

DIMENSIONALITY REDUCTION AND THE FENCHEL GAME

LUCAS TUCKER

ABSTRACT. In this paper, we survey the dimensionality reduction algorithms Principal Component Analysis (PCA), Laplacian Eigenmaps (LE), and Isomap. We propose an improved Laplacian Eigenmaps algorithm along with a gradient-descent based PCA algorithm. Finally, we examine the novel Fenchel Game framework – a framework from which we may derive classical first-order methods such as gradient descent – and prove a bound for gradient descent with averaging. This paper is meant to serve as a reference for the theoretical underpinning of such popular dimensionality reduction techniques, as well as an attempt to clarify and build upon Wang-Abernathy-Levy’s paper.

CONTENTS

1. PCA	1
2. Experiment: PCA	3
Non-linear Methods	4
3. Laplacian Eigenmaps	4
4. Experiment: Laplacian Eigenmaps	6
5. Isomap	8
6. Fenchel Game No-Regret Dynamics (FGNRD)	10
Acknowledgments	17
7. Bibliography	17
References	17
Appendix A.	17
Appendix B.	18
Appendix C.	19

1. PCA

We first examine a commonly used linear method to reduce the dimensionality of a data matrix. This method, Principal Component Analysis (PCA), is most helpful for data with linearly related, highly correlated features. This discussion will lead to a comparison of a novel gradient descent-based PCA solution algorithm with the baseline approach outlined in Theorem 1.6.

Definition 1.1. A projection on a vector space V is a linear operator $P : V \rightarrow V$ such that $P^2 = P$. A projection on a Hilbert space V is an orthogonal projection if $\langle Px, y \rangle = \langle x, Py \rangle$.

Date: September 15, 2023.

Definition 1.2. The “Frobenius norm”, denoted by $\|\cdot\|_F$ is a matrix norm defined over $\mathbb{R}^{m \times n}$ as

$$\|\mathbf{M}\|_F := \sqrt{\sum_{i=1}^m \sum_{j=1}^n \mathbf{M}_{ij}^2}$$

Definition 1.3. For a sample $S = (x_1, \dots, x_m)$ and feature mapping $\Phi : \mathcal{X} \rightarrow \mathbb{R}^n$, we define the data matrix $(\Phi(x_1), \dots, \Phi(x_m)) =: \mathbf{X} \in \mathbb{R}^{n \times m}$. We deem \mathbf{X} a “centered data matrix” if $\sum_{i=1}^m \Phi(x_i) = \mathbf{0}$. Let \mathcal{P}_k denote the set of n -dimensional rank- k orthogonal projection matrices. PCA (Principal Component Analysis) is defined by the orthogonal projection matrix

$$(1.4) \quad \mathbf{P}^* := \operatorname{argmin}_{\mathbf{P} \in \mathcal{P}_k} \|\mathbf{P}\mathbf{X} - \mathbf{X}\|_F^2$$

Moreover, the sample covariance matrix corresponding to \mathbf{X} is given by $\frac{1}{m}\mathbf{X}\mathbf{X}^T$ since

$$\frac{1}{m}(\mathbf{X}\mathbf{X}^T)_{ij} = \frac{1}{m} \sum_{\ell=1}^m X_{i\ell} X_{\ell j}^T = \frac{1}{m} \sum_{\ell=1}^m \Phi(x_\ell)_i \Phi(x_\ell)_j$$

$$= E[\Phi(x)_i \Phi(x)_j] = E[\Phi(x)_i \Phi(x)_j] - E[\Phi(x)_i] E[\Phi(x)_j] = \operatorname{Cov}(\Phi(x)_i, \Phi(x)_j)$$

where the right hand term is the covariance between i -th and j -th coordinates of the feature output based on m samples.

Definition 1.5. The “top singular vector” of a matrix \mathbf{M} is the vector \mathbf{x} which maximizes the Rayleigh quotient

$$r(\mathbf{x}, \mathbf{M}) = \frac{\mathbf{x}^T \mathbf{M} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$$

Theorem 1.6. Let $\mathbf{P}^* \in \mathcal{P}_k$ be the PCA solution for a centered data matrix $\mathbf{X} = (\Phi(x_1), \dots, \Phi(x_m)) \in \mathbb{R}^{n \times m}$. Then, $\mathbf{P}^* = \mathbf{U}_k \mathbf{U}_k^T$, where $\mathbf{U}_k \in \mathbb{R}^{n \times k}$ is the matrix formed by the top k singular vectors of $\mathbf{C} := \frac{1}{m} \mathbf{X}\mathbf{X}^T$. Moreover, the associated k -dimensional representation of \mathbf{X} is given by $\mathbf{Y} = \mathbf{U}_k^T \mathbf{X}$.

Proof. See Appendix A.1 □

Remark 1.7. The top singular vectors of \mathbf{C} are the directions of maximal variance in the data, and the \mathbf{u}_i are the variances, so that PCA may be understood as projection onto the subspace of maximal variance. We wish to compare a gradient descent-based algorithm for solving \mathbf{P}^* from Definition 1.3 to the eigenvector approach outlined by Theorem 1.5, hence Theorem 1.7.

Theorem 1.8. The following holds for the PCA gradient:

$$\frac{\partial}{\partial \mathbf{U}_k} \|\mathbf{U}_k \mathbf{U}_k^T \mathbf{X} - \mathbf{X}\|_F^2 = 2(\mathbf{X}\mathbf{X}^T \mathbf{U} \mathbf{U}^T \mathbf{U} + \mathbf{U} \mathbf{U}^T \mathbf{X}\mathbf{X}^T \mathbf{U} - 2\mathbf{X}\mathbf{X}^T \mathbf{U})$$

Proof. See Appendix A.2 □

Remark 1.9. While

$$\operatorname{argmin}_{\mathbf{P} \in \mathbb{R}^{n \times n}} \|\mathbf{P}\mathbf{X} - \mathbf{X}\|_F^2$$

minimizes a convex function, the assumption $\mathbf{P} = \mathbf{U}_k \mathbf{U}_k^T$ leads to a non-convex minimization

$$\operatorname{argmin}_{\mathbf{U}_k \in \mathbb{R}^{n \times k}} \|\mathbf{U}_k \mathbf{U}_k^T \mathbf{X} - \mathbf{X}\|_F^2$$

Further, after finding \mathbf{U}_{\min} using gradient descent, we must find the closest orthogonal matrix $\mathbf{U} := \mathbf{A}\mathbf{B}^T$ for a Singular Value Decomposition $\mathbf{U}_{\min} = \mathbf{A}\mathbf{\Sigma}\mathbf{B}^T$. We now compare a gradient-descent based solution to Equation 1.4 with the baseline eigenvector approach outlined by Theorem 1.6.

2. EXPERIMENT: PCA

Define the t -similarity score to be the following: Given a dataset $\mathbf{X} \in \mathbb{R}^{n \times m}$ of m points in \mathbb{R}^n , and a point $\mathbf{x} \in \mathbf{X}$, let $n_t(\mathbf{x})$ be the t closest points $y \in X$ to \mathbf{x} under ℓ_2 -distance. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^k$ be the dimensionality reduction map from full dimension n to low-dimension k . Define the t -similarity score for $x \in X$ to be $\ell_t(\mathbf{x})$ given by

$$\ell_{t,f}(\mathbf{x}) := |n_t(\mathbf{x}) \cap n_t(f(\mathbf{x}))|$$

Let the t -similarity score of a dimensionality reduction technique f to be

$$\operatorname{score}_t(f) := \frac{1}{m} \cdot \sum_{\mathbf{x} \in X} \ell_{t,f}(\mathbf{x})$$

The experiment is as follows:

- Take $m = 2000$ points from MNIST ($n = 784$ pixels)
- Compute f_{random}^k , f_{pca}^k , and $f_{\text{gradient-pca}}^k$ given by random projection of \mathbb{R}^n into \mathbb{R}^k , baseline PCA implementation, and gradient-based PCA implementation.
- Plot the three curves: where the x -axis is the dimension k ranging (in logarithmic steps from $k = 2$ to $k = 784$, and the y -axis is given by $\operatorname{score}_t(f_{\text{random}}^k)$ (blue), $\operatorname{score}_t(f_{\text{pca}}^k)$ (red), and $\operatorname{score}_t(f_{\text{gradient-pca}}^k)$ (green) using $t = 10$. Note that the code for this experiment is available in the associated repository.

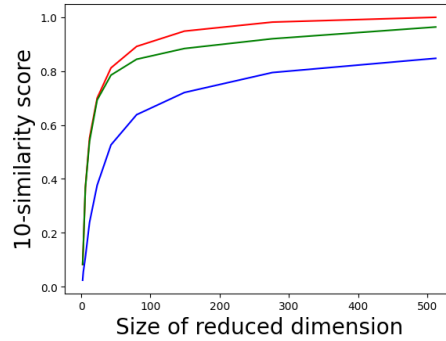


FIGURE 1. Comparison of similarity scores among the various PCA techniques as a function of dimension size

Size of reduced dimension (k-value)	Normal PCA	Random Projection	Gradient-based PCA
2	0.0825	0.0246	0.08315
3	0.14675	0.055	0.1451
6	0.36865	0.1097	0.3672
12	0.55215	0.23895	0.5395
23	0.6997	0.3765	0.69245
43	0.8116	0.5263	0.7848
80	0.8911	0.63805	0.84355
149	0.9478	0.72025	0.88335
276	0.98125	0.7943	0.91955
512	0.9991	0.847	0.9631

FIGURE 2. Values of the points in Figure 1, where the dimension size input is logarithmically spaced

Remark 2.1. As visualized above, the similarity scores (with respect to the data matrix X) of both the gradient-based PCA method and baseline implementation remained within 0.02 of one another up to a 40 dimensional representation, while the random projection’s similarity score remained on average within 0.25 from the baseline PCA implementation.

NON-LINEAR METHODS

Rather than linearly transforming our data matrix \mathbf{X} (through matrix multiplication), we may be interested in optimizing our representation with respect to a specific distance or weight matrix, or in preserving local structural information.

3. LAPLACIAN EIGENMAPS

Definition 3.1. The laplacian eigenmaps algorithm aims to find a k -dimensional representation of the data matrix \mathbf{X} which best preserves the weighted neighborhood relations specified by a matrix \mathbf{W} :

Algorithm 1 Laplacian Eigenmaps

Require: points $\mathbf{x} \in \mathbb{R}^n$

Require: scaling parameter σ

Require: t-nearest neighbors function $N_t(\mathbf{x})$

Define: $\mathbf{W} \in \mathbb{R}^{m \times m}$ as $W_{ij} := \begin{cases} 0 & \mathbf{x}_i \notin N_t(\mathbf{x}_j), \mathbf{x}_j \notin N_t(\mathbf{x}_i) \\ e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}} & \text{otherwise} \end{cases}$

Define: $\mathbf{D} \in \mathbb{R}^{m \times m}$ as $D_{ij} = \begin{cases} \sum_{s=1}^m W_{is} & j = i \\ 0 & j \neq i \end{cases}$

Evaluate: $\mathbf{Y} \in \mathbb{R}^{k \times m}$ as $\mathbf{Y} = \operatorname{argmin}_{\mathbf{Y}'} \left\{ \sum_{i,j} W_{ij} \|\mathbf{y}'_i - \mathbf{y}'_j\|_2^2 \right\}$

Intuitively, the above minimization penalizes k -dimensional representations of neighbors that differ largely under the ℓ_2 norm.

Proposition 3.2. *The solution to the laplacian eigenmaps minimization is $\mathbf{U}_{\mathbf{L},k}^T$, where $\mathbf{L} = \mathbf{D} - \mathbf{W}$ is the “graph Laplacian” and $\mathbf{U}_{\mathbf{L},k}^T$ are the bottom k singular vectors of \mathbf{L} (excluding 0 if the underlying neighborhood graph has connections).*

Proof. See Appendix B.1 □

Remark 3.3. Note that when the data is not uniformly distributed, cluster size remains the same. We may instead want to better learn the local manifold structure by increasing the number of neighbors measured for denser areas and reducing this number in sparser areas. We explore this idea in our novel proposed algorithm in Definition 3.5, which uses the notion of “average distance” defined by W. Jiang et al.’s laplacian eigenmaps algorithm defined below.

Definition 3.4. (W. Jiang et al.) Fix a sample size $t \in \mathbb{N}$ and t -nearest neighbors function $N_t(\mathbf{x}_i)$, where the $\{\mathbf{x}_i\}_{i=1}^m \in \mathbb{R}^n$ comprise the data matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$. Define the *average distance matrix* $A \in \mathbb{R}^{m \times m}$ as

$$A_i = \frac{\sqrt{\sum_{j=1}^t \|\mathbf{x}_i - \mathbf{x}_j\|_2^2}}{t}$$

where the $\mathbf{x}_j \in N_t(\mathbf{x}_i)$. Then, define the *distance matrix* $D \in \mathbb{R}^{m \times m}$ as

$$D_{ij} = \begin{cases} \frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2}{\sqrt{A_i A_j}}, & \text{if } \mathbf{x}_i \in N_t(\mathbf{x}_j) \text{ or } \mathbf{x}_j \in N_t(\mathbf{x}_i) \\ 0, & \text{else} \end{cases}$$

Finally, the “Averaged Laplacian Eigenmaps” weight matrix is defined as $W^{\text{ALE}} \in \mathbb{R}^{m \times m}$ as

$$W_{ij}^{\text{ALE}} = e^{-\frac{D_{ij}^2}{t}}$$

Note that the solution to the Laplacian Eigenmaps minimization problem given by Proposition 3.1 still holds in this case.

Definition 3.5. We now extend Algorithm 1 by including variable nearest neighbors for each point. In particular, rather than a fixed t -nearest neighbors we fix a “maximum” radius t_{\max} and define the variable radius $t(\mathbf{x}_i)$ as

$$A_{\text{avg}} = \frac{1}{m} \sum_{i=1}^m A_i$$

$$t(\mathbf{x}_i) = \min \left(\left\lfloor t_{\max} \frac{A_{\text{avg}}}{A_i} \right\rfloor, t_{\max} \right)$$

Hence, the radius is shortened in sparser areas and extended in denser areas, where we use the A_i values defined in Definition 3.4. We define this “Variable Radius Laplacian Eigenmaps” (VLE) by redefining our weight matrix as

$$(3.6) \quad W_{ij}^{\text{VLE}} := \begin{cases} e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}} & \mathbf{x}_j \in N_{t(\mathbf{x}_i)}(\mathbf{x}_i) \text{ or } \mathbf{x}_i \in N_{t(\mathbf{x}_j)}(\mathbf{x}_j) \\ 0 & \text{otherwise} \end{cases}$$

Note that the “or” condition is essential to preserve the symmetry of W . We may then solve the minimization problem by way of Proposition 3.1.

4. EXPERIMENT: LAPLACIAN EIGENMAPS

Remark 4.1. We now compare the the standard laplacian eigenmaps with both the “Averaged Laplacian Eigenmaps” (ALE) algorithm given by W. Jiang et al. and the “Variable Radius Laplacian Eigenmaps” algorithm we propose. We visually analyze the dimension reduction of a colored swiss roll (2d manifold) and an oversampled colored swiss roll. The following figures are available and can be recreated in the associated Python repository.

The figures below are labeled as follows:

- LE: Graph after applying baseline laplacian eigenmaps
- VLE: Graph after applying our proposed variable radius laplacian eigenmaps
- ALE: Graph after applying the algorithm given by W. Jiang et al.

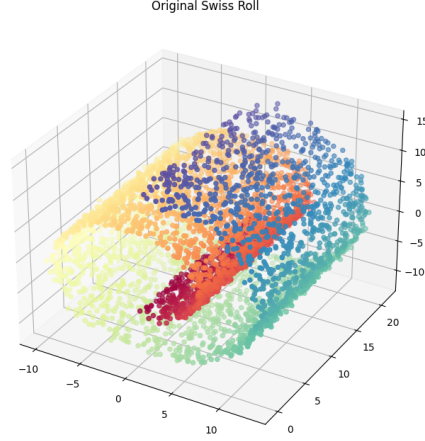


FIGURE 3. Original swiss roll with color gradient

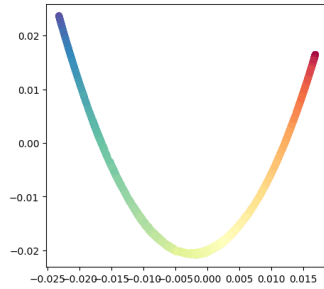


FIGURE 4. LE

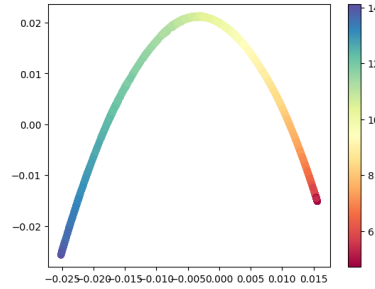


FIGURE 5. VLE

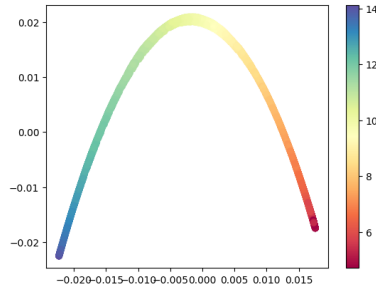


FIGURE 6. ALE

Remark 4.2. As depicted above, the three algorithms appear to unravel the roll correctly given uniformly sampled data. We now attempt to reduce the dimension of an oversampled swiss roll with the same algorithms.

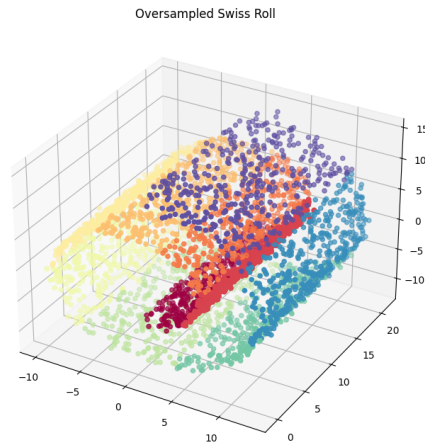


FIGURE 7. Oversampled swiss roll with color gradient

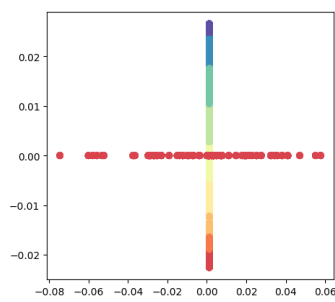


FIGURE 8. LE

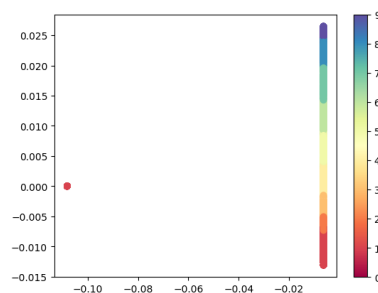


FIGURE 9. VLE

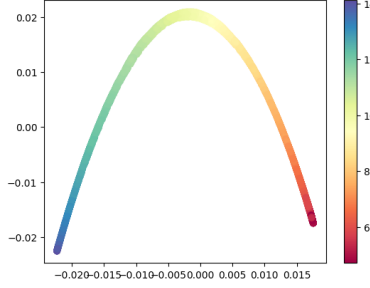


FIGURE 10. ALE

Remark 4.3. We find that our proposed algorithm (VLE) appears to adapt to the oversampling, as does the density-scaled algorithm (ALE) of W. Jiang et al., while the standard LE displays the oversampling. Note that the following hyperparameters were chosen (and can be adjusted): $t = 10$ neighbors for LE, $t_{\max} = 18$ for VLE, and $t = 20$ for ALE.

5. ISOMAP

Definition 5.1. Isomap extracts the low-dimensional data that best preserves pairwise distances between inputs based on their geodesic distances along a manifold. The algorithm is specified as follows:

Algorithm 2 Isomap

Require: Points $\mathbf{x}_i \in \mathbb{R}^n$

Evaluate: $N_t(\mathbf{x}_i) \forall \mathbf{x}_i$

Construct: Undirected neighborhood graph \mathcal{G}

Evaluate: Approximate Δ_{ij} as shortest distance in \mathcal{G}

Evaluate: $\mathbf{K}_{\text{Iso}} := -\frac{1}{2}(\mathbf{I}_m - \frac{1}{m}\mathbf{1}\mathbf{1}^T)\Delta(\mathbf{I}_m - \frac{1}{m}\mathbf{1}\mathbf{1}^T) := -\frac{1}{2}\mathbf{H}\Delta\mathbf{H}$

Evaluate: $\mathbf{Y} \in \mathbb{R}^{k \times m}$, $\mathbf{Y} := \operatorname{argmin}_{\mathbf{Y}' \in \mathbb{R}^{k \times m}} \sum_{i,j} (\|\mathbf{y}'_i - \mathbf{y}'_j\|_2 - \Delta_{ij})^2$ as $\mathbf{Y} = (\Sigma_{\text{Iso}, k})^{\frac{1}{2}} \mathbf{U}_{\text{Iso}, k}^T$ (see Proposition 5.1)

Remark 5.2. Note that $\mathbf{1} \in \mathbb{R}^n$ is the all-ones vector, Δ is the squared distance matrix of the \mathbf{x}_i , and $N_t(\mathbf{x}_i)$ are calculated based on the ℓ_2 norm difference in this case (the ℓ_2 norm also determines edge length in \mathcal{G}). Additionally, the calculation

$$(5.3) \quad \mathbf{K}_{\text{Iso}} = -\frac{1}{2}\mathbf{H}\Delta\mathbf{H} \approx \mathbf{X}^{*T}\mathbf{X}^*$$

(where \mathbf{X}^* is mean-centered \mathbf{X}) is helpful when geodesic distances Δ_{ij} can be approximated and $\mathbf{X}^{*T}\mathbf{X}^*$ is expensive to calculate directly. The Isomap minimization method is proven in Proposition 5.1, and Lemma 5.2 shows that (*) is an equality when measuring distances with the Euclidean ℓ_2 norm.

Proposition 5.4. *The optimal k -dimensional representation $\mathbf{Y} \in \mathbb{R}^{k \times m}$ due to the Isomap algorithm is given by*

$$\mathbf{Y} = (\boldsymbol{\Sigma}_{Iso, k})^{\frac{1}{2}} \mathbf{U}_{Iso, k}^T$$

where $\boldsymbol{\Sigma}_{Iso, k}$ is the diagonal matrix of the top k singular values of \mathbf{K}_{Iso} and $\mathbf{U}_{Iso, k}$ are the corresponding singular vectors.

Remark 5.5. Note that the top k eigenvectors of \mathbf{K}_{Iso} give the optimal coordinates of the points in a lower k -dimensional space that preserve pairwise geodesic distances specified by $\boldsymbol{\Delta}$. We then scale the matrix according to the corresponding eigenvalues with $\sqrt{\boldsymbol{\Sigma}_{Iso, k}}$. Hence, $\boldsymbol{\Delta}$ in this case behaves as a covariance matrix for the space whose dimensions are defined by the data points (i.e. $m \times m$ in this case), and non-linearity is introduced through calculations of the interpoint distances in $\boldsymbol{\Delta}$.

Definition 5.6. In Lemmas 5.7-6.0 we prove the correctness of double centering (i.e that $-\frac{1}{2}\mathbf{H}\boldsymbol{\Delta}\mathbf{H} = \mathbf{X}^{*T}\mathbf{X}^*$) using Euclidean distance. Hence, we define \mathbf{X} as in Theorem 1.5, and define \mathbf{X}^* to have $\mathbf{x}_i^* := \mathbf{x}_i - \bar{\mathbf{x}}$ as its i -th column. Let $\mathbf{K} := \mathbf{X}^T\mathbf{X}$ and let \mathbf{D} denote the Euclidean distance matrix with $\mathbf{D}_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|$ so that $\mathbf{D} = \boldsymbol{\Delta}$ in this case.

Lemma 5.7. *For \mathbf{K} and \mathbf{D} as defined in Definition 5.6, we have $\mathbf{K}_{ij} = \frac{1}{2}(\mathbf{K}_{ii} + \mathbf{K}_{jj} + \mathbf{D}_{ij}^2)$*

Proof. See Appendix C.1 □

Lemma 5.8. *For \mathbf{K} and \mathbf{X}^* as defined in Definition 5.6, we show that $\mathbf{K}^* := \mathbf{X}^{*T}\mathbf{X}^*$ satisfies*

$$\mathbf{K}^* = \mathbf{K} - \frac{1}{m}\mathbf{K}\mathbf{1}\mathbf{1}^T - \frac{1}{m}\mathbf{1}\mathbf{1}^T\mathbf{K} + \frac{1}{m^2}\mathbf{1}\mathbf{1}^T\mathbf{K}\mathbf{1}\mathbf{1}^T$$

Proof. See Appendix C.2 □

Lemma 5.9. *For \mathbf{K}^* defined in Lemma 5.8 and \mathbf{D} defined in Definition 5.6 we have*

$$\mathbf{K}_{ij}^* = -\frac{1}{2} \left(\mathbf{D}_{ij}^2 - \frac{1}{m} \sum_{t=1}^m (\mathbf{D}_{it}^2 + \mathbf{D}_{tj}^2) + \frac{1}{m^2} \sum_{t=1}^m \sum_{\ell=1}^m \mathbf{D}_{t\ell}^2 \right)$$

Proof. See Appendix C.3 □

Theorem 5.10. *For \mathbf{X}^*, \mathbf{D} from Definition 5.6 we confirm the correctness of the Isomap algorithm for Euclidean distance, namely that $\frac{1}{m}\mathbf{X}^{*T}\mathbf{X}^* = -\frac{1}{2}\mathbf{H}\mathbf{D}\mathbf{H}$ for $\mathbf{H} = \mathbf{I}_m - \frac{1}{m}\mathbf{1}\mathbf{1}^T$.*

Proof. We have

$$(\mathbf{\Delta}(\mathbf{I}_m - \frac{1}{m}\mathbf{1}\mathbf{1}^T))_{\ell j} = \Delta_{\ell j} - \frac{1}{m} \sum_{t=1}^m \Delta_{\ell t}$$

hence we may solve for $(\mathbf{H}\mathbf{\Delta}\mathbf{H})_{ij}$ as

$$((\mathbf{I}_m - \frac{1}{m}\mathbf{1}\mathbf{1}^T)\mathbf{\Delta}(\mathbf{I}_m - \frac{1}{m}\mathbf{1}\mathbf{1}^T))_{ij} = \Delta_{ij} - \frac{1}{m} \sum_{t=1}^m \Delta_{it} - \frac{1}{m} \sum_{\ell=1}^m (\Delta_{\ell j} - \frac{1}{m} \sum_{t=1}^m \Delta_{\ell t})$$

so from Lemma 5.9

$$= -2\mathbf{K}_{ij}^* \Rightarrow \mathbf{K}^* = -\frac{1}{2}\mathbf{H}\mathbf{\Delta}\mathbf{H}$$

□

6. FENCHEL GAME NO-REGRET DYNAMICS (FGNRD)

We now seek to understand such algorithms in the context of the Fenchel Game No-Regret Dynamics framework (FGNRD) introduced by Wang-Abernethy-Levy.

Definition 6.1. For a function $f : \mathcal{K} \rightarrow \mathbb{R} \cup \infty$ where $\mathcal{K} \subset \mathbb{R}^d$, we define its conjugate $f^* : \mathbb{R}^d \rightarrow \mathbb{R} \cup \infty$ as

$$f^*(y) := \sup_{x \in D} \{\langle y, x \rangle - f(x)\}$$

Proposition 6.2. *Conjugate functions of convex functions are convex.*

Proof. For $f : \mathcal{K} \rightarrow \mathbb{R}$ convex where $\mathcal{K} \subset \mathbb{R}^d$, we find that

$$\begin{aligned} f^*(\lambda x + (1-\lambda)y) &= \sup_{x' \in \mathcal{K}} \{\langle x', \lambda x + (1-\lambda)y \rangle - f(x')\} \\ &= \sup_{x' \in \mathcal{K}} \{\langle x', \lambda x + (1-\lambda)y \rangle - f(x')\} \\ &= \sup_{x' \in \mathcal{K}} \{\langle x', \lambda x \rangle + \langle x', y \rangle - \lambda \langle x', y \rangle - f(x')\} \\ &= \sup_{x' \in \mathcal{K}} \{\lambda \langle x, x' \rangle - \lambda f(x') + \langle y, x' \rangle - f(x') - \lambda \langle y, x' \rangle + \lambda f(x')\} \\ &= \sup_{x' \in \mathcal{K}} \{\lambda (\langle x, x' \rangle - f(x')) + (1-\lambda) (\langle y, x' \rangle - f(x'))\} \\ &\leq \lambda \sup_{x' \in \mathcal{K}} \{\langle x, x' \rangle - f(x')\} + (1-\lambda) \sup_{x'' \in \mathcal{K}} \{\langle y, x'' \rangle - f(x'')\} \\ &= \lambda f^*(x) + (1-\lambda) f^*(y) \end{aligned}$$

□

Definition 6.3. The subdifferential $\partial f(x)$ is the set of all subgradients of f at x , i.e.

$$\partial f(x) = \{f_x : f(z) \geq \langle f_x, z - x \rangle + f(x), \forall z \in \mathcal{K}\}$$

Theorem 6.4. *For a closed convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, the following are equivalent:*

- I. $y \in \partial f(x)$
- II. $x \in \partial f^*(y)$
- III. $\langle x, y \rangle = f(x) + f^*(y)$

Proof. We first show I \Rightarrow II. Suppose $y \in \partial f(x)$. Then, for any $z \in \mathcal{K}$ we have

$$f(z) - f(x) \geq \langle y, z - x \rangle \Rightarrow \langle y, x \rangle - f(x) \geq \langle y, z \rangle - f(z)$$

so that

$$\langle y, x \rangle - f(x) \geq f^*(y)$$

Then,

$$\begin{aligned} \langle z, x \rangle - f(x) &\leq f^*(z) \Rightarrow \langle y, x \rangle - f(x) \leq f^*(z) - \langle x, z - y \rangle \\ &\Rightarrow f^*(y) \leq f^*(z) - \langle x, z - y \rangle \Rightarrow x \in \partial f^*(y) \end{aligned}$$

To show II \Rightarrow III, we find that, for any $z \in \mathcal{K}$

$$\langle x, y \rangle - f^*(y) \geq \langle x, z \rangle - f^*(z)$$

hence

$$\langle x, y \rangle - f^*(y) \geq f^{**}(x) = \sup_{y' \in \mathcal{K}} \langle x, y' \rangle - \sup_{x' \in \mathcal{K}} (\langle y', x' \rangle - f(x'))$$

so since f is closed, $\sup_{x' \in \mathcal{K}} (\langle y', x' \rangle - f(x'))$ is attained by some x' as

$$\begin{aligned} &\geq \langle x, \nabla f(x) \rangle - \langle \nabla f(x), x' \rangle + f(x') = f(x') - \langle \nabla f(x), x' - x \rangle \geq f(x) \\ &\Rightarrow f^*(y) \leq \langle x, y \rangle - f(x) \Rightarrow f^*(y) = \langle x, y \rangle - f(x) \end{aligned}$$

Finally, III \Rightarrow I as

$$\begin{aligned} \langle x, y \rangle &\geq f(x) + f^*(y) \Rightarrow \langle x, y \rangle - f(x) \geq \langle y, z \rangle - f(z), \forall z \in \mathcal{K} \\ &\Rightarrow f(z) - f(x) - \langle y, z - x \rangle \geq 0, \forall z \in \mathcal{K} \end{aligned}$$

so that all three statements are equivalent. \square

Definition 6.5. We define our two-input “payoff” function $g : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ as

$$g(x, y) := \langle x, y \rangle - f^*(y)$$

We will understand this function as a zero-sum game in which, if player 1 selects action x and player 2 selects action y , $g(x, y)$ is the “cost” for player 1 and the “gain” for player 2.

Definition 6.6. Given a zero-sum game with a payoff function $g(x, y)$ which is convex in x and concave in y , we define

$$V^* := \inf_{x \in \mathcal{X}} \sup_{y \in \mathcal{Y}} g(x, y)$$

We further define an “ ϵ -equilibrium” of $g(\cdot, \cdot)$ as a pair \hat{x}, \hat{y} for which

$$V^* - \epsilon \leq \inf_{x \in \mathcal{X}} g(x, \hat{y}) \leq V^* \leq \sup_{y \in \mathcal{Y}} g(\hat{x}, y) \leq V^* + \epsilon$$

where \mathcal{X} and \mathcal{Y} are convex decision spaces of the x -player and y -player respectively.

Definition 6.7. To solve for $\inf_{x \in D} f(x)$, we define $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ as

$$g(x, y) := \langle x, y \rangle - f^*(y) = \langle x, y \rangle - \sup_{x' \in \mathcal{K}} \{\langle x', y \rangle - f(x')\}$$

and attempt to find an ϵ -equilibrium for $g(x, y)$.

Proposition 6.8. *An equilibrium for the Fenchel Game function solves the minimization problem $\inf_{x \in D} f(x)$.*

Proof. For an ϵ -equilibrium \hat{x}, \hat{y} of g defined as above, we have

$$\inf_{x \in \mathcal{K}} f(x) = -\sup_{x \in \mathcal{K}} \{-f(x)\} = -\sup_{x' \in \mathcal{K}} \{\langle x', y \rangle - \langle x', y \rangle - f(x')\} =: h(y)$$

so that

$$\inf_{x \in \mathcal{K}} \left\{ \langle x, \hat{y} \rangle - \sup_{x' \in \mathcal{K}} \{\langle x', \hat{y} \rangle - f(x')\} \right\} \leq h(\hat{y}) \leq \sup_{y \in \mathcal{Y}} \left\{ \langle \hat{x}, y \rangle - \sup_{x' \in \mathcal{K}} \{\langle x', y \rangle - f(x')\} \right\}$$

hence

$$(*) \quad |V^* - h(y)| \leq 2\epsilon$$

where

$$V^* = \inf_{x \in \mathcal{K}} \sup_{y \in \mathcal{Y}} \left\{ \langle x, y \rangle - \sup_{x' \in \mathcal{K}} \{\langle x', y \rangle - f(x')\} \right\}$$

and as $\epsilon \rightarrow 0$ we have

$$\begin{aligned} V^* &= \sup_{y \in \mathcal{Y}} \left\{ \langle \hat{x}, y \rangle - \sup_{x' \in \mathcal{K}} \{\langle x', y \rangle - f(x')\} \right\} \\ &= \sup_{y \in \mathcal{Y}} \{\langle \hat{x}, y \rangle - f^*(y)\} = f(\hat{x}) \end{aligned}$$

which follows from Theorem 7.4. \square

Corollary 6.9. *If (\hat{x}, \hat{y}) is an ϵ -equilibrium of the Fenchel Game as defined above, then*

$$|f(\hat{x}) - \inf_{x \in \mathcal{K}} f(x)| \leq \epsilon$$

Proof. Follows from $(*)$ above for $\epsilon' := \frac{\epsilon}{2}$. \square

Definition 6.10. “Online convex optimization” works as follows. At each round t (of T many), the learner selects a point $z_t \in \mathcal{Z}$ and suffers a loss $\alpha_t \ell_t(z_t)$ for this selection, where α is the weight vector and $\mathcal{Z} \subset \mathbb{R}^d$ is a convex decision set of actions.

Remark 6.11. In general it is assumed that, upon selecting z_t during round t , the learner has observed all loss functions $\alpha_1 \ell_1(\cdot), \dots, \alpha_{t-1} \ell_{t-1}(\cdot)$ up to but not including time t . An exception to this are the “prescient” learners, whose algorithms, marked with a “+” superscript, have access to the loss ℓ_t prior to selecting z_t .

Algorithm 3 Protocol for weighted online convex optimization

Require: convex decision set $\mathcal{Z} \subset \mathbb{R}^d$
Require: number of rounds T
Require: weights $\alpha_1, \alpha_2, \dots, \alpha_T > 0$
Require: algorithm OAlg
for $t = 1, 2, \dots, T$ **do**
 Return: $z_t \leftarrow \text{OAlg}$
 Receive: $\alpha_t, \ell_t(\cdot) \rightarrow \text{OAlg}$
 Evaluate: $\text{Loss} \leftarrow \text{Loss} + \alpha_t \ell_t(z_t)$
end for

Remark 6.12. The “OAlg” referenced above refers to an algorithm performed within the current algorithm, and “OAlg^X” will refer to the algorithm updating the x coordinate in the Fenchel Game No Regret Dynamics.

Definition 6.13. We define a learner’s “regret” as

$$\alpha\text{-REG}^z(z^*) := \sum_{t=1}^T \alpha_t \ell_t(z_t) - \sum_{t=1}^T \alpha_t \ell_t(z^*)$$

where $z^* \in \mathcal{Z}$ is the “comparator” to which the online learner is compared. We further define “average regret” as that normalized by the time weight $A_T : \sum_{t=1}^T \alpha_t$ and denote it by

$$\overline{\alpha\text{-REG}}^z(z^*) := \frac{\alpha\text{-REG}^z(z^*)}{A_T}$$

Finally, “no-regret algorithms” guarantee $\overline{\alpha\text{-REG}}^z(z^*) \rightarrow 0$ as $A_T \rightarrow \infty$

Remark 6.14. The following batch-style online-learning strategies modify the central algorithm Follow The Leader (FTL)

Algorithm 4 Online Learning Strategies

Require: convex set \mathcal{Z} , initial point $z_{\text{init}} \in \mathcal{Z}$
Require: $\alpha_1, \dots, \alpha_T > 0$, $\ell_1, \dots, \ell_T : \mathcal{Z} \rightarrow \mathbb{R}$

FTL[z_{init}]:
 $z_t \leftarrow z_{\text{init}}$ **if** $t = 1$, **else**
 $z_t \leftarrow \arg\min_{z \in \mathcal{Z}} \left(\sum_{s=1}^{t-1} \alpha_s \ell_s(z) \right)$
 FTL⁺:
 $z_t \leftarrow \arg\min_{z \in \mathcal{Z}} \left(\sum_{s=1}^t \alpha_s \ell_s(z) \right)$
 FTRL[$R(\cdot), \eta$]:
 $z_t \leftarrow \arg\min_{z \in \mathcal{Z}} \left(\sum_{s=1}^t \alpha_s \ell_s(z) + \frac{1}{\eta} R(z) \right)$

Definition 6.15. A first-order oracle for a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a primitive that, given $x \in \mathbb{Q}^n$, outputs the value $f(x) \in \mathbb{Q}$ and a vector $h(x) \in \mathbb{Q}^n$ such that, for any $z \in \mathbb{R}^n$,

$$f(z) \geq f(x) + \langle h(x), z - x \rangle$$

so $h(x) = \nabla f(x)$ for f differentiable, else it is a subgradient of f at x .

We now examine online mirror descent and its prescient counterpart before recovering gradient descent from the Fenchel Game framework.

Definition 6.16. For fixed $\epsilon > 0$ and norm $\|\cdot\|$, a differentiable function $f : \mathcal{K} \rightarrow \mathbb{R}$ where $\mathcal{K} \subset \mathbb{R}^d$ is convex, is considered “ ϵ -strongly convex with respect to $\|\cdot\|$ ” if

$$f(y) - f(x) \geq \langle \nabla f(x), y - x \rangle + \frac{\epsilon}{2} \|y - x\|^2$$

Definition 6.17. We define the Bregman divergence $D_z^\phi(\cdot)$ centered at z with respect to a β -strongly convex distance generating function $\phi(\cdot)$ as

$$D_z^\phi(x) := \phi(x) - \langle \nabla \phi(z), x - z \rangle - \phi(z)$$

Algorithm 5 Update-style online learning strategies

Require: convex set \mathcal{Z} , initial point $z_0 \in \mathcal{Z}$

Require: $\alpha_1, \dots, \alpha_T > 0$, $\ell_1, \dots, \ell_T : \mathcal{Z} \rightarrow \mathbb{R}$

OMD $[\phi(\cdot), z_0, \gamma]$:

$$z_t \leftarrow \operatorname{argmin}_{z \in \mathcal{Z}} \left(\alpha_{t-1} \ell_{t-1}(z) + \frac{1}{\gamma} D_{z_{t-1}}^\phi(z) \right)$$

OMD $^+[\phi(\cdot), z_0, \gamma]$:

$$z_t \leftarrow \operatorname{argmin}_{z \in \mathcal{Z}} \left(\alpha_t \ell_t(z) + \frac{1}{\gamma} D_{z_{t-1}}^\phi(z) \right)$$

Remark 6.18. To introduce the FGNRD framework, we show vanilla gradient descent (Algorithm 6 below) can be understood as a two player Fenchel Game. For $f(\cdot)$ convex, let

$$G = \max_{y \in \partial f(w), w \in \mathcal{K}} \|y\|^2$$

and let $R \in \mathbb{R}$ be an upper bound to $\|w_0 - w^*\|$ where $w^* := \operatorname{argmin}_{w \in \mathcal{K}} f(w)$.

Remark 6.19. Note that the use of Bregman divergence in the FGNRD Equivalent (which helps specify the subgradient y_t selected) involves the term $\frac{1}{\gamma} D_{z_{t-1}}^\phi(z)$ as in Algorithm 5 to ensure we remain within a neighborhood of our previous iteration upon descent.

Theorem 6.20. *The FGNRD formulation is equivalent to vanilla gradient descent, i.e. $w_t = x_t$ at every time step, hence $\bar{w}_t = \frac{1}{t} \sum_{s=1}^t w_s = \frac{1}{t} \sum_{s=1}^t x_s = \bar{x}_t$.*

Algorithm 6 Vanilla gradient descent and its FGNRD equivalent

Require: Convex function $f(\cdot)$ and iterations T

Initialize: $w_0 = x_0 \in \mathcal{K} \subseteq \mathbb{R}^d$

Initialize: $\gamma = \begin{cases} \frac{R}{G\sqrt{T}}, & \text{if } f(\cdot) \text{ is non-smooth} \\ \frac{1}{2L}, & \text{if } f(\cdot) \text{ is } L\text{-smooth} \end{cases}$

Gradient Descent[w_0]:

$$w_t := w_{t-1} - \gamma \delta_{t-1} \text{ for } \delta_{t-1} \in \partial f(w_{t-1})$$

$$\bar{w}_t := \frac{1}{t} \sum_{s=1}^t w_s$$

FGNRD Equivalent:

$$g(x, y) := \langle x, y \rangle - f^*(y)$$

$$\alpha_t := 1 \text{ for } t \in \{1, \dots, T\}$$

$$\text{OAlg}^X := \text{OMD}[\frac{1}{2} \|\cdot\|_2^2, x_0, \gamma]$$

$$\text{OAlg}^Y := \text{BESTRESP}^+$$

Proof. We find that

$$\begin{aligned} \text{OAlg}^x &:= \operatorname{argmin}_{x \in \mathcal{K}} \alpha_{t-1} \ell_{t-1}(x) + \frac{G\sqrt{T}}{R} \left(\frac{1}{2} \|x\|_2^2 - \frac{1}{2} \|x_{t-1}\|_2^2 - \langle x_{t-1}, x - x_{t-1} \rangle \right) \\ &= \operatorname{argmin}_{x \in \mathcal{K}} \langle x, y_{t-1} \rangle - f^*(y) + \frac{G\sqrt{T}}{R} \left(\frac{1}{2} \|x\|_2^2 - \langle x_{t-1}, x - x_{t-1} \rangle \right) \\ &= \operatorname{argmin}_{x \in \mathcal{K}} \langle x, y_{t-1} - \frac{G\sqrt{T}}{R} x_{t-1} \rangle + \frac{G\sqrt{T}}{2R} \|x\|_2^2 := \operatorname{argmin}_{x \in \mathcal{K}} F(x) \end{aligned}$$

so since \mathcal{K} is convex, the minimum is reached when $\nabla F = \mathbf{0}$, i.e.

$$\mathbf{0} = y_{t-1} - \frac{G\sqrt{T}}{R} x_{t-1} + \frac{G\sqrt{T}}{R} x_t \Rightarrow x_t = x_{t-1} - \frac{R}{G\sqrt{T}} y_{t-1}$$

Now it suffices to show $y_t \in \partial f(x_t)$:

$$\begin{aligned} y_t &= \operatorname{argmin}_{y \in \mathcal{K}} -(\alpha_t \langle x_t, y \rangle - f^*(y)) \\ &= \operatorname{argmax}_{y \in \mathcal{K}} (\langle x_t, y \rangle - f^*(y)) \end{aligned}$$

Then, since, for any $z \in \mathcal{K}$ we have

$$\begin{aligned} \langle x_t, y_t \rangle - f^*(y_t) &\geq \langle x_t, z \rangle - f^*(z) \\ \Rightarrow f^*(z) &\geq f^*(y_t) + \langle x_t, z - y_t \rangle \Rightarrow x_t \in \partial f^*(y_t) \end{aligned}$$

By Theorem 6.4 we thus have

$$y_t \in \partial f(x_t)$$

□

Lemma 6.21. *Let $\phi(\cdot)$ be a β -strongly convex function with respect to the norm $\|\cdot\|_*$, and consider a sequence of lower semi-continuous convex loss functions $\{\alpha_t \ell_t(\cdot)\}_{t=1}^T$. Then, for any comparator $z^* \in \mathcal{Z}$, $\text{OMD}[\phi(\cdot), z_0, \gamma]$ satisfies*

$$\alpha\text{-REG}^z(z^*) \leq \frac{1}{\gamma} D_{z_1}^\phi(z^*) + \frac{\gamma}{2\beta} \sum_{t=1}^T \|\alpha_t \delta_t\|_*^2$$

for $\delta_t \in \partial \ell_t(z_t)$

Proof. We wish to show that

$$\alpha_t \ell_t(\operatorname{argmin}_{z \in \mathcal{Z}} \{\alpha_{t-1} \ell_{t-1}(z) + D_{z_{t-1}}^\phi(z)\}) - \alpha_t \ell_t(z^*) \leq \frac{1}{\gamma} D_{z_1}^\phi(z^*) + \frac{\gamma}{2\beta} \|\alpha_t \delta_t\|_*^2$$

By the minimality of $\alpha_{t-1} \ell_{t-1}(z_{t+1}) + D_{z_{t-1}}^\phi(z_{t+1})$ we have

$$\alpha_t (\ell_t(z_{t+1}) - \ell_t(z^*)) \leq \frac{1}{\gamma} (\phi(z^*) - \phi(z_{t+1}) + \langle \nabla \phi(z_t), z_{t+1} - z^* \rangle)$$

hence

$$(*) \quad \langle \alpha_t \delta_t, z_{t+1} - z^* \rangle \leq \langle \frac{1}{\gamma} (\nabla \phi(z_t) - \nabla \phi(z_{t+1})), z_{t+1} - z^* \rangle$$

since the convexity of ℓ_t ensures the existence of a $\delta_t \in \partial \ell_t(z_t)$ with $\delta_t \in \partial \ell_t(z^*)$.

We now find that

$$\begin{aligned} \alpha_t \ell_t(z_t) - \alpha_t \ell_t(z^*) &\leq \langle \alpha_t \delta_t, z_t - z^* \rangle \\ &= \langle \frac{1}{\gamma} (\nabla \phi(z_{t+1}) - \nabla \phi(z_t)), z^* - z_{t+1} \rangle + \langle \frac{1}{\gamma} (\nabla \phi(z_t) - \nabla \phi(z_{t+1})) - \alpha_t \delta_t, z^* - z_{t+1} \rangle \\ &\quad + \langle \alpha_t \delta_t, z_t - z_{t+1} \rangle \end{aligned}$$

so by (*) we have

$$\begin{aligned} &\leq \langle \frac{1}{\gamma} (\nabla \phi(z_{t+1}) - \nabla \phi(z_t)), z^* - z_{t+1} \rangle + \langle \alpha_t \delta_t, z_t - z_{t+1} \rangle \\ &= \frac{1}{\gamma} (D_{z_t}^\phi(z^*) - D_{z_{t+1}}^\phi(z^*) - D_{z_t}^\phi(z_{t+1})) + \langle \alpha_t \delta_t, z_t - z_{t+1} \rangle \end{aligned}$$

Then, since $\langle \alpha_t \nabla \ell_t(z_t), z_t - z_{t+1} \rangle \leq \frac{\gamma}{2\beta} \|\alpha_t \delta_t\|_*^2 + \frac{\beta}{2\gamma} \|z_t - z_{t+1}\|^2$,

$$\leq \frac{1}{\gamma} (D_{z_t}^\phi(z^*) - D_{z_{t+1}}^\phi(z^*)) + \frac{\gamma}{2\beta} D_{z_t}^\phi(z^*)$$

We now sum from $t = 1$ to $t = T$ and the result follows. \square

Lemma 6.22. *Algorithm 6 satisfies $f(\bar{w}_T) - \min_{w \in \mathcal{K}} f(w) = O(\frac{GR}{\sqrt{T}})$.*

Proof. By Lemma 6.17, OAlg^X satisfies

$$\alpha\text{-REG}^x(x^*) \leq \frac{1}{\gamma} D_{x_0}^\phi(x^*) + \frac{\gamma}{2} \sum_{t=1}^T \|\alpha_t y_t\|_2^2$$

while OAlg^Y suffers no regret (prescient). Further, since

$$\|y\|_2^2 - \|x\|_2^2 - 2\langle x, y - x \rangle = \|y - x\|_2^2$$

we have that $\phi(x) = \frac{1}{2} \|x\|_2^2$ is 1-strongly convex with respect to the Euclidean norm $\|\cdot\|_2$, hence

$$\begin{aligned} \overline{\alpha\text{-REG}}^x[\text{OMD}] + \overline{\alpha\text{-REG}}^x[\text{BESTRESP}^+] &\leq \frac{1}{A_t} \left(\frac{1}{\gamma} D_{x_0}^\phi(x^*) + \sum_{t=1}^T \frac{\gamma}{2} \|\alpha_t y_t\|_2^2 \right) \\ &\leq \frac{1}{T} \left(\frac{R^2}{\gamma} + \frac{\gamma T G^2}{2} \right) \\ &= \frac{1}{T} \left(R G \sqrt{T} + \frac{R G \sqrt{T}}{2} \right) = \frac{3 R G}{2 \sqrt{T}} \end{aligned}$$

$$= O\left(\frac{GR}{\sqrt{T}}\right)$$

□

Remark 6.23. The gradient descent algorithm presented in Algorithm 6 is that used in the gradient descent implementation of PCA. In particular, the generated $n \times k$ matrices used in the Python experiment (see Figures 1 and 2) have norm at most 1. Further, the gradient y is kept normalized in the experiment, and $\|\operatorname{argmin}_{w \in K} f(w)\| \leq 24$ so that $G = 1$ and $R = 25$ over the $T = 10$ iterations, which guarantees a convergence rate of

$$f(\bar{w}_1 0) - \min_{w \in K} f(w) \leq \frac{3 * 24}{2\sqrt{10}} \leq \frac{25}{2} = 12.5$$

However, empirically this bound is only approximate since

$$\operatorname{argmin}_{\mathbf{U}_k \in \mathbb{R}^{n \times k}} \|\mathbf{U}_k \mathbf{U}_k^T \mathbf{X} - \mathbf{X}\|_F^2$$

is a non-convex minimization as in Remark 1.9.

ACKNOWLEDGMENTS

I would like to thank my mentor, Antares Chen, for helping direct my research throughout the summer. I would also like to thank Peter May for organizing this REU.

7. BIBLIOGRAPHY

REFERENCES

- [1] <http://www.ams.org/publications/authors/tex/amslatex>
- [2] Jun-Kun Wang, Jacob Abernethy, Kfir Y. Levy. No-Regret Dynamics in the Fenchel Game: A Unified Framework for Algorithmic Convex Optimization <https://arxiv.org/abs/2111.11309>
- [3] Mehryar Mohri. Foundations of Machine Learning (second edition) MIT Press. 2018.
- [4] Wei Jiang, Nan Li, Hongpeng Yin and Yi Chai. An Improved Laplacian Eigenmaps Algorithm for Nonlinear Dimensionality Reduction. Proceedings of the 2015 Chinese Intelligent Systems Conference

APPENDIX A.

Theorem A.1. Let $\mathbf{P}^* \in \mathcal{P}_k$ be the PCA solution for a centered data matrix $\mathbf{X} = (\Phi(x_1), \dots, \Phi(x_m)) \in \mathbb{R}^{n \times m}$. Then, $\mathbf{P}^* = \mathbf{U}_k \mathbf{U}_k^T$, where $\mathbf{U}_k \in \mathbb{R}^{n \times k}$ is the matrix formed by the top k singular vectors of $\mathbf{C} := \frac{1}{m} \mathbf{X} \mathbf{X}^T$. Moreover, the associated k -dimensional representation of \mathbf{X} is given by $\mathbf{Y} = \mathbf{U}_k^T \mathbf{X}$.

Proof. For $\mathbf{P} = \mathbf{P}^T$ an orthogonal projection matrix, we seek to minimize

$$\begin{aligned} \|\mathbf{P}\mathbf{X} - \mathbf{X}\|_F^2 &= \sum_{i=1}^n \sum_{j=1}^n ((\mathbf{P}\mathbf{X} - \mathbf{X})_{ij})^2 = \operatorname{Tr}[(\mathbf{P}\mathbf{X} - \mathbf{X})^T (\mathbf{P}\mathbf{X} - \mathbf{X})] \\ &= \operatorname{Tr}[\mathbf{X}^T \mathbf{P}^2 \mathbf{X} - \mathbf{X}^T \mathbf{P}^T \mathbf{X} - \mathbf{X}^T \mathbf{P} \mathbf{X} + \mathbf{X}^T \mathbf{X}] = \operatorname{Tr}[\mathbf{X}^T \mathbf{P} \mathbf{X} - 2\mathbf{X}^T \mathbf{P} \mathbf{X} + \mathbf{X}^T \mathbf{X}] \\ &= \operatorname{Tr}[\mathbf{X}^2] - \operatorname{Tr}[\mathbf{X}^T \mathbf{P} \mathbf{X}] \end{aligned}$$

hence we would like to maximize

$$\operatorname{Tr}[\mathbf{X}^T \mathbf{P} \mathbf{X}] = \operatorname{Tr}[\mathbf{X}^T \mathbf{U}_k \mathbf{U}_k^T \mathbf{X}] = \operatorname{Tr}[\mathbf{U}_k^T \mathbf{X} \mathbf{X}^T \mathbf{U}_k]$$

$$= \sum_{i=1}^k \left(\sum_{j=1}^n (\mathbf{U}_k^T \mathbf{X} \mathbf{X}^T)_{ij} (\mathbf{U}_k)_{ji} \right) = \sum_{i=1}^k \left(\sum_{j=1}^n \left(\sum_{\ell=1}^n (\mathbf{U}_k^T)_{i\ell} (\mathbf{X} \mathbf{X}^T)_{\ell j} \right) (\mathbf{U}_k)_{ji} \right)$$

so for $\mathbf{u}_i := ((\mathbf{U}_k)_{1i}, \dots, (\mathbf{U}_k)_{ni})$,

$$= \sum_{i=1}^n \left(\mathbf{u}_i^T \mathbf{X} \mathbf{X}^T \mathbf{u}_i \right)$$

where

$$\mathbf{P} \mathbf{X} = \mathbf{U}_k \mathbf{U}_k^T \mathbf{X}$$

so that $\mathbf{Y} := \mathbf{U}_k^T \mathbf{X}$ is a k -dimensional representation of \mathbf{X} . \square

Theorem A.2. *The following holds for the PCA gradient:*

$$\frac{\partial}{\partial \mathbf{U}_k} \|\mathbf{U}_k \mathbf{U}_k^T \mathbf{X} - \mathbf{X}\|_F^2 = 2(\mathbf{X} \mathbf{X}^T \mathbf{U} \mathbf{U}^T \mathbf{U} + \mathbf{U} \mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U} - 2\mathbf{X} \mathbf{X}^T \mathbf{U})$$

Proof. We find that

$$\frac{\partial}{\partial \mathbf{U}_k} \|\mathbf{P} \mathbf{X} - \mathbf{X}\|_F^2 = \frac{\partial}{\partial \mathbf{U}_k} \left(\text{Tr}[\mathbf{X}^T \mathbf{P}^2 \mathbf{X}] - 2\text{Tr}[\mathbf{X}^T \mathbf{P} \mathbf{X}] + \text{Tr}[\mathbf{X}^T \mathbf{X}] \right)$$

We decompose the right hand side and relax notation as $\mathbf{U} := \mathbf{U}_k$ (since we know our desired dimension k):

$$\begin{aligned} \frac{\partial}{\partial \mathbf{U}} \text{Tr}[\mathbf{X}^T \mathbf{P}^2 \mathbf{X}] &= \frac{\partial}{\partial \mathbf{U}} \text{Tr}[\mathbf{X}^T \mathbf{U} \mathbf{U}^T \mathbf{U} \mathbf{U}^T \mathbf{X}] \\ &= (\mathbf{X}^T)^T (\mathbf{U}^T \mathbf{U} \mathbf{U}^T \mathbf{X})^T + (\mathbf{U} \mathbf{U}^T \mathbf{X})(\mathbf{X}^T \mathbf{U}) + (\mathbf{X}^T \mathbf{U} \mathbf{U}^T)(\mathbf{U}^T \mathbf{X})^T + (\mathbf{X})(\mathbf{X}^T \mathbf{U} \mathbf{U}^T \mathbf{U}) \\ &= 2\mathbf{X} \mathbf{X}^T \mathbf{U} \mathbf{U}^T \mathbf{U} + 2\mathbf{U} \mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U} \end{aligned}$$

so that

$$\begin{aligned} &\frac{\partial}{\partial \mathbf{U}} \left(\text{Tr}[\mathbf{X}^T \mathbf{P}^2 \mathbf{X}] - 2\text{Tr}[\mathbf{X}^T \mathbf{P} \mathbf{X}] \right) \\ &= 2\mathbf{X} \mathbf{X}^T \mathbf{U} \mathbf{U}^T \mathbf{U} + 2\mathbf{U} \mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U} - 2 \left((\mathbf{X}^T)^T (\mathbf{U}^T \mathbf{X})^T + (\mathbf{X})(\mathbf{X}^T \mathbf{U}) \right) \\ &= 2(\mathbf{X} \mathbf{X}^T \mathbf{U} \mathbf{U}^T \mathbf{U} + \mathbf{U} \mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U} - 2\mathbf{X} \mathbf{X}^T \mathbf{U}) \end{aligned}$$

\square

APPENDIX B.

Proposition B.1. *The solution to the laplacian eigenmaps minimization is $\mathbf{U}_{\mathbf{L},k}^T$, where $\mathbf{L} = \mathbf{D} - \mathbf{W}$ is the “graph Laplacian” and $\mathbf{U}_{\mathbf{L},k}^T$ are the bottom k singular vectors of \mathbf{L} (excluding 0 if the underlying neighborhood graph has connections).*

Proof. We find that, for $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{Y} \in \mathbb{R}^{k \times m}$ we have

$$\begin{aligned} (\mathbf{Y} \mathbf{L} \mathbf{Y}^T)_{ij} &= \sum_{\ell=1}^m \sum_{t=1}^m Y_{i\ell} L_{\ell t} Y_{tj}^T = \sum_{\ell,t} Y_{i\ell} (\mathbf{D} - \mathbf{W})_{\ell t} Y_{jt} \\ &= \sum_{\ell, t \neq \ell} Y_{i\ell} (-W_{\ell t}) Y_{jt} + \sum_{\ell=t, s \neq \ell} Y_{i\ell} W_{\ell s} Y_{js} \\ (B.2) \quad &= \sum_{\ell, t \neq \ell} W_{\ell t} (Y_{i\ell} Y_{jt} - Y_{it} Y_{j\ell}) \end{aligned}$$

while

$$\begin{aligned}
 \sum_{\ell,t} W_{\ell t} \|\mathbf{y}'_{\ell} - \mathbf{y}'_t\|_2^2 &= \sum_{\ell,t} W_{\ell t} (\mathbf{y}'_{\ell} - \mathbf{y}'_t)^T (\mathbf{y}'_{\ell} - \mathbf{y}'_t) \\
 &= \sum_{\ell,t} W_{\ell t} ((\mathbf{y}'_{\ell})^2 - 2(\mathbf{y}'_t^T \mathbf{y}'_{\ell}) + (\mathbf{y}'_t)^2) \\
 &= \sum_{\ell,t} W_{\ell t} \left(\sum_{j=1}^m (\mathbf{y}'_{\ell})_j^2 - 2(\mathbf{y}'_t)_j (\mathbf{y}'_{\ell})_j + (\mathbf{y}'_t)_j^2 \right) \\
 &= \sum_{\ell,t} W_{\ell t} \left(\sum_{j=1}^m Y'_{j\ell}{}^2 - 2Y'_{jt} Y'_{j\ell} + Y'_{jt}{}^2 \right)
 \end{aligned}$$

hence by (3.3)

$$= \sum_{j=1}^k 2(\mathbf{Y}' \mathbf{L} \mathbf{Y}'^T)_{jj}$$

so for $\mathbf{Y} := \mathbf{Y}'^T$, by the final simplification used in Theorem 1.5,

$$= 2 \sum_{j=1}^k \mathbf{y}_j^T \mathbf{L} \mathbf{y}_j$$

thus $\mathbf{Y} = \mathbf{U}_{\mathbf{L},k}^T$ are the bottom k singular vectors of \mathbf{L} . \square

APPENDIX C.

Lemma C.1. *For \mathbf{K} and \mathbf{D} as defined in Definition 5.6, we have $\mathbf{K}_{ij} = \frac{1}{2}(\mathbf{K}_{ii} + \mathbf{K}_{jj} + \mathbf{D}_{ij}^2)$*

Proof. We find that

$$\begin{aligned}
 \mathbf{K}_{ij} &= \sum_{\ell=1}^m \mathbf{x}_{i\ell}^T \mathbf{x}_{\ell j} = \frac{1}{2} \left(\sum_{\ell=1}^m \mathbf{x}_{\ell i}^2 - \mathbf{x}_{\ell i}^2 + \mathbf{x}_{\ell j}^2 - \mathbf{x}_{\ell j}^2 + 2\mathbf{x}_{\ell i} \mathbf{x}_{\ell j} \right) \\
 &= \frac{1}{2} \left(\sum_{\ell=1}^m \mathbf{x}_{\ell i}^2 + \mathbf{x}_{\ell j}^2 - (\mathbf{x}_{\ell j} - \mathbf{x}_{\ell i})^2 \right) = \frac{1}{2} (\mathbf{K}_{ii} + \mathbf{K}_{jj} - \|\mathbf{x}_i - \mathbf{x}_j\|^2) \\
 &= \frac{1}{2} (\mathbf{K}_{ii} + \mathbf{K}_{jj} - \mathbf{D}_{ij}^2)
 \end{aligned}$$

\square

Lemma C.2. *For \mathbf{K} and \mathbf{X}^* as defined in Definition 5.6, we show that $\mathbf{K}^* := \mathbf{X}^{*T} \mathbf{X}^*$ satisfies*

$$\mathbf{K}^* = \mathbf{K} - \frac{1}{m} \mathbf{K} \mathbf{1} \mathbf{1}^T - \frac{1}{m} \mathbf{1} \mathbf{1}^T \mathbf{K} + \frac{1}{m^2} \mathbf{1} \mathbf{1}^T \mathbf{K} \mathbf{1} \mathbf{1}^T$$

Proof. Let $\mathbf{K}^* := \mathbf{X}^{*T} \mathbf{X}^*$. We have

$$\begin{aligned}
 \frac{1}{m} (\mathbf{K} \mathbf{1} \mathbf{1}^T)_{ij} &= \frac{1}{m} \sum_{t=1}^m \mathbf{K}_{it} = \frac{1}{m} \sum_{t=1}^m \sum_{\ell=1}^m \mathbf{x}_{\ell i} \mathbf{x}_{\ell t} = \sum_{\ell=1}^m (\bar{\mathbf{x}})_{\ell} (\mathbf{x}_i)_{\ell} \\
 \frac{1}{m} (\mathbf{1} \mathbf{1}^T \mathbf{K})_{ij} &= \frac{1}{m} \sum_{t=1}^m \mathbf{K}_{tj} = \frac{1}{m} \sum_{t=1}^m \sum_{\ell=1}^m \mathbf{x}_{\ell t} \mathbf{x}_{\ell j} = \sum_{\ell=1}^m (\bar{\mathbf{x}})_{\ell} (\mathbf{x}_j)_{\ell}
 \end{aligned}$$

and

$$\frac{1}{m^2}(\mathbf{1}\mathbf{1}^T\mathbf{K}\mathbf{1}\mathbf{1}^T)_{ij} = \frac{1}{m^2} \sum_{t=1}^m (\mathbf{1}\mathbf{1}^T)_{it} (\mathbf{K}\mathbf{1}\mathbf{1}^T)_{tj} = \frac{1}{m} \sum_{t=1}^m \sum_{\ell=1}^m (\bar{\mathbf{x}})_{\ell} (\mathbf{x}_t)_{\ell} = \sum_{\ell=1}^m (\bar{\mathbf{x}}_{\ell})^2$$

Then,

$$\begin{aligned} \mathbf{K}_{ij}^* &= \sum_{\ell=1}^N \mathbf{X}_{i\ell}^{*T} \mathbf{X}_{j\ell}^* = \sum_{\ell=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})_{\ell} (\mathbf{x}_j - \bar{\mathbf{x}})_{\ell} \\ &= \sum_{\ell=1}^N (\mathbf{x}_i)_{\ell} (\mathbf{x}_j)_{\ell} - (\mathbf{x}_i)_{\ell} (\bar{\mathbf{x}})_{\ell} - (\mathbf{x}_j)_{\ell} (\bar{\mathbf{x}})_{\ell} + (\bar{\mathbf{x}})_{\ell}^2 \\ &= \mathbf{K}_{ij} - \frac{1}{m} (\mathbf{K}\mathbf{1}\mathbf{1}^T)_{ij} - \frac{1}{m} (\mathbf{1}\mathbf{1}^T\mathbf{K})_{ij} + \frac{1}{m^2} (\mathbf{1}\mathbf{1}^T\mathbf{K}\mathbf{1}\mathbf{1}^T)_{ij} \end{aligned}$$

so that

$$\mathbf{K}^* = \mathbf{K} - \frac{1}{m} \mathbf{K}\mathbf{1}\mathbf{1}^T - \frac{1}{m} \mathbf{1}\mathbf{1}^T\mathbf{K} + \frac{1}{m^2} \mathbf{1}\mathbf{1}^T\mathbf{K}\mathbf{1}\mathbf{1}^T$$

□

Lemma C.3. For \mathbf{K}^* defined in Lemma 5.8 and \mathbf{D} defined in Definition 5.6 we have

$$\mathbf{K}_{ij}^* = -\frac{1}{2} \left(\mathbf{D}_{ij}^2 - \frac{1}{m} \sum_{t=1}^m (\mathbf{D}_{it}^2 + \mathbf{D}_{tj}^2) + \frac{1}{m^2} \sum_{t=1}^m \sum_{\ell=1}^m \mathbf{D}_{t\ell}^2 \right)$$

Proof. From Lemma 5.8 we have

$$\begin{aligned} \mathbf{K}_{ij}^* &= \mathbf{K}_{ij} - \frac{1}{m} (\mathbf{K}\mathbf{1}\mathbf{1}^T)_{ij} - \frac{1}{m} (\mathbf{1}\mathbf{1}^T\mathbf{K})_{ij} + \frac{1}{m^2} (\mathbf{1}\mathbf{1}^T\mathbf{K}\mathbf{1}\mathbf{1}^T)_{ij} \\ &= \frac{1}{2} (\mathbf{K}_{ii} + \mathbf{K}_{jj} - \mathbf{D}_{ij}^2) - \frac{1}{m} (\mathbf{K}\mathbf{1}\mathbf{1}^T)_{ij} - \frac{1}{m} (\mathbf{1}\mathbf{1}^T\mathbf{K})_{ij} + \frac{1}{m^2} (\mathbf{1}\mathbf{1}^T\mathbf{K}\mathbf{1}\mathbf{1}^T)_{ij} \\ &= \frac{1}{2} (\mathbf{K}_{ii} + \mathbf{K}_{jj} - \mathbf{D}_{ij}^2) - \frac{1}{m} \sum_{t=1}^m \mathbf{K}_{it} - \frac{1}{m} \sum_{t=1}^m \mathbf{K}_{tj} + \frac{1}{m^2} \sum_{t=1}^m \sum_{\ell=1}^m \mathbf{K}_{t\ell} \end{aligned}$$

so that applying Lemma 5.7 to \mathbf{K}_{it} and \mathbf{K}_{tj} we have

$$\begin{aligned} &= \frac{1}{2} (\mathbf{K}_{ii} + \mathbf{K}_{jj} - \mathbf{D}_{ij}^2) - \frac{1}{2m} \sum_{t=1}^m \left((\mathbf{K}_{ii} + \mathbf{K}_{tt} - \mathbf{D}_{it}^2) + (\mathbf{K}_{tt} + \mathbf{K}_{jj} - \mathbf{D}_{tj}^2) - \frac{1}{m} \sum_{\ell=1}^m (\mathbf{K}_{t\ell} + \mathbf{K}_{\ell\ell} - \mathbf{D}_{t\ell}^2) \right) \\ &= \frac{1}{2} (-\mathbf{D}_{ij}^2) - \frac{1}{2m} \sum_{t=1}^m \left((\mathbf{K}_{tt} - \mathbf{D}_{it}^2) - \mathbf{D}_{tj}^2 - \frac{1}{m} \sum_{\ell=1}^m (\mathbf{K}_{\ell\ell} - \mathbf{D}_{t\ell}^2) \right) \\ &= \frac{1}{2} \left(-\mathbf{D}_{ij}^2 - \frac{1}{m} \sum_{t=1}^m (\mathbf{K}_{tt} - \mathbf{D}_{it}^2 - \mathbf{D}_{tj}^2) + \frac{1}{m^2} \sum_{t=1}^m \sum_{\ell=1}^m (\mathbf{K}_{\ell\ell} - \mathbf{D}_{t\ell}^2) \right) \\ &= -\frac{1}{2} \left(\mathbf{D}_{ij}^2 - \frac{1}{m} \sum_{t=1}^m (\mathbf{D}_{it}^2 + \mathbf{D}_{tj}^2) + \frac{1}{m^2} \sum_{t=1}^m \sum_{\ell=1}^m \mathbf{D}_{t\ell}^2 \right) \end{aligned}$$

□