

Fast KDAC

April 4, 2017

Given the objective function :

$$\begin{aligned}
 \min \quad & -\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{tr}(W^T A_{i,j} W)}{2\sigma^2}} \\
 \text{W} \quad & \\
 \text{s.t.} \quad & W^T W = I \\
 & W \in \mathbb{R}^{d \times q} \\
 & A \in \mathbb{R}^{d \times d} \\
 & \gamma_{i,j} \in \mathbb{R}
 \end{aligned} \tag{1}$$

To optimize this cost function, the original KDAC uses an optimization technique called Dimensional Growth (DG). It rewrites the cost function into separate columns of the W matrix, and solve the problem one column at a time in a Greedy fashion.

$$\begin{aligned}
 \min \quad & -\sum_{i,j} \gamma_{i,j} e^{-\frac{w_1^T A_{i,j} w_1}{2\sigma^2}} e^{-\frac{w_2^T A_{i,j} w_2}{2\sigma^2}} \dots e^{-\frac{w_q^T A_{i,j} w_q}{2\sigma^2}} \quad w_i = i \text{ th column of } W \\
 \text{W} \quad & \\
 \text{s.t.} \quad & W^T W = I
 \end{aligned}$$

For example, to solve the first column w_1 , it ignores the rest of the columns and simplify the problem into :

$$\begin{aligned}
 \min \quad & -\sum_{i,j} \gamma_{i,j} e^{-\frac{w_1^T A_{i,j} w_1}{2\sigma^2}} \\
 w_1 \quad & \\
 \text{s.t.} \quad & w_1^T w_1 = 1
 \end{aligned}$$

This problem could be solved using standard Gradient methods. Once w_1 is computed, DG then treats the exponential term as a constant $g(w_1)$ and solve for the next variable.

$$\begin{aligned}
 f(w_2) &= -\sum_{i,j} \gamma_{i,j} e^{-\frac{w_1^T A_{i,j} w_1}{2\sigma^2}} e^{-\frac{w_2^T A_{i,j} w_2}{2\sigma^2}} \\
 f(w_2) &= -\sum_{i,j} \gamma_{i,j} g(w_1) e^{-\frac{w_2^T A_{i,j} w_2}{2\sigma^2}}
 \end{aligned}$$

Without the orthogonality constraint, each stage of the optimization process could be solved using Gradient methods. However, the orthogonality constraint of $W^T W = I$ requires further complication to ensure compliance. To start, the initialization of each new column w_i must go through the Gram Schmit method to ensure its orthogonality against all previous columns. Further more, the gradient direction calculated during each iterations also must undergo Gram Schmit. By removing components from previous vectors, it ensures that each update of w_i maintains feasibility.

Dimension Growth was the original approach used to solve the optimization problem of equation (1), and it achieved its objective of demonstrating the viability of KDAC. However, as the technology approach its next developmental stage, the implementation of this technology on large scale data requires KDAC to adapt for a more implementable algorithm.

The complexity of Dimension Growth algorithm heavily increases time of code development. This issue is especially prominent when speed requirement forces the development to be done in C or on the GPU. In these cases, simpler algorithm using off the shelf techniques could significantly reduce the developmental time and therefore the cost of its implementation.

The convergence speed of Dimension Growth in KDAC is slow. As we know from optimization theory, the convergence rate for gradient methods heavily depend on the conditional value of the Hessian matrix. The conditional value is defined as the ratio between the maximum and the minimum eigenvalue of the Hessian matrix.

$$\text{condition value} = \frac{\text{eig}_{\max}(\nabla^2 f(x))}{\text{eig}_{\min}(\nabla^2 f(x))}$$

The ideal condition value is when the ratio is equal to 1 and the convergence rate slows down very quickly as we increase the condition value beyond 10. Given this fact, it would be instructive to study the Hessian matrix to potentially explain the slow convergence of gradient methods. The Hessian for each column r has the following form :

$$\nabla^2 f(w) = \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{w^T A_{i,j} w_r}{2\sigma^2}} \left[A_{i,j} - \frac{1}{\sigma^2} A_{i,j} w_r w_r^T A_{i,j} \right]$$

From the Hessian matrix, we see that the condition value depends on the summation of matrix $A_{i,j}$ and $A_{i,j} w$ multiplied by some constant term, $\frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{w^T A_{i,j} w_r}{2\sigma^2}}$. From this form, it is intuitive to find clues to size of the conditional value from the individual behaviors of $A_{i,j}$ and $A_{i,j} w$.

Due to the complexity of the Hessian form, it may be sufficiently instructive to simply look for an approximation of the Hessian to study its conditional value behavior. Given the problem :

$$f(w) = -\sum_{i,j} \gamma_{i,j} e^{-\frac{w^T A_{i,j} w}{2\sigma^2}} \quad (2)$$

We could approximate the equation (2), by using the Taylor Expansion around 0. Since gradient methods are techniques using the 1st order approximation, we could similarly make a 1st order approximate of the original function.

$$f(w) \approx -\sum_{i,j} \gamma_{i,j} \left(1 - \frac{w^T A_{i,j} w}{2\sigma^2} \right) \quad (3)$$

At this point, we find the approximate Hessian by taking the 2nd derivative.

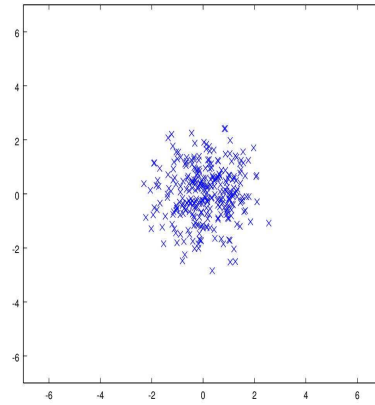
$$\begin{aligned} \nabla f(w) &\approx \sum \frac{\gamma_{i,j}}{\sigma^2} A_{i,j} w \\ \nabla^2 f(w) &\approx \sum \frac{\gamma_{i,j}}{\sigma^2} A_{i,j} \end{aligned}$$

From this form, we further simplified the Hessian matrix. And the simplified Hessian suggests a dominant influence of the summation of the $A_{i,j}$ matrix with some constant value $\gamma_{i,j}/\sigma^2$.

At this point, let's take a step back and ask how the eigenvalues of $A_{i,j}$ and $A_{i,j} w_r$ influence the conditional value depending on the type of data we handle. We first note that the $A_{i,j}$ matrix is formed with the following equation.

$$A_{i,j} = (x_i - x_j)(x_i - x_j)^T$$

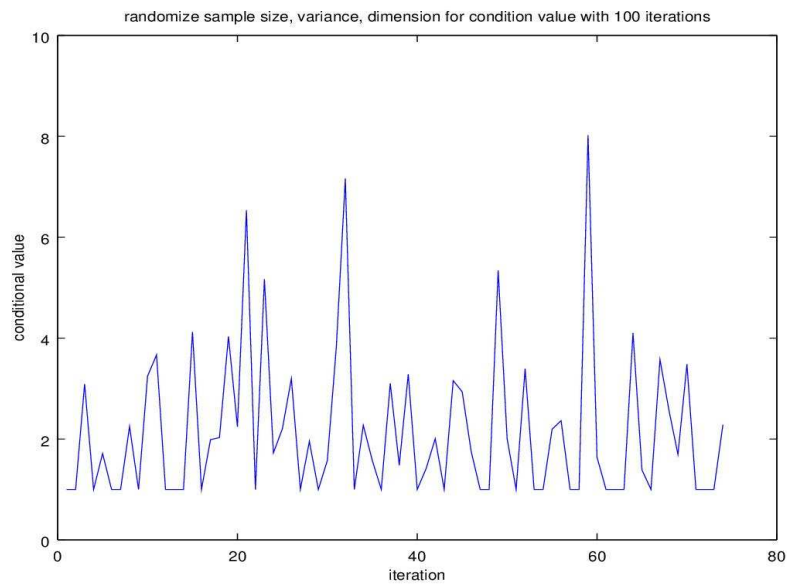
Given a single gaussian cluster of data :

**Figure 1.**

If we calculate

$$A = \sum A_{i,j}$$

If we randomly generate 100 similar distributions with randomize mean, variance, sample size and dimension, the condition value stays within a small bounded range.

**Figure 2.**

This plot is generated with the following code :

```
cv = []
for k = 1:100
    n = 20; %floor(40*rand());
    d = floor(5*rand());
    A = zeros(d,d);
    p = floor(5*rand())*randn(n,d);
```

```

for m = 1:n
    for n = 1:n
        v = p(m,:) - p(n,:);
        A = A + v'*v;
    end
end

[U,S,V] = svd(A);
D = diag(S) + 1;

%max(D)/min(D)
cv = [cv max(D)/min(D)];

end

plot(cv)
xlabel('iteration')
ylabel('conditional value')
title('randomize sample size, variance, dimension for condition value with 100
iterations')

```

From the plot of the conditional value, we can conclude that the conditional value stays bounded regardless of the sample size, variance, mean and dimensionality. This is for the case of a single cohesive cluster. However, what would happen if there are multiple clusters? As we vary the number of clusters as well as their distance apart, a clear pattern emerges. The following plot shows how the conditional value of 2 and 3 gaussian distributions. As we move them further apart from each other, the condition value explodes very quickly.

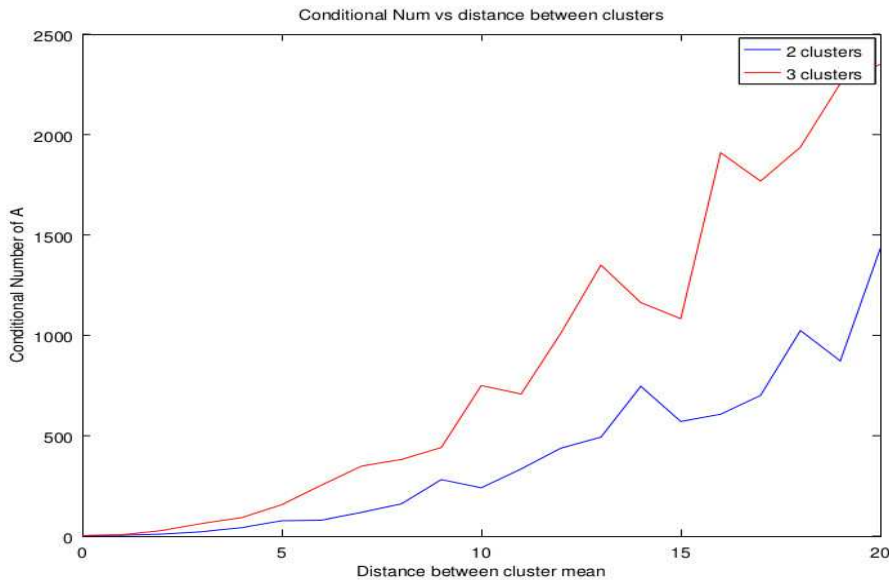


Figure 3.

In conjunction to the plot of a single cluster, this plot suggests that the conditional value is not bounded when the data type are separate into different clusters. As a matter of fact, the more clusters are involved the faster the conditional value grows. The purpose of showing ill-conditionality of this problem provide a reasoning to avoid Gradient method as an optimization approach. This allows us to lead into using a different optimization all together.

1 Fast KDAC

FKDAC is an alternative approach to Dimension Growth that estimates the optimal result without using gradient methods. We restate the optimization problem here :

$$\begin{aligned}
\min \quad & -\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} \\
s.t \quad & W^T W = I \\
& W \in \mathbb{R}^{d \times q} \\
& A \in \mathbb{R}^{d \times d} \\
& \gamma_{i,j} \in \mathbb{R}
\end{aligned} \tag{4}$$

1.1 Optimality Conditions

Due to the non-convex nature of the problem, FKDAC cannot guarantee to discover the global minimum. However, a sufficient optimality condition exists for a local minimum. Given the problem :

$$\begin{aligned}
\min \quad & f(x) \\
& x \\
s.t \quad & h(x) = 0
\end{aligned}$$

There exists a sufficiency condition according to Bertsekas in [1], Proposition 3.2.1. It states that x^* is a optimal solution of it satisfies the following proposition.

Proposition 1. (*Second Order Sufficiency Conditions*)

Assume that f and h are twice continuously differentiable, and let $x^* \in \mathbb{R}^n$ and $\lambda^* \in \mathbb{R}^m$ satisfy the following 3 conditions :

$$\nabla_x \mathcal{L}(x^*, \lambda^*) = 0 \quad \text{Condition 1}$$

$$\nabla_\lambda \mathcal{L}(x^*, \lambda^*) = 0 \quad \text{Condition 2}$$

$$z^T \nabla_{xx}^2 \mathcal{L}(x^*, \lambda^*) z > 0 \quad \text{for all } z \neq 0 \text{ with } \nabla h(x^*)^T z = 0 \quad \text{Condition 3}$$

Then x^* is a strict local minimum of f subjected to $h(x) = 0$.

Notice that this is a very general sufficiency assumption that does not assume the form of the functions f or h . Using this requirement as a proof of optimality, we define how our cost function can satisfy the definition. We start by finding the gradient of the Lagrangian with respect to w .

$$\begin{aligned}
\mathcal{L} &= -\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} - \frac{1}{2} \text{Tr}[\Lambda (W^T W - I)] \\
\frac{\partial \mathcal{L}}{\partial W} &= \left[\sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} A_{i,j} W \right] - W \Lambda = 0 \\
\frac{\partial \mathcal{L}}{\partial W} &= \left[\sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} A_{i,j} \right] W = W \Lambda
\end{aligned}$$

If we let :

$$\Phi(W) = \left[\sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} A_{i,j} \right]$$

We can rewrite the equation :

$$\Phi(W^*)W^* = W^* \Lambda \quad (5)$$

The equation above implies that the gradient of the Lagrangian is equal to 0 when W^* have columns of eigenvectors of $\Phi(w^*)$. Meeting this condition will satisfy the 1st condition of Proposition 1.

The 2nd condition of Proposition 1 requires that : $\nabla_{\Lambda} \mathcal{L} = 0$. From the cost function, we get :

$$\frac{\partial \mathcal{L}}{\partial \Lambda} = 1 - W^T W = 0 \quad (6)$$

With the conclusions from equation (5), it is obvious that if the optimal W^* have columns as unit norm eigenvectors, this condition is automatically satisfied.

To satisfy the final condition, we first define :

$$\mathcal{D}f(X)[Z] := \lim_{t \rightarrow 0} \frac{f(X + tZ) - f(X)}{t}$$

We want to define a set of Z such that :

$$\nabla h(x^*)^T Z = 0$$

We know that $h(W) = W^T W - 1$, therefore :

$$\mathcal{D}h(W)[Z] = \lim_{t \rightarrow 0} \frac{(W + tZ)^T (W + tZ) - W^T W}{t} = Z^T W + W^T Z$$

From this, we define as set of Z such that :

$$S_Z = \{Z : Z^T W + W^T Z = 0\}$$

The 3rd optimality condition is therefore met when :

$$\text{Tr}(Z^T \mathcal{D}f(W)[Z]) - \text{Tr}(\Lambda Z^T Z) > 0 \quad \forall Z \in S_Z$$

In Simpler terms, if the 2nd order gradient of the Lagrangian is Positive Definite, then the condition is satisfied. Due to the complexity of computing the Hessian in real time, a relaxation approach is discussed in a later section.

1.2 1st Order Relaxation

Given that $\Phi(W)$ is a function of W , it is difficult to take advantage of the eigenvalue/eigenvector relationship to find the solution.

$$\Phi(W^*)W^* = W^* \Lambda$$

A standard approach to circumvent this restriction is to treat W^* inside Φ as a separate variable and solve the problem iteratively.

$$\Phi(W_{k-1})W = W \Lambda$$

The key is, therefore, finding a reasonable initialization value for W_0 . This section provide a theoretical suggestion for a logical initialization point. Given the following problem :

$$\begin{aligned} \min_W \quad & - \sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} \\ \text{s.t} \quad & W^T W = I \end{aligned}$$

Writing out the Larangian :

$$\mathcal{L} = - \sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} + \frac{1}{2} \text{Tr}(\Lambda (I - W^T W))$$

We simplify the Lagrangian by using the Taylor approximation on the problematic exponential term. The Taylor approximation is expanded up to the 1st order centering around 0.

$$\mathcal{L} \approx - \sum_{i,j} \gamma_{i,j} \left(1 - \frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2} \right) + \frac{1}{2} \text{Tr}(\Lambda (I - W^T W))$$

Following the similar procedure by finding the derivative of \mathcal{L} and setting it to zero.

$$\nabla \mathcal{L} \approx \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} A_{i,j} W - W \Lambda = 0$$

We arrange the problem into standard eigenvalue/eigenvector problem.

$$\left[\sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} A_{i,j} \right] W = W \Lambda$$

From this, we see that Φ is no longer a function of W . If $W \in \mathbb{R}^{d \times q}$, the W_0 is therefore, q eigenvectors of the matrix Φ_0 :

$$\Phi_0 = \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} A_{i,j}$$

And Λ is the eigenvalue matrix.

$$\Lambda = \text{eig} \left[\sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} A_{i,j} \right]$$

Using the W_0 discovered from this initialization point, we could similarly form an updated Φ_{n+1} using the original cost derivative without any approximation.

$$\Phi_{n+1} = \left[\sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W_n^T A_{i,j} W_n)}{2\sigma^2}} A_{i,j} \right]$$

Given that only q eigenvectors are chosen out of d total eigenvectors, and $q \ll d$. The total possible combinations could be extremely large depending on the values of d and q , more specifically it is equal to $\binom{d}{q}$. Solving this problem becomes a difficult and separate combinatorial optimization.

One reasonable approach is to use the greedy algorithm to find q eigenvectors that produces the lowest cost. However, as we explore the theoretic motivation for choosing the optimal eigenvectors, we demonstrate the possibility of knowing the optimal eigenvectors without resorting to greedy methods.

1.3 Convergence

The convergence of the sequence generated from using the 1st order relaxation may not be possible, however through the Bolzano-Weierstrass Theorem, a convergent subsequence can be shown. According to [1], it states that :

Theorem 2. *A bounded sequence of real numbers has convergent subsequence.*

According to this theorem, if we can show that the sequences generated from the 1st order relaxation is bounded, it has a convergent subsequence. If we study the equation more closely :

$$\Phi(W_{k-1}) W_k = W_k \Lambda$$

The key driver of the sequence of W_k is the matrix Φ , therefore, if we can show that if this matrix is bounded, the sequence itself is also bounded. We look inside the construction of the matrix itself.

$$\Phi_{n+1} = \left[\sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W_n^T A_{i,j} W_n)}{2\sigma^2}} A_{i,j} \right]$$

From this, equation, start with the matrix $A_{i,j} = (x_i - x_j)(x_i - x_j)^T$. Since x_i, x_j are data points that are always centered and scaled to a variance of 1, the size of this matrix is always constrained.

It also implies that $A_{i,j}$ is a PSD matrix. From this, the exponential term $e^{-\frac{\text{Tr}(W_n^T A_{i,j} W_n)}{2\sigma^2}}$ is always limited between the value of 0 and 1. The value of σ is a constant given from the initialization stage. Lastly, we have the $\gamma_{i,j}$ term. Looking at it closely, it is equal to :

$$\gamma_{i,j} = \frac{u_i^T u_j}{d_i d_j} - \lambda (H Y Y^T H)_{i,j}$$

The vectors u_i will always be a vectors with a norm of 1, d_i is always less than \sqrt{N} , where N is the number of samples. λ is a constant given from the initialization stage. H and Y can be considered as constants as well. From this, it is clear that the infinity norm of the Φ is always bounded. The eigenvalue matrix of Λ is therefore also bounded. Using the Bolzano-Weierstrass Theorem, the sequence contains a convergent subsequence.

1.4 2nd Order Condition

As mentioned in the previous section, knowing that picking the eigenvectors of $\Phi(W^*)$ satisfies the 1st order condition is not sufficient. Ideally, we want to pick the optimal eigenvectors. The purpose of this section is to provide the 2nd order condition and prove the following theorem.

Theorem 3. *Given a large enough σ , picking the smallest q eigenvectors of $\Phi(W^*)$ satisfies the 2nd order condition.*

Proof.

□

Given :

$$\begin{aligned}\mathcal{L} &= -\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} - \frac{1}{2} \text{Tr}(\Lambda(W^T W - I)) \\ \nabla \mathcal{L} &= \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} A_{i,j} W - W \Lambda \\ \mathcal{D}\nabla \mathcal{L}[Z] &= \lim_{t \rightarrow 0} \frac{\partial}{\partial t} \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}((W+tZ)^T A_{i,j} (W+tZ))}{2\sigma^2}} A_{i,j} (W+tZ) - (W+tZ) \Lambda\end{aligned}$$

Let's divide this problem into 3 terms.

$$\mathcal{D}\nabla \mathcal{L}[Z] = T_1 + T_2 - T_3$$

Let's now solve each term separately.

$$\begin{aligned}T_1 &= \lim_{t \rightarrow 0} \frac{\partial}{\partial t} \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}((W+tZ)^T A_{i,j} (W+tZ))}{2\sigma^2}} A_{i,j} W \\ T_1 &= \lim_{t \rightarrow 0} \frac{\partial}{\partial t} \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W + t Z^T A_{i,j} W + t W^T A_{i,j} Z + t^2 Z^T A_{i,j} Z)}{2\sigma^2}} A_{i,j} W \\ T_1 &= -\sum_{i,j} \frac{\gamma_{i,j}}{2\sigma^4} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} \text{Tr}(Z^T A_{i,j} W + W^T A_{i,j} Z) A_{i,j} W \\ T_2 &= \lim_{t \rightarrow 0} \frac{\partial}{\partial t} \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} t e^{-\frac{\text{Tr}((W+tZ)^T A_{i,j} (W+tZ))}{2\sigma^2}} A_{i,j} Z \\ T_2 &= \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} A_{i,j} Z \\ T_3 &= \lim_{t \rightarrow 0} \frac{\partial}{\partial t} (W+tZ) \Lambda \\ T_3 &= Z \Lambda\end{aligned}$$

Putting together all 3 terms :

$$\mathcal{D}\nabla \mathcal{L}[Z] = \left\{ \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} \left[A_{i,j} Z - \frac{1}{2\sigma^2} \text{Tr}(Z^T A_{i,j} W + W^T A_{i,j} Z) A_{i,j} W \right] \right\} - Z \Lambda$$

Since $A_{i,j}$ is symmetric, we know that $\text{Tr}(Z^T A_{i,j} W) = \text{Tr}(W^T A_{i,j} Z)$, we can combine the two terms.

$$\mathcal{D}\nabla\mathcal{L}[Z] = \left\{ \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}((W^T A_{i,j} W))}{2\sigma^2}} \left[A_{i,j} Z - \frac{1}{\sigma^2} \text{Tr}(Z^T A_{i,j} W) A_{i,j} W \right] \right\} - Z\Lambda$$

We know that the 2nd order condition is satisfied if :

$$\text{Tr}(Z^T \mathcal{D}\nabla\mathcal{L}[Z]) > 0 \quad (7)$$

For all Z satisfying the following condition.

$$Z^T W + W^T Z = 0 \quad (8)$$

Therefore :

$$\text{Tr}(Z^T \mathcal{D}\nabla\mathcal{L}[Z]) = \left\{ \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}((W^T A_{i,j} W))}{2\sigma^2}} \left[\text{Tr}(Z^T A_{i,j} Z) - \frac{1}{\sigma^2} \text{Tr}(Z^T A_{i,j} W)^2 \right] \right\} - \text{Tr}(Z^T Z \Lambda_W)$$

Or it is easier to read them in 3 parts.

$$\text{Tr}(Z^T \mathcal{D}\nabla\mathcal{L}[Z]) = \text{Part 1} + \text{Part 2} + \text{Part 3}$$

$$\text{Part 1} = \text{Tr} \left(Z^T \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}((W^T A_{i,j} W))}{2\sigma^2}} A_{i,j} Z \right)$$

$$\text{Part 2} = - \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^4} e^{-\frac{\text{Tr}((W^T A_{i,j} W))}{2\sigma^2}} \text{Tr}(Z^T A_{i,j} W)^2$$

$$\text{Part 3} = - \text{Tr}(Z^T Z \Lambda_W)$$

Because Z has the same dimension as W , each column of Z resides in the space of \mathbb{R}^d . Since the eigenvectors of $\Phi(W)$ spans the \mathbb{R}^d , each column of Z can be represented as a linear combination of the eigenvalues. For the sake of clarity, we reiterate that W as the q eigenvectors chosen during each iteration, and \bar{W} as the $d - q$ eigenvectors not chosen. Each column of Z can therefore be written as :

$$Z = [z_1 \ z_2 \ \dots \ z_q]$$

$$z_i = \sum_{\alpha}^{d-q} \rho_{\bar{W}_\alpha}^{(i)} \bar{W}_\alpha + \sum_{\beta}^q \rho_{W_\beta}^{(i)} W_\beta$$

$$z_i = \bar{W} P_{\bar{W}}^{(i)} + W P_W^{(i)}$$

\bar{W}_α denotes the α th column of \bar{W} matrix, and W_β denotes the β th column of the W matrix. $\rho_{\bar{W}_\alpha}^{(i)}$ denotes the coefficient for the vector \bar{W}_α for the z_i vector. $\rho_{W_\beta}^{(i)}$ denotes the coefficient for the vector W_β for the z_i vector. To further simplify the notation, we eliminated the summation and turned $\rho_{\bar{W}_\alpha}^{(i)}$ into a vector $P_{\bar{W}}^{(i)}$.

Let plug this definition of Z into each part.

Part 1

$$\begin{aligned}
& \text{Tr} \left(Z^T \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^2} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} A_{i,j} Z \right) \\
& \quad \text{Tr}(Z^T \Phi(W) Z) \\
& = \text{Tr}(Z^T \Phi(W) [z_1 \ z_2 \ \dots \ z_q]) \\
& = \text{Tr} \left(Z^T \Phi(W) \begin{bmatrix} (\bar{W} P_W^{(1)} + W P_W^{(1)}) & (\bar{W} P_W^{(2)} + W P_W^{(2)}) & \dots & (\bar{W} P_W^{(q)} + W P_W^{(q)}) \end{bmatrix} \right) \\
& = \text{Tr} \left(Z^T \begin{bmatrix} (\Phi(W) \bar{W} P_W^{(1)} + \Phi(W) W P_W^{(1)}) & \dots & (\Phi(W) \bar{W} P_W^{(q)} + \Phi(W) W P_W^{(q)}) \end{bmatrix} \right) \\
& = \text{Tr} \left(Z^T \begin{bmatrix} (\bar{W} \bar{\Lambda} P_W^{(1)} + W \Lambda P_W^{(1)}) & \dots & (\bar{W} \bar{\Lambda} P_W^{(q)} + W \Lambda P_W^{(q)}) \end{bmatrix} \right) \\
& = \text{Tr} \left(\begin{bmatrix} (\bar{W} P_W^{(1)} + W P_W^{(1)}) & \dots & (\bar{W} P_W^{(q)} + W P_W^{(q)}) \end{bmatrix}^T \dots \right. \\
& \quad \left. \begin{bmatrix} (\bar{W} \bar{\Lambda} P_W^{(1)} + W \Lambda P_W^{(1)}) & \dots & (\bar{W} \bar{\Lambda} P_W^{(q)} + W \Lambda P_W^{(q)}) \end{bmatrix} \right) \\
& = \text{Tr} \left(\begin{bmatrix} (P_W^{(1)T} \bar{W}^T + P_W^{(1)T} W^T) \\ \dots \\ (P_W^{(q)T} \bar{W}^T + P_W^{(q)T} W^T) \end{bmatrix} \dots \right. \\
& \quad \left. \begin{bmatrix} (\bar{W} \bar{\Lambda} P_W^{(1)} + W \Lambda P_W^{(1)}) & \dots & (\bar{W} \bar{\Lambda} P_W^{(q)} + W \Lambda P_W^{(q)}) \end{bmatrix} \right) \\
& = (P_W^{(1)T} \bar{\Lambda} P_W^{(1)} + P_W^{(2)T} \bar{\Lambda} P_W^{(2)} + \dots + P_W^{(q)T} \bar{\Lambda} P_W^{(q)}) + (P_W^{(1)T} \Lambda P_W^{(1)} + P_W^{(2)T} \Lambda P_W^{(2)} + \dots + P_W^{(q)T} \Lambda P_W^{(q)}) \\
& = \text{Tr} \left(\begin{bmatrix} P_W^{(1)T} \\ P_W^{(2)T} \\ \dots \\ P_W^{(q)T} \end{bmatrix} \bar{\Lambda} \begin{bmatrix} P_W^{(1)} & P_W^{(2)} & \dots & P_W^{(q)} \end{bmatrix} \right) + \text{Tr} \left(\begin{bmatrix} P_W^{(1)T} \\ P_W^{(2)T} \\ \dots \\ P_W^{(q)T} \end{bmatrix} \Lambda \begin{bmatrix} P_W^{(1)} & P_W^{(2)} & \dots & P_W^{(q)} \end{bmatrix} \right) \\
& = \text{Tr}(P_W^T \bar{\Lambda} P_W) + \text{Tr}(P_W^T \Lambda P_W)
\end{aligned}$$

There's not much to simplify for Part 2, we leave it alone.

$$\text{Part 2} = - \sum_{i,j} \frac{\gamma_{i,j}}{\sigma^4} e^{-\frac{\text{Tr}(W^T A_{i,j} W)}{2\sigma^2}} \text{Tr}(Z^T A_{i,j} W)^2$$

Par 3 also have a Z term that could be simplified, using the same logic as Part 1, we can also simplify it.

$$\begin{aligned}\text{Part 3} &= -\text{Tr}(Z^T Z \Lambda_W) \\ \text{Part 3} &= -\text{Tr}(Z \Lambda_W Z^T) \\ &= \text{Tr}(P_W^T \Lambda P_W) + \text{Tr}(P_W^T \Lambda P_W)\end{aligned}$$

Together :

$$\text{Tr}(Z^T \mathcal{D} \nabla \mathcal{L}[Z]) = \text{Tr}(P_W^T \bar{\Lambda} P_W) + \text{Tr}(P_W^T \Lambda P_W) - \text{part 2} - \text{Tr}(P_W^T \Lambda P_W) - \text{Tr}(P_W^T \Lambda P_W)$$

We can cancel out the term $\text{Tr}(P_W^T \Lambda P_W)$ and combine the term $\text{Tr}(P_W^T \Lambda P_W)$.

$$\text{Tr}(Z^T \mathcal{D} \nabla \mathcal{L}[Z]) = \text{Tr}(P_W^T (\bar{\Lambda} - \Lambda) P_W) - \text{part 2}$$

Notice that we can shift the trace term :

$$\text{Tr}(Z^T \mathcal{D} \nabla \mathcal{L}[Z]) = \text{Tr}((\bar{\Lambda} - \Lambda) (P_W P_W^T)) - \text{part 2}$$

Putting it into this form allows us to look at the trace as an element wise product of the 2 matrices, $(\bar{\Lambda} - \Lambda)$ and $P_W P_W^T$. Since $(\bar{\Lambda} - \Lambda)$ is a diagonal matrix with all other elements equal to 0, only the diagonal terms of $P_W P_W^T$ will factor into element-wised sum. Since the diagonal of $P_W P_W^T$ will always be positive, if $(\bar{\Lambda} - \Lambda)$ is a positive diagonal matrix, the resulting trace will always be positive. The 2nd order condition is, therefore, satisfied, when :

$$\text{Tr}((\bar{\Lambda} - \Lambda) (P_W P_W^T)) > \text{part 2} \tag{9}$$

From the formula above, to satisfy the 2nd order condition, our goal is to make $\text{Tr}(Z^T \mathcal{D} \nabla \mathcal{L}[Z])$ as large as possible. It would make sense, therefore, to pick Λ_W to be as small as possible to maximize the left hand side of the equation. From this relationship, we conclude that if the ratio of the eigenvalues not chosen is sufficiently greater than the chosen eigenvalues, the 2nd order condition is satisfied. Further more, notice that part 2 has a $\frac{1}{\sigma^4}$ term. For a largen enough sigma, part 2 approach 0. Therefore, since the left hand side is always positive, picking the smallest eigenvectors automatically satisfies the 2nd order condition. This proof's theorem 3.

The inequality implies that if the gap between the eigenvalues chosen and the eigenvalues not chosen is larger than part 2, the 2nd order condition is satisfied. To characterize this gap, it is reasonable to measure the difference between the largest eigenvalue from Λ_W against the smallest eigenvalue from $\bar{\Lambda}$. We will call this difference the threshold eigengap. To further study the inequality, experiments were conducted to identify the threshold eigengap percent increase in each experiment.

The following table lists the threshold eigengap percent increase for multiple experiments. The percent increase is calculated by dividing the threshold gap by the largest eigenvalue of Λ_W .

Eigen Gap % increase

Experiment	Percent increase
Small 4 Gaussians	559,604%
Large 4 Gaussians	1,479%
Moon no noise	6,761%
Moon with noise	660%
Breast Cancer Data	211%
Facial image Data	113.8%

It is from this derivation, that each iterative process, the smallest eigenvectors are chosen as the new W matrix.

It is interesting to note that this formulation points to potential logical values for q . Locations where the eigengap percent increase is large have the highest chance of satisfying the 2nd order sufficiency condition.

2 Results

The following plot is a speed test generated by running various experiment using the different bench mark algorithms. Since all the bench mark algorithms are gradient based, they are prone to be stuck in a local minimum. To achieve better result, random initialization is required to search for a global optimal. For these experiemnts, 30 random initializations were required. The benefit of FKDAC is that since a theoretical initialization is provided, it achieves the same optimal result without multiple random initializations.

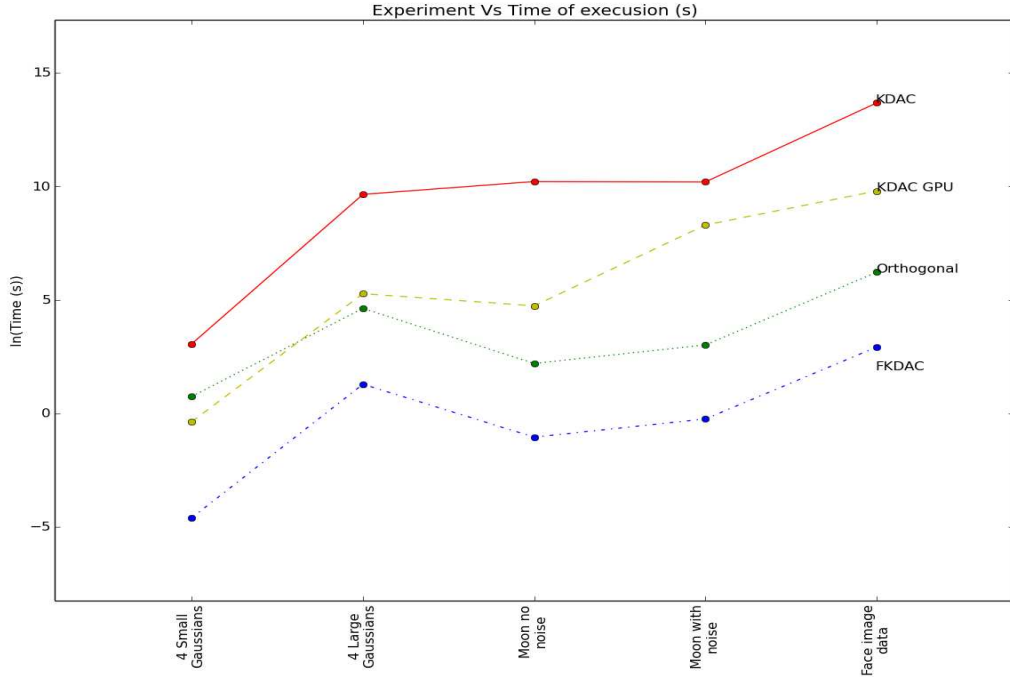


Figure 4.

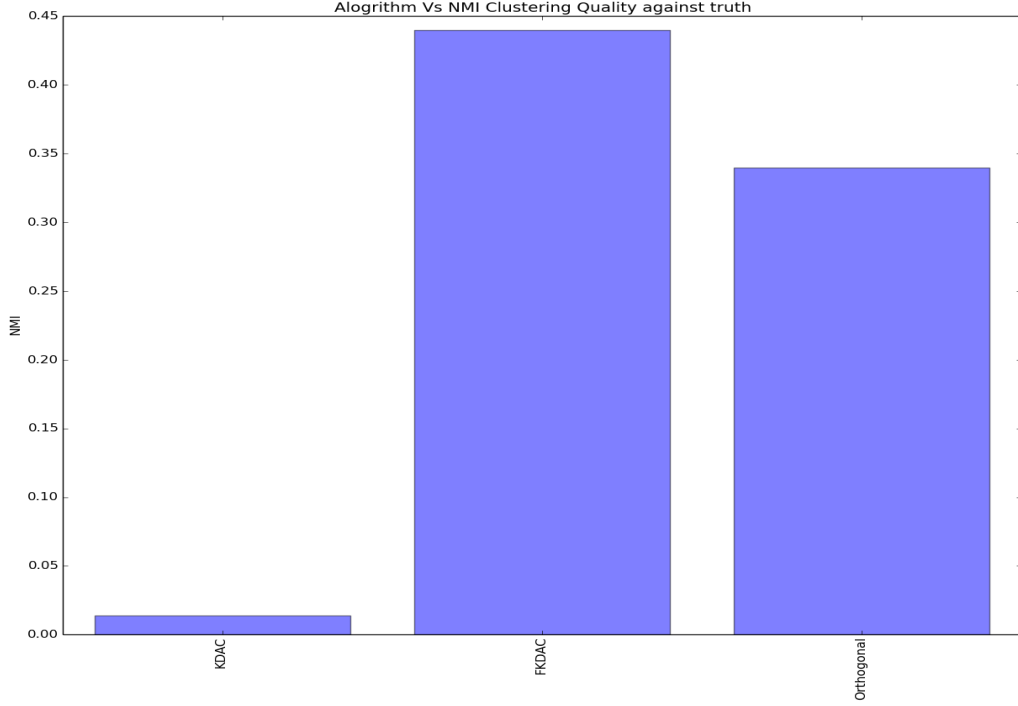


Figure 5.

3 Approximating w_{k+1}

Let q be the reduced the dimension and d as the original dimension. We have seen from the previous section that the stationary points could be found by picking q eigenvectors from the Φ marix. However, since q could be significantly smaller than d , using the appropriate eigenvectors could effect the outcome. Without any prior knowledge, the problems requires the optimal solution to choose from potenial large outcomes of $\binom{d}{q}$. To simplify the selection process, this section provides a theoretical argument on how these vectors could be chosen intelligently. To demonstrate the idea, we start by simplifying the problem into a single column w vector. This is without much loss of generality since the idea could be generalized into higher dimensions. Again, assuming that the Lipschitz condition holds, we start by taking the 1st order Taylor Expansion around some w_k .

$$f(w) = -\sum_{i,j} \gamma_{i,j} e^{-\frac{w_k^T A_{i,j} w_k}{2\sigma^2}} + 2 \left[\sum_{i,j} \gamma_{i,j} e^{-\frac{w_k^T A_{i,j} w_k}{2\sigma^2}} A_{i,j} w_k \right]^T (w_{k+1} - w_k) \quad (10)$$

The approximation generally assume that the constraint space is convex. For our case, since our constraint space is not convex, we must assume at least that the constraint space is convex within a certain radius.

$$w \in S \quad s.t \quad S \text{ is convex} \quad \forall \quad \|w - w_k\| \leq \varepsilon$$

Looking at the right side of equation 6, the first term is identical to the current cost at w_k . If we assume that higher order terms are insignificant, we can achieve a lower cost as long as we pick w_{k+1} such that the 2nd term is a negative value. Let's look at the 2nd term more closely.

$$2 \left[\sum_{i,j} \gamma_{i,j} e^{-\frac{w_k^T A_{i,j} w_k}{2\sigma^2}} A_{i,j} w_k \right]^T (w_{k+1} - w_k) \leq 0$$

Note that $A_{i,j}$ are symmetric. We can rewrite this expression :

$$w_k^T \left[\sum_{i,j} \gamma_{i,j} e^{-\frac{w_k^T A_{i,j} w_k}{2\sigma^2}} A_{i,j} \right] w_{k+1} \leq w_k^T \left[\sum_{i,j} \gamma_{i,j} e^{-\frac{w_k^T A_{i,j} w_k}{2\sigma^2}} A_{i,j} \right] w_k \quad (11)$$

Looking at the equation above, variously conclusions could be drawn. First, by realizing that the equation :

$$v^T = w_k^T \left[\sum_{i,j} \gamma_{i,j} e^{-\frac{w_k^T A_{i,j} w_k}{2\sigma^2}} A_{i,j} \right]$$

is a vector, we could rewrite equation (7) as :

$$v^T w_{k+1} \leq v^T w_k$$

To draw the first conclusion, we note that $v^T w_{k+1}$ is a product of two vectors. To ensure that $v^T w_{k+1}$ is always less than or equal to $v^T w_k$, we simply pick w_{k+1} that minimizes the term $v^T w_{k+1}$.

$$\min_{w_{k+1}} v^T w_{k+1}$$

We know from geometry that the product of two vectors is minimized when they are opposite directions of each other. Therefore w_{k+1} is minimized when $w_{k+1} = -v$. From this approach, we can iteratively pick w_{k+1} by setting $w_{k+1} = -v$ to achieve a descending direction. However, this approach encounters the problem of discovering the step length accompanying the descending direction. Since the constraint space is not convex, the step length cannot be discovered by methods of interpolation.

Although finding a definite descent direction that stays on the Stiefel manifold is difficult, another insight could be drawn by making some obvious observations of the following equation.

$$w_k^T \Phi w_{k+1} \leq w_k^T \Phi w_k \quad (12)$$

Observations :

1. We have performed Taylor expansion around w_k
2. Any w chosen must have a norm of 1
3. If we have chosen a proper w_{k+1} that satisfies the inequality, it forms a descent direction.
4. If we continue to follow the descent direction, we would hit w^* at convergence
5. At the point of convergence $w_{k+1} = w_k$
6. At the point of convergence :

$$w_{k+1}^T \Phi w_{k+1} = w_k^T \Phi w_k \quad (13)$$

and no vector v exists such that :

$$w^{*T} \Phi w^* < v^T \Phi v \quad (14)$$

Conclusion :

Given the point of expansion w_k , there is no need to iteratively approach convergence because we already know w^* . The only vector w^* that satisfies equation (10) and has a norm of 1, is the least dominant eigenvector of Φ . Therefore, the least dominant eigenvector of Φ is the best approximation of w^* at the Taylor expansion of w_k .

Linking this equation to FKDAC, we have previously shown that the eigenvectors of Φ provide a stationary point such that $\nabla \mathcal{L} = 0$. However, it was not certain which eigenvectors are most appropriate. From the conclusion we have just drawn, it shows that the least dominant eigenvectors are the most reasonable.

4 Expand the same idea to multiple columns

After understanding the basic concept through the usage of a single column example, we can now expand the same idea to a more general multiple column case. Given the first order Taylor expansion.

$$f(W) = f(W_k) + \nabla f(W_k)^T (W - W_k) \quad (15)$$

Applying this equation to our problem, we get :

$$f(W) = f(w) = -\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W_k^T A_{i,j} W_k)}{2\sigma^2}} + \left[\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W_k^T A_{i,j} W_k)}{2\sigma^2}} \text{Vec}(A_{i,j} W_k) \right]^T \text{Vec}(W - W_k)$$

Note that the notation $\text{Vec}(\cdot)$ is the vectorization of a matrix. Similar to the single column case, we only need to concentrate on making the 2nd term negative.

$$\left[\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{Tr}(W_k^T A_{i,j} W_k)}{2\sigma^2}} \text{Vec}(A_{i,j} W_k) \right]^T \text{Vec}(W - W_k) \leq 0$$

Let :

$$\beta_{i,j} = \gamma_{i,j} e^{-\frac{\text{Tr}(W_k^T A_{i,j} W_k)}{2\sigma^2}}$$

$$\left[\sum_{i,j} \beta_{i,j} (I \otimes A_{i,j}) \text{Vec}(W_k) \right]^T [\text{Vec}(W) - \text{Vec}(W_k)] \leq 0$$

We now rearrange the terms and get :

$$\left[\sum_{i,j} \beta_{i,j} (I \otimes A_{i,j}) \text{Vec}(W_k) \right]^T \text{Vec}(W) \leq \left[\sum_{i,j} \beta_{i,j} (I \otimes A_{i,j}) \text{Vec}(W_k) \right]^T \text{Vec}(W_k)$$

$$\text{Vec}(W_k)^T \left[\sum_{i,j} \beta_{i,j} (I \otimes A_{i,j}) \right] \text{Vec}(W) \leq \text{Vec}(W_k)^T \left[\sum_{i,j} \beta_{i,j} (I \otimes A_{i,j}) \right] \text{Vec}(W_k)$$

Using the same logic to lower the upper bound, the W_k that minimizes the upper bound consists of the q least dominant eigenvectors of the matrix Φ .

5 Optimality condition

The Taylor expansion approach corresponds directly with the optimality condition with a convex constrained space. Given a problem of :

$$\min_{x \in X} f(x)$$

We know that given x^* as a local minimum in a convex constrained space of X , the following condition must be satisfied.

$$\nabla f(x^*)^T (x - x^*) \geq 0 \quad \forall x \in X$$

Given that we have a non-convex constraint space, we must make more strict assumptions. Suppose x^* is a local minimum within a ball defined as $\mathcal{B}(x, \varepsilon) := \{x : \|x - x^*\| < \varepsilon\}$. Assume $\mathcal{B}(x, \varepsilon) \cap X$ is convex and f is convex within $\mathcal{B}(x, \varepsilon) \cap X$, then we state the following as the optimality condition.

$$\nabla f(x^*)^T (x - x^*) \geq 0 \quad \forall x \in X \quad \text{and} \quad \|x - x^*\| \leq \varepsilon$$

Given a function of :

$$f(w) = -\sum \gamma_{i,j} e^{-\frac{w^T A_{i,j} w}{2\sigma^2}}$$

We create an approximate of the function :

$$f(w) \approx -\sum \gamma_{i,j} \left(1 - \frac{1}{2\sigma^2} w^T A_{i,j} w \right)$$

From the approximation, we find the gradient :

$$\nabla f(w) \approx \sum \frac{\gamma_{i,j}}{\sigma^2} A_{i,j} w$$

Using the approximated gradient, we could now check for the optimality condition.

$$\begin{aligned} w^{*T} \left[\sum \gamma_{i,j} A_{i,j} \right]^T (w - w^*) &\geq 0 \\ w^{*T} \left[\sum \gamma_{i,j} A_{i,j} \right]^T w &\geq w^{*T} \left[\sum \gamma_{i,j} A_{i,j} \right]^T w^* \end{aligned}$$

From the equation above, we see that in order for w^* to be a local minimum, w^* must be chosen such that the left hand side is always larger than the right hand side for any w . This is only possible when w^* is the least dominant eigenvector of $\sum \gamma_{i,j} A_{i,j}$ matrix. From this perspective, we conclude that picking the least dominant eigenvector is a reasonable approximation for the local minimum of the original cost function.

6 Implementation Details of Cost function

The computation of cost function presented in this paper, is a complicated equation that slows down both the implementation and speed of the results.

$$\begin{aligned} \min_{W} \quad & -\sum_{i,j} \gamma_{i,j} e^{-\frac{\text{tr}(W^T A_{i,j} W)}{2\sigma^2}} \\ \text{s.t.} \quad & W^T W = I \\ & W \in \mathbb{R}^{d \times q} \\ & A \in \mathbb{R}^{d \times d} \\ & \gamma_{i,j} \in \mathbb{R} \end{aligned} \tag{16}$$

Instead of solving the function itself, it could be mostly easily done with the following equation.

$$\text{cost} = \text{HSIC}(XW, U) - \lambda \text{HSIC}(XW, Y)$$

The simplest way is to write a HSIC function, and pass XW, U , and Y to compute the final cost. Although easy, this approach is not the fastest in terms of separating out the portion of the code that requires constant update, and the portion that remains constant. In this section, a faster approach to implement the cost function is outlined.

Starting with the original cost function :

$$\text{cost} = \text{HSIC}(XW, U) - \lambda \text{HSIC}(XW, Y)$$

Convert it into trace format.

$$\text{cost} = \text{Tr}(\tilde{K} H U U^T H) - \lambda \text{Tr}(\tilde{K} H Y Y^T H)$$

Where \tilde{K} is the normalized kernel of XW , which could also be written as $\tilde{K} = D^{-\frac{1}{2}} K_{XW} D^{-\frac{1}{2}}$. Putting this into the cost function.

$$\text{cost} = \text{Tr}\left(D^{-\frac{1}{2}} K_{XW} D^{-\frac{1}{2}} H U U^T H\right) - \lambda \text{Tr}\left(D^{-\frac{1}{2}} K_{XW} D^{-\frac{1}{2}} H Y Y^T H\right)$$

When optimizing U , it is obvious that the 2nd portion does not effect the optimization. Therefore, U can be solved using the following form.

$$U = \underset{U}{\text{argmin}} \text{Tr}(U^T H D^{-1/2} K D^{-1/2} H U)$$

The situation get a bit more complicated if we are optimization for W . Using the combination of the rotation property and the combination of the 2 traces, the cost can be written as :

$$\text{cost} = \text{Tr}([D^{-1/2} H (U U^T - \lambda Y Y^T) H D^{-1/2}] K)$$

In this form, it can be seen that the update of W matrix will only affect the kernel K and the degree matrix D . Therefore, it makes sens to treat the middle portion as a constant which we refer as Ψ .

$$\text{cost} = \text{Tr}([D^{-1/2} \Psi D^{-1/2}] K)$$

Given that $[D^{-1/2} \Psi D^{-1/2}]$ is a symmetric matrix, from this form, we can convert the trace into an element wise product \odot .

$$\text{cost} = [D^{-1/2} \Psi D^{-1/2}] \odot K$$

To further reduction the amount of operation, we let d be a vector of the diagonal elements of $D^{-1/2}$, hence $d = \text{diag}(D^{-1/2})$, this equality hold.

$$D^{-1/2} \Psi D^{-1/2} = [d d^T] \odot \Psi$$

Therefore, the final cost function can be written in its simplest form as :

$$\text{cost} = \Gamma = \Psi \odot [d d^T] \odot K$$

During update, as W update during each iteration, the matrix Ψ stays as a constant while $d d^T$ and K update. The benefit of this form minimize the complexity of the equation, while simplify cost into easily parallelizable matrix multiplications. The equation also clearly separates the elements into portions that require an update and portions that does not.

7 Implementation Details of the Derivative

As it was shown from previous sections, the gradient of our cost function using the Gaussian Kernel has the following form.

$$\nabla f(W) = \left[\frac{1}{\sigma^2} \sum \gamma_{i,j} K_{i,j} A_{i,j} \right] W$$

It is often shown as :

$$\nabla f(W) = \Phi W$$

The key is therefore to find Φ .

$$\Phi = \frac{1}{\sigma^2} \sum \gamma_{i,j} K_{i,j} A_{i,j}$$

If we note that $A_{i,j} = (x_i - x_j)(x_i - x_j)^T$. It can be seen that the inner portion is identical to the cost function. The difference is the addition of the $A_{i,j}$ matrix and a constant of $\frac{1}{\sigma^2}$. These extra factors can be incorporate in the following form.

$$\Phi = \frac{1}{\sigma^2} Q^T \text{diag}(\text{Vec}(\Gamma)) Q$$

Where :

$$Q = (X \otimes \mathbf{1}_n) - (\mathbf{1}_n \otimes X)$$

Note that \otimes is a tensor product and $\mathbf{1}_n$ is a 1 vector with a length of n .

And :

$$\text{cost} = \Gamma = \Psi \odot [dd^T] \odot K$$

Since Q is a constant that never changes during the optimization, it could be calculated at the beginning and cached. During each W update using the gradient, Ψ and Q are considered as constants while K and D require a constant update. However, each time the U matrix is updated, Ψ must also be updated.

Here we outline the Algorithm for the W optimization update scheme.

1. Initialize $W_0 = 0$ for the first time, and W_k if U has been updated.
2. Calculate Q, Ψ and store them as constants
3. Calculate K, D, Φ
4. $W_{k+1} = \vec{\text{eig}}_{\min}(\Phi)$, pick q least dominant eigenvectors as W_{k+1} .
5. Repeat 3,4 until W convergence

[1] Bartle, Robert G.; Sherbert, Donald R. (2000). *Introduction to Real Analysis* (3rd ed.). New York: J. Wiley.