathbf{x}^T \mathbf{x} = \|\mathbf{x}\|_2^2$ for any vector \mathbf{x} .

Let us now consider these notions in higher dimensions, that is, the samples of $\mathbf{A}_{(i)}$ are no longer numbers, but vectors in \mathbb{R}^d . The empirical mean translates straightforwardly and is also commonly known as the centroid. The notion of variance is not as easy to generalize. Ideally, we would like to retain the notion that the variance quantifies the spread of the data set with respect to the mean (or centroid). The difficulty of extending this notion is that the spread is different along different directions. This is properly captured by the *covariance matrix*. Our notion of generalization will be simpler, as we are looking for a single number, rather than the more complex spectral structure included in the covariance matrix. Instead, we define the *directional variance* along an arbitrary unit vector \mathbf{v} as

$$\mathbf{Var}_{\mathbf{v}}[\mathbf{A}_{(i)}] = \mathbf{E}[(\mathbf{A}_{(i)} - \mathbf{E}[\mathbf{A}_{(i)}])^T \mathbf{v}]^2.$$

Again, for normalized inputs with $\mathbf{E}[\mathbf{A}_{(i)}] = \mathbf{0}$, this reduces to

$$\mathbf{Var}_{\mathbf{v}}[\mathbf{A}_{(i)}] = \mathbf{E}[(\mathbf{A}_{(i)}^T \mathbf{v})^2].$$

Geometrically, this expression means that we project all points along the direction \mathbf{v} and compute the variance of a (now) 1-dimensional set of samples. To capture the entire variance of the point set, we pick an arbitrary orthogonal basis $V = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d\}$ of \mathbb{R}^d and compute

$$\mathbf{Var}[\mathbf{A}_{(i)}] \triangleq \sum_{j=1}^d \mathbf{Var}[(\mathbf{A}_{(i)}^T \mathbf{v}_j)^2].$$

We note that this definition is invariant under any choice of orthogonal basis, that is, for two distinct candidates $V = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d\}$ and $W = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d\}$, the evaluation of

$\mathbf{Var}[\mathbf{A}_{(i)}]$ is still identical. Let us prove this. First, observe that vector \mathbf{v}_k can be expressed as a linear combination $\mathbf{v}_k = \sum_{j=1}^d \alpha_{k,j} \cdot \mathbf{w}_j$ and any vector \mathbf{w}_j can be expressed as a linear combination $\mathbf{w}_j = \sum_{k=1}^d \beta_{j,k} \cdot \mathbf{v}_k$. Notice that, because $\mathbf{w}_r^T \mathbf{w}_r = 1$ and $\mathbf{w}_r^T \mathbf{w}_j = 0$ for $r \neq j$, we have that:

$$\mathbf{v}_k^T \mathbf{w}_j = \sum_{r=1}^d \alpha_{k,r} \cdot \mathbf{w}_r^T \mathbf{w}_j = \alpha_{k,j}$$

and, similarly, that

$$\mathbf{w}_j^T \mathbf{v}_k = \sum_{r=1}^d \beta_{j,r} \cdot \mathbf{v}_r^T \mathbf{v}_k = \beta_{j,k}.$$

Therefore, $\beta_{j,k} = \alpha_{k,j}$, so

$$\mathbf{w}_j = \sum_{k=1}^d \alpha_{k,j} \cdot \mathbf{v}_k.$$

Notice also that

$$\begin{aligned} 1 &= \mathbf{w}_j^T \mathbf{w}_j = \left(\sum_{k=1}^d \alpha_{k,j} \cdot \mathbf{v}_k \right)^T \cdot \left(\sum_{k=1}^d \alpha_{k,j} \cdot \mathbf{v}_k \right) \\ &= \sum_{k=1}^d \alpha_{k,j}^2 \mathbf{v}_k^T \mathbf{v}_k + \sum_{k=1}^d \sum_{\substack{r=1 \\ r \neq k}}^d \alpha_{k,r} \mathbf{v}_k^T \mathbf{v}_r = \sum_{j=1}^d \alpha_{k,j}^2, \end{aligned} \tag{1}$$

and for $j \neq r$

$$\begin{aligned} 0 &= \mathbf{w}_j^T \mathbf{w}_r = \left(\sum_{k=1}^d \alpha_{k,j} \cdot \mathbf{v}_k \right)^T \cdot \left(\sum_{k=1}^d \alpha_{k,r} \cdot \mathbf{v}_k \right) \\ &= \sum_{k=1}^d \alpha_{k,j} \alpha_{k,r} \mathbf{v}_k^T \mathbf{v}_k + \sum_{k=1}^d \sum_{\substack{\ell=1 \\ \ell \neq k}}^d \alpha_{k,j} \alpha_{\ell,r} \mathbf{v}_k^T \mathbf{v}_\ell = \sum_{k=1}^d \alpha_{k,j} \alpha_{k,r}. \end{aligned} \tag{2}$$

Then,

$$\begin{aligned} \mathbf{Var}_{\mathbf{v}_k}[\mathbf{A}_{(i)}] &= \mathbf{E}[(\mathbf{A}_{(i)} - \mathbf{E}[\mathbf{A}_{(i)}])^T \mathbf{v}_k]^2 = \mathbf{E} \left[\left((\mathbf{A}_{(i)} - \mathbf{E}[\mathbf{A}_{(i)}])^T \sum_{j=1}^d \alpha_{k,j} \cdot \mathbf{w}_j \right)^2 \right] \\ &= \mathbf{E} \left[\sum_{j=1}^d ((\mathbf{A}_{(i)} - \mathbf{E}[\mathbf{A}_{(i)}])^T \cdot \mathbf{w}_j)^2 \alpha_{k,j}^2 \right] \\ &\quad + \mathbf{E} \left[\sum_{j=1}^d \sum_{\substack{r=1 \\ r \neq j}}^d (\mathbf{A}_{(i)} - \mathbf{E}[\mathbf{A}_{(i)}])^T \cdot \mathbf{w}_j \cdot (\mathbf{A}_{(i)} - \mathbf{E}[\mathbf{A}_{(i)}])^T \cdot \mathbf{w}_r \cdot \alpha_{k,j} \alpha_{k,r} \right] \\ &= \sum_{j=1}^d \mathbf{Var}_{\mathbf{w}_j}[\mathbf{A}_{(i)}] \cdot \alpha_{k,j}^2 + \mathbf{E} \left[\sum_{j=1}^d \sum_{\substack{r=1 \\ r \neq j}}^d \mathbf{x}^T \mathbf{w}_j \mathbf{x}^T \mathbf{w}_r \cdot \alpha_{k,j} \alpha_{k,r} \right], \end{aligned}$$

if we define $\mathbf{x} = \mathbf{A}_{(i)} - \mathbf{E}[\mathbf{A}_{(i)}]$. Summing up over all \mathbf{v}_k , we then obtain

$$\sum_{k=1}^d \mathbf{Var}_{\mathbf{v}_k}[\mathbf{A}_{(i)}] = \sum_{j=1}^d \mathbf{Var}_{\mathbf{w}_j}[\mathbf{A}_{(i)}] \cdot \sum_{k=1}^d \alpha_{k,j}^2 + \mathbf{E} \left[\sum_{j=1}^d \sum_{\substack{r=1 \\ r \neq j}}^d \mathbf{x}^T \mathbf{w}_j \mathbf{x}^T \mathbf{w}_r \cdot \sum_{k=1}^d \alpha_{k,j} \alpha_{k,r} \right] = \sum_{j=1}^d \mathbf{Var}_{\mathbf{w}_j}[\mathbf{A}_{(i)}],$$

using Equations (1) and (2).

As in the one dimensional case, our notion of high dimensional variance has an algebraic interpretation. The *Frobenius norm* of a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ is defined as

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d A_{i,j}^2}.$$

If the centroid is equal to the origin, the squared Frobenius norm is, up to scale, equal to the multidimensional variance, as well as the 1-means cost. To see the former, consider the basis $\{\mathbf{e}_k\}_{k=1}^n$, where \mathbf{e}_k is the vector that is equal to 1 at the k th coordinate and 0 everywhere else. We have

$$\mathbf{Var}_{\mathbf{e}_k}[\mathbf{A}_{(i)}] = \mathbf{E}[\mathbf{A}_{(i)}^T \mathbf{e}_k] = \frac{1}{n} \sum_{i=1}^n \mathbf{A}_{(i)k}^2 = \frac{1}{n} \sum_{i=1}^n \mathbf{A}_{i,k}^2$$

and

$$\|\mathbf{A}\|_F^2 = \sum_{k=1}^d \sum_{i=1}^n \mathbf{A}_{i,k}^2 = n \sum_{k=1}^d \mathbf{Var}_{\mathbf{e}_k}[\mathbf{A}_{(i)}] = \mathbf{Var}[\mathbf{A}]$$

Principal Component Analysis is all about dimensionality reduction. As a tentative step, let us consider reducing the dimension down to 1. The main question is which direction is the most important one. Our notion of directional variance helps us in this regard. If a direction has extremely low directional variance, we can confidently say that the centroid (or origin if our data are normalized) will approximate the point set well enough. The most uncertainty is with respect to directions of high directional variance. Hence, if we are only allowed to choose a single direction, we should choose the one with maximum directional variance. Phrased as an optimization problem, we aim to solve the following.

$$\max_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|=1} \mathbf{Var}_{\mathbf{v}}[\mathbf{A}_{(i)}].$$

Again, this has an algebraic interpretation. Specifically, the maximum directional variance is (up to scaling) known as the squared spectral norm, where for any n by d matrix \mathbf{A} the spectral norm is defined as

$$\|\mathbf{A}\|_2 = \max_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2 = \max_{\mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_2=1} \sqrt{\sum_{i=1}^n (\mathbf{A}_{(i)}^T \mathbf{v})^2}.$$

The connection to algebra is deeper. We will see shortly that the vector \mathbf{v} that induces the spectral norm of \mathbf{A} is an eigenvector of $\mathbf{A}^T \mathbf{A}$ and a right singular vector of \mathbf{A} . The directional variance itself is (up to scaling) also the largest eigenvalue of $\mathbf{A}^T \mathbf{A}$. Let us first recall the following definitions.

Definition 1 (Singular Vectors and Values). *Let $\mathbf{A} \in \mathbb{R}^{n \times d}$. Two vectors $\mathbf{u} \in \mathbb{R}^n$ and $\mathbf{v} \in \mathbb{R}^d$ with unit Euclidean norm are respectively called left and right singular vectors of \mathbf{A} if the following two equations hold*

- $\mathbf{A}\mathbf{v} = \sigma\mathbf{u}$
- $\mathbf{u}^T\mathbf{A} = \sigma\mathbf{v}^T$.

σ is known as a singular value of \mathbf{A} .

Definition 2 (Eigenvectors and Eigenvalues). Let $\mathbf{A} \in \mathbb{R}^{d \times d}$. A vector $v \in \mathbb{R}^d$ with unit Euclidean norm is an eigenvector with eigenvalue e if

- $\mathbf{A}\mathbf{v} = \sigma\mathbf{v}$ and
- $\mathbf{v}^T\mathbf{A} = \sigma\mathbf{v}^T$.

Proposition 3. Let \mathbf{A} be a matrix with right singular vector \mathbf{v} and singular value σ . Then \mathbf{v} is an eigenvector of $\mathbf{A}^T\mathbf{A}$ with eigenvalue σ^2 .

Proof. $\mathbf{A}^T\mathbf{A}\mathbf{v} = \mathbf{A}^T\mathbf{u}\sigma = \mathbf{v}^T\sigma^2$ and $\mathbf{v}^T\mathbf{A}^T\mathbf{A} = \sigma\mathbf{u}^T\mathbf{A} = \sigma^2\mathbf{v}^T$. □

This shows that the largest eigenvalue of $\mathbf{A}^T\mathbf{A}$ and the squared largest singular value of \mathbf{A} are equivalent. Let us further show that these are equivalent to the squared spectral norm.

Theorem 4. Let \mathbf{A} be a matrix. Then the spectral norm $\|\mathbf{A}\|_2$ is equal to the largest singular value of \mathbf{A} .

Proof. Let \mathbf{v} be the vector that induces the spectral norm of \mathbf{A} . Then $\|\mathbf{A}\mathbf{v}\|_2^2 = \mathbf{v}^T\mathbf{A}^T\mathbf{A}\mathbf{v}$. Let us consider \mathbf{v} as a linear combination of eigenvectors of $\mathbf{A}^T\mathbf{A}$, that is, $\mathbf{v} = \sum_{i=1}^d \alpha_i \mathbf{v}_i$ with $\sum_{i=1}^d \alpha_i^2 = 1$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_j\}$ being an orthogonal basis of eigenvectors of $\mathbf{A}^T\mathbf{A}$. Then

$$\begin{aligned} \mathbf{v}^T\mathbf{A}^T\mathbf{A}\mathbf{v} &= \left(\sum_{i=1}^d \alpha_i \mathbf{v}_i^T \right) \mathbf{A}^T\mathbf{A} \left(\sum_{j=1}^d \alpha_j \mathbf{v}_j \right) \\ &= \left(\sum_{i=1}^d \alpha_i \sigma_i \mathbf{v}_i^T \right) \left(\sum_{j=1}^d \alpha_j \sigma_j \mathbf{v}_j \right) \\ &= \sum_{i=1}^d \alpha_i^2 \mathbf{v}_i^T \mathbf{v}_i \sigma_i^2 + \sum_{i \neq j} \alpha_i \alpha_j \mathbf{v}_i^T \mathbf{v}_j \sigma_i \sigma_j. \end{aligned}$$

Since V is an orthogonal basis $\mathbf{v}_i^T \mathbf{v}_j = 0$ and $\mathbf{v}_i^T \mathbf{v}_i = 1$. Hence

$$\mathbf{v}^T\mathbf{A}^T\mathbf{A}\mathbf{v} = \sum_{i=1}^d \alpha_i^2 \sigma_i^2 \leq \sum_{i=1}^d \alpha_i^2 \max_{1 \leq j \leq d} \sigma_j^2 = \max_{1 \leq j \leq d} \sigma_j^2.$$

□

Theorem 4 and Proposition 3 tell us that it is sufficient to find the largest eigenvalue of $\mathbf{A}^T\mathbf{A}$ and associated eigenvector to determine the maximum directional variance. Obviously, this is a very well studied problem and we will not discuss this in detail. We will only present a simple algorithm known as the power method.

1. Pick a random unit vector \mathbf{x}
2. Repeat $\mathbf{x} \leftarrow \frac{\mathbf{A}^T\mathbf{A}\mathbf{x}}{\|\mathbf{A}^T\mathbf{A}\mathbf{x}\|_2}$

To analyze this algorithm, we require two steps. First, we will show that a random unit vector \mathbf{x} is not too far off from the target eigenvector \mathbf{v} . Second, we show that under certain assumptions, this method will quickly converge.

To pick a random unit vector in d dimensions, we choose d random and independent Gaussian variables X_1, \dots, X_d with mean 0 and variance 1. We then set $\mathbf{x}^T = (X_1, \dots, X_d) \cdot \frac{1}{\sqrt{\sum_{i=1}^d X_i^2}}$.

Then the following holds

Lemma 5. *Fix a unit vector \mathbf{v} . For any random unit vector \mathbf{x} , we have $\mathbf{E}[(\mathbf{x}^T \mathbf{v})^2] \geq \frac{1}{d}$.*

Proof. We first observe that because \mathbf{x} is chosen uniformly at random from all unit vectors, we may assume \mathbf{v} to be the unit vector along the first axis, that is, $\mathbf{v}^T = (1, 0, \dots, 0)$. Then

$$\mathbf{E}[(\mathbf{x}^T \mathbf{v})^2] = \mathbf{E}\left[X_1^2 \cdot \frac{1}{\sum_{i=1}^d X_i^2}\right] = \mathbf{E}\left[\frac{1}{1 + \sum_{i=2}^d \frac{X_i^2}{X_1^2}}\right] \geq \frac{1}{\mathbf{E}\left[1 + \sum_{i=2}^d \frac{X_i^2}{X_1^2}\right]} = \frac{1}{1 + \sum_{i=2}^d \frac{\mathbf{E}[X_i^2]}{\mathbf{E}[X_1^2]}} = \frac{1}{d}$$

where the inequality follows from Jensen's inequality¹ and the last two equalities from the fact that X_1 and X_i (for $i \geq 2$) are independent and identically distributed. \square

We cannot give good converge speeds of the power method in general. But if the gap between first and second eigenvalue is small, the method works pretty well. To see this, let $\mathbf{x} = \sum_{i=1}^d \alpha_i \mathbf{v}_i$ be the initial choice of \mathbf{x} . Let us now consider the vector $\mathbf{x}^{(k)}$, that is, the current candidate \mathbf{x} after k iterations. We have (up to scaling)

$$\mathbf{x}^{(k)} = (\mathbf{A}^T \mathbf{A})^k \mathbf{x} = \sum_{i=1}^d \alpha_i \mathbf{v}_i \sigma_i^{2k} = \sigma_1^{2k} \cdot \left(\alpha_1 \mathbf{v}_1 + \sum_{i=2}^d \alpha_i \mathbf{v}_i \left(\frac{\sigma_i}{\sigma_1}\right)^{2k} \right).$$

We note that since σ_1^2 is the largest eigenvalue, the term $\sum_{i=2}^d \alpha_i \mathbf{v}_i \left(\frac{\sigma_i}{\sigma_1}\right)^{2k}$ will eventually converge to zero. The time required for this to happen depends on the ratio between σ_1^2 and the next largest eigenvalue σ_2^2 . In practice, we may assume that $\sigma_2 \cdot (1 + \varepsilon) < \sigma_1$ for a few reasons. First, data are often noisy, so the bad case of having an extremely small gap between the two eigenvalues is unlikely. Second, if the values are sufficiently close, we may simply be satisfied with any vector that is a linear combination of \mathbf{v}_1 and \mathbf{v}_2 (and possibly $\mathbf{v}_3, \mathbf{v}_4, \dots$, if more eigenvalues are extremely close). We note that α_1 must not be equal to 0 for this algorithm to converge to the largest eigenvector. Hence, the random initialization.

Best-Fit Subspaces We now consider dimension reduction onto more than just a single dimension. Here, we reuse many of the notions we have seen before.

Again, formulated as an optimization problem, we are looking for k vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ such that the directional variances $\sum_{j=1}^k \mathbf{Var}_{\mathbf{v}_j}[\mathbf{A}_{(i)}]$ are maximized. Algebraically, we aim to find a $d \times k$ orthogonal matrix \mathbf{V} , such that $\|\mathbf{AV}\|_F^2$ is maximized. To see this, consider

$$n \cdot \sum_{j=1}^k \mathbf{Var}_{\mathbf{v}_j}[\mathbf{A}_{(i)}] = \sum_{j=1}^k \|\mathbf{Av}_j\|_2^2 = \|\mathbf{AV}\|_F^2.$$

The solution to this problem essentially boils down to finding a the maximum directional variance, and then projecting the entire point set onto the orthogonal subspace.

¹Jensen's inequality says that if a function $f(x)$ is convex, then for a random variable X we have that $\mathbf{E}[f(X)] \geq f(\mathbf{E}[X])$.

Theorem 6. Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix. Then the best fit subspace is characterized as follows.

$$\begin{aligned} \mathbf{v}_1 &= \operatorname{argmax}_{\|\mathbf{v}\|_2=1} \|\mathbf{A}\mathbf{v}\|_2 \\ \mathbf{v}_2 &= \operatorname{argmax}_{\|\mathbf{v}\|_2=1, \mathbf{v}^T \mathbf{v}_1=0} \|\mathbf{A}\mathbf{v}\|_2 \\ &\vdots \\ \mathbf{v}_k &= \operatorname{argmax}_{\substack{\|\mathbf{v}\|_2=1, \\ \mathbf{v}^T \mathbf{v}_i=0, i \in \{1, \dots, k-1\}}} \|\mathbf{A}\mathbf{v}\|_2 \end{aligned}$$

Further, any matrix may be written as

$$\mathbf{U}\mathbf{D}\mathbf{V}^T$$

where \mathbf{U} and \mathbf{V} are orthogonal matrices containing the singular vectors and \mathbf{D} is a diagonal matrix whose entries contain the singular values. This decomposition is known as the singular value decomposition.

Proof. We prove this by induction. The base case was proven in Theorem 4. Assume that W_k is an optimal subspace with basis vectors $\mathbf{w}_1, \dots, \mathbf{w}_k$, and let \mathbf{W}_k be the matrix in $\mathbb{R}^{k \times d}$ with rows $\mathbf{w}_1, \dots, \mathbf{w}_k$.

We then have

$$\|\mathbf{A}\mathbf{W}_k\|_F^2 = \sum_{i=1}^k \|\mathbf{A}\mathbf{w}_i\|_2^2.$$

By assumption $\|\mathbf{A}\mathbf{w}_i\|_2^2 \leq \|\mathbf{A}\mathbf{v}_i\|_2^2$ for all $i \in \{1, \dots, k-1\}$. Hence

$$\|\mathbf{A}\mathbf{W}_k\|_F^2 \leq \sum_{i=1}^{k-1} \|\mathbf{A}\mathbf{v}_i\|_2^2 + \|\mathbf{A}\mathbf{w}_k\|_2^2.$$

By optimality of \mathbf{W}_k , we can conclude that \mathbf{W}_k is in the subspace spanned by $\{\mathbf{v}_1, \dots, \mathbf{v}_{k-1}, \mathbf{w}_k\}$. Then there exists a vector \mathbf{x} in that subspace orthogonal to $\mathbf{v}_1, \dots, \mathbf{v}_{k-1}$ such that

$$\|\mathbf{A}\mathbf{W}_k\|_F^2 - \sum_{i=1}^{k-1} \|\mathbf{A}\mathbf{v}_i\|_2^2 = \|\mathbf{A}\mathbf{x}\|_2^2.$$

Since \mathbf{v}_k optimized over the entire space (and not just the space \mathbf{W}_k , we know that $\|\mathbf{A}\mathbf{v}_k\|_2^2 \geq \|\mathbf{A}\mathbf{x}\|_2^2$, which concludes the proof of the first statement.

For the second statement, we simply let $k = d$. □

To compute the best subspace, we may run the power method k times, projecting onto the orthogonal subspace whenever we converge.

Applications to k -means We first would like to consider the sister minimization problem. We first observe that by the Pythagorean theorem, for any unit vector v we have

$$\|\mathbf{A}\|_F^2 = \|\mathbf{A}\mathbf{v}\mathbf{v}^T\|_2^2 + \|\mathbf{A}(I - \mathbf{v}\mathbf{v}^T)\|_F^2 = \|\mathbf{A}\mathbf{v}\|_2^2 + \|\mathbf{A} - \mathbf{A}\mathbf{v}\mathbf{v}^T\|_F^2.$$

Hence maximizing $\|\mathbf{A}\mathbf{v}\|_2^2$ is equivalent to minimizing $\|\mathbf{A} - \mathbf{A}\mathbf{v}\mathbf{v}^T\|_F^2$. These notions straightforwardly extend to subspaces, that is,

$$\|\mathbf{A}\|_F^2 = \|\mathbf{A}\mathbf{W}\|_2^2 + \|\mathbf{A} - \mathbf{A}\mathbf{W}\mathbf{W}^T\|_F^2.$$

Minimizing

$$\|\mathbf{A} - \mathbf{A}\mathbf{W}\mathbf{W}^T\|_F^2$$

for a rank k subspace \mathbf{W} is known in literature as finding the best rank- k approximation.

Let us now study the relationship between \mathbf{A} and $\mathbf{A}\mathbf{V}_k\mathbf{V}_k^T$. We first note that $\mathbf{A}\mathbf{V}_k\mathbf{V}_k^T = \mathbf{U}\mathbf{D}\mathbf{V}^T\mathbf{V}_k\mathbf{V}_k^T = \mathbf{U}\mathbf{D}\mathbf{V}_k^T = \mathbf{U}_k\mathbf{D}_k\mathbf{V}_k^T = \mathbf{U}_k\mathbf{U}_k^T\mathbf{A}$, where \mathbf{U}_k is obtained from \mathbf{U} by just taking the first k singular vectors and setting everything else to 0 and \mathbf{D}_k is the diagonal matrix obtained from \mathbf{D} such that all but the first k diagonal entries of \mathbf{D} were set to 0. Note that we may also remove the 0 entries from \mathbf{U}_k just like we do in \mathbf{V}_k , that is, consider \mathbf{U}_k to be a $n \times k$ matrix. Maximizing $\|\mathbf{A}\mathbf{W}\|_F^2$ for some rank k column subspace \mathbf{W} is therefore equivalent to maximizing $\|\mathbf{T}\mathbf{A}\|_F^2$ for some rank k row subspace \mathbf{T} .

Now consider the 1-means objective function, where we aim to find a point μ such that $\sum_{i=1}^n \|\mathbf{A}_{(i)} - \mathbf{c}\|^2$ is minimized. We know that $\mathbf{c} = \mu = \frac{1}{n} \sum_{i=1}^n \mathbf{A}_{(i)}$ is optimal. Can we express this problem algebraically? Indeed, we can. Let us rewrite the one means objective as

$$\sum_{i=1}^n \|\mathbf{A}_{(i)} - \mathbf{c}\|^2 = \|\mathbf{A} - \mathbf{C}\|_F^2$$

with the constraint that every row of \mathbf{C} is identical. Consider the vector $\mathbf{X} = \frac{1}{\sqrt{n}} \cdot \mathbf{1}$. Then the optimal matrix \mathbf{C} , that is, the matrix where every row is μ can be expressed as $\mathbf{X}\mathbf{X}^T\mathbf{A}$. Moreover, \mathbf{X} is a unit vector. 1-means is therefore nothing but a constrained low-rank approximation problem.

For k -means we have a similar picture. Consider the n by k clustering matrix \mathbf{X} defined as

$$X_i = \begin{cases} \frac{1}{\sqrt{|C_j|}} & \text{if point } A_i \text{ is in cluster } C_j \\ 0 & \text{otherwise} \end{cases}.$$

The columns of \mathbf{X} are orthogonal, that is, the j th column \mathbf{X}^j satisfies $\|\mathbf{X}^j\|_2 = 1$ and any column \mathbf{X}^i has a zero entry whenever a column X^j has a non-zero entry. Notice how $\mathbf{X}^i(\mathbf{X}^i)^T\mathbf{A}$ is mapped to the centroid of the cluster C_i . k -means can be therefore viewed as

$$\min_{\text{rank } k \text{ clustering matrix } \mathbf{X}} \|\mathbf{A} - \mathbf{X}\mathbf{X}^T\mathbf{A}\|_F^2.$$

By lifting the constraint that \mathbf{X} need be a clustering matrix, we are back to solving the low-rank approximation. Hence if we only cluster in the best k -dimensional subspace instead of the original d -dimensional space, we preserve most of the cost. This is made formal in the following theorem.

Theorem 7. *Let k be an integer and $A \in \mathbb{R}^{n \times d}$. Suppose we have an algorithm Alg that computes an α -approximation. Let $\mathbf{A}_k = \mathbf{U}\mathbf{A}\mathbf{V}_k^T$ be the best rank- k approximation. Then running Alg on \mathbf{A}_k yields a $\alpha + 1$ approximation.*

Proof. Let \mathbf{X} be the optimal clustering matrix. We observe that the optimal k -means cost $\|\mathbf{A} - \mathbf{X}\mathbf{X}^T\mathbf{A}\|_F^2$ is lower bounded by $\|\mathbf{A} - \mathbf{A}_k\|_F^2$. Let \mathbf{Y} be the clustering matrix obtained by Alg. Then

$$\begin{aligned} \|\mathbf{A} - \mathbf{Y}\mathbf{Y}^T\mathbf{A}\|_F^2 &\leq \|\mathbf{A} - \mathbf{Y}\mathbf{Y}^T\mathbf{A}_k\|_F^2 \leq \|\mathbf{A}_k - \mathbf{Y}\mathbf{Y}^T\mathbf{A}_k\|_F^2 + \|\mathbf{A} - \mathbf{A}_k\|_F^2 \\ &\leq \alpha \cdot \|\mathbf{A}_k - \mathbf{X}\mathbf{X}^T\mathbf{A}_k\|_F^2 + \|\mathbf{A} - \mathbf{X}\mathbf{X}^T\mathbf{A}\|_F^2 \\ &\leq \alpha \cdot \|\mathbf{A} - \mathbf{X}\mathbf{X}^T\mathbf{A}\|_F^2 + \|\mathbf{A} - \mathbf{X}\mathbf{X}^T\mathbf{A}\|_F^2 \leq (\alpha + 1) \|\mathbf{A} - \mathbf{X}\mathbf{X}^T\mathbf{A}\|_F^2. \end{aligned}$$

□

We remark that using \mathbf{A}_m instead of \mathbf{A}_k for $m > k/\varepsilon$, we obtain an $(\alpha + \varepsilon)$ approximation.