

## Course Project: Subspace Clustering

Due: March 24, 2023, 11:59PM PT

*Student Names: Matthew Sullivan and Will Mijangos*

# 1 Paper Descriptions

## 1.1 Paper 1

**Paper Title:** Sparse Subspace Clustering: Algorithm, Theory, and Applications**Student Name:** Matthew Sullivan

### 1.1.1 Problem Description & Formulation

*Sparse Subspace Clustering: Algorithm, Theory, and Applications* contributed a significant framework for clustering high dimensional data. The algorithm, Sparse Subspace Clustering (SSC), relied on the notion of high dimensional data lying in or near some low dimensional structure. In this case, SSC searched for a set of subspaces to cluster data. More importantly, SSC asserted that the correct subspace clustering of data shall be the sparsest [1]. As a result, SSC focused on selecting a sparse representation for the data. This is important as it is robust to common errors in data such as missing and noisy data.

Clustering is typically unsupervised, and this case is no different. SSC relied on the property of self-expressiveness, formally defined as, "each data point in a union of subspaces can be efficiently reconstructed by a combination of other points in the dataset" (Section 2.1) [1]. This can be compared to a linear dependency: a data point that is self-expressive is recoverable through a linear combination of exclusively other data points, Eq. 2.

In the paper, the data is defined to have a set of  $N$  noise-free samples,  $\{y_1, \dots, y_N\}$ , with each sample of data being  $D$  dimensional,  $y_i \in \mathbb{R}^D$ . The collection of data is slotted into a matrix,  $\mathbf{Y}$ , as columns altered by an unknown permutation matrix  $\mathbf{\Gamma}$ , seen in Eq. 1. From Eq. 1, one sees how the samples,  $y_i$ , are the columns of the larger  $\mathbf{Y}$  matrix. One critical component is that  $\mathbf{Y}$  needs to be over-determined  $\text{rank}(\mathbf{Y}) < N$ , to allow linear dependency between the data. If we do not attain this, some data will lie in a nullspace and won't be self-expressive.

At the same time, there exists a set of  $n$ -many linear subspaces,  $\{S_\ell\}_{\ell=1}^n$ , with each subspace being in  $d$ -dimensional,  $S_i \in \mathbb{R}^D$ . All of the noise-free data points must lie in a union of these linear subspaces. It's important to note that  $n$  is a controllable variable. We often do not know how many linear subspaces the data lies in. The paper suggested one can make a few assumptions on the dimension of the subspaces using the weights of the singular values.

$$\mathbf{Y} \triangleq [\mathbf{y}_1 \cdots \mathbf{y}_N] = [\mathbf{Y} \cdots \mathbf{Y}]\mathbf{\Gamma} \quad (1)$$

Given that  $\mathbf{Y}$  has been made, SSC is a sparse optimization algorithm to find the fewest amount of coefficients to calculate the self-expressive representation. Eq. 3 visualizes self-expressiveness as data that is linearly dependent. The sparse optimization problem locates the coefficients,  $c$ , necessary to reconstruct  $\mathbf{y}_i$ . Ideally, the number of coefficients is equal to the dimension of the subspace, allowing the other data points to be the basis vectors. In practice, this number is often unknown and overestimated as a result.

$$\mathbf{y}_i = \mathbf{Y}\mathbf{c}_i \quad c_{ii} = 0 \quad (2)$$

$$\mathbf{y}_i = c_1 \mathbf{y}_1 + c_2 \mathbf{y}_2 + \cdots + 0 \mathbf{y}_i + \cdots + c_n \mathbf{y}_n \quad (3)$$

To solve SSC, a sparse optimization problem, Eq. 4 must be solved. It can be solved using convex optimization techniques [1].

### 1.1.2 Algorithm Description

Sparse Subspace Clustering (SSC) is an optimization problem that yields a sparse solution – the coefficients necessary to reconstruct any data point using exclusively other data points. It can be compared to PCA. PCA components are subspaces that capture the most variance in the data in descending order. Instead of capturing variance, SSC focuses on capturing the least amount of coefficients to build subspaces with the most self-expressive points. That is to say: SSC finds the least amount of coefficients to reconstruct  $\mathbf{y}_i$  and  $\mathbf{y}_j, i \neq j$ , given they lie in the same subspace. It does this by locating coefficients for linear combinations in the optimization problem shown in Eq. 4 and Eq. 5. Both of these are convex and can use convex optimization techniques like alternating direction method of multipliers. After the coefficients have been found and stored in a matrix,  $\mathbf{C}$ , and made sparse – either by a threshold or keeping the top-K largest – a weight matrix is formed through  $W = |\mathbf{C}| + |\mathbf{C}^T|$  – which yields a Laplacian. This Laplacian is clustered using spectral clustering with the results returned from the algorithm.

$$\min \|\mathbf{c}_i\|_1 \text{ such that } \mathbf{y}_i = \mathbf{Y}\mathbf{c}_i, \quad c_{ii} = 0 \quad (4)$$

$$\hat{\mathbf{c}}_i = \arg \min_{\mathbf{c}_i} \|\mathbf{c}_i\|_1 + \lambda \|\mathbf{y}_i - \mathbf{Y}\mathbf{c}_i\|_1, \quad c_{ii} = 0 \quad (5)$$

---

#### Algorithm 1 Sparse Subspace Clustering

---

**Initialize:** Set of data  $\{y_1, \dots, y_n\}$  that lies in a union of  $n$  linear subspaces

Solve Eq. 4 sparsely. We did this by solving 5.

Normalize  $\mathbf{c}_i = \frac{\mathbf{c}_i}{\|\mathbf{c}_i\|_1}, \forall i \in N$

Form Weight Matrix  $\mathbf{W} = \|\mathbf{C}\| + \|\mathbf{C}\|^T$

Spectral Cluster  $\mathbf{W}$

**Output:** Clusters, Clustered Data

---

The paper also has a section on "Practical Extensions". In clustering, there seems to be a few recurring problems: data with noisy outliers, incomplete data, and data that lies in affine subspaces instead of linear subspaces. Data that lies in an affine subspaces means the data does not appear to be self-expressive – there is no linear combination because there is no origin for the subspace. Loosely speaking, without an origin, there is no starting point for basis vectors, therefore, no basis vectors exist. The paper produces solutions for these predicament that require altering your sparse optimization problem and constraints. It works very well empirically. However, its computational cost is exponential on the dimension of the subspaces.

### 1.1.3 Theoretical Results

The foundational assumption is that the sparse optimization program recovers coefficients that can recover linearly dependent data. For example, given  $N - 1$  even valued data and 1 prime data point, SSC would have to recover odd, even, and negative valued coefficients to reconstruct the prime data point. If it does not, the prime data point will be lost. There is no assumption on the distribution of data in each subspace [1].

There are certain conditions that must be met for sparse subspace recovery. The main condition is that the point to be recovered must lie within some subspace sparse representation to begin with – it cannot lie in the nullspace. It can lie in intersection of subspaces. Furthermore, these conditions must hold for two

cases: the smallest coefficients,  $\mathbf{c}_i$ , to recover a data point and another condition I do not understand, Eq. 6. SSC converges to a local minimum.

$$\mathbf{c}_i = \arg \min_{\mathbf{c}_i} \|\mathbf{c}_i\|_1 \text{ such that } \mathbf{y}_i = Y_{-i} \mathbf{c}^4 \quad (6)$$

#### 1.1.4 Relation to Course Material

The algorithm leverages convex optimization – a topic loosely mentioned in class. The true sparse optimization problem is a  $\ell_0$ -norm instead of an  $\ell_1$ -norm in Eq. 4. Similar to class, the paper performs a convex relaxation to make it operable. Because it is convex, our initial programming solution used gradient descent. It was slow. I spent a few hours trying to get an ADMM solution, however, I could not produce it. The spectral clustering is the exact same as introduced in the first demo, and the solution code is used in our replication of SSC. The idea of self-expressiveness is just a linear combination of data where basis vectors are the coefficients of other data points – which could be viewed as a new origin to define the subspace.

I do not understand the graph theory very well. I understand more of the sparse theory than I anticipated. I think these should both be possible within the next few months. Looking at the ADMM solution, I understand parts of it, but not as much as I would like.

In my opinion, the theory exceeds the content of the course. While part are recognizable, the interconnect is missing, starting from Section 4.2.

## 1.2 Paper 2

**Paper Title:** Subspace Clustering using Ensembles of  $K$ -Subspaces

**Student Name:** Matthew Sullivan

*Subspace Clustering using Ensembles of  $K$ -Subspaces* (EKSS) is accurate to the description of the algorithm: use an ensemble of  $K$ -subspaces to vote on the correct clustering. This works on "evidence accumulation clustering": the idea that something about the random clustering must be partially correct [2]. If random clusterings have correct information, the correct clustering is recoverable if one can extract the partially correct information from the random clusters. EKSS extracts the correct information by performing a PCA on the clustered points and counting the number of times a point is assigned to a cluster.

The  $K$ -subspaces algorithm wants to minimize the sum of residuals from the original points to their projected subspace, Eq. 7. The estimated clusters,  $=\{c_1, \dots, c_K\}$ , lie in subspaces with assigned orthonormal bases,  $\mathcal{U} = \{U_1, \dots, U_K\}$ . Eq. 7 is unfortunately very computationally expensive, and is commonly solved using an alternating algorithm to find a global minima. However, if ran, KSS does get some information correct, and this is the foundation for consensus clustering and EKSS.

$$\min \|y - Ax\|_2^2 \longrightarrow \min_{\mathcal{C}, \mathcal{U}} \sum_{k=1}^K \sum_{i: x_i \in c_k} \|x_i - U_k U_k^T x_i\|_2^2, \quad A = U_k U_k^T, x_i \in y \quad (7)$$

EKSS uses consensus clustering to extract information from KSS and form a similarity matrix. Consensus clustering is randomly clustering points and assigning them points to each clustering such that, if repeated many times, some pattern of information emerges.

Generally, subspace clustering is an important field as it is a form of unsupervised classification.  $K$ -subspaces (KSS) was a previous algorithm discussed for its efficiency and performance. One of its failures, however, was a lack of theoretical guarantees – something this paper provides for specific cases. Furthermore, this paper produced EKSS – a altered form of KSS that has theoretical guarantees and a competitive theoretical run time. To simplify things: KSS was a desirable but lacked theoretical guarantees; EKSS has a specific approach that is competitive to KSS and has theoretical guarantees.

### 1.2.1 Algorithm Description

EKSS is a geometric algorithm – it performs clustering by utilizing properties of the data. It uses iterations to capture information within the structure of the data to cluster it. Furthermore, it assumes that the potential subspaces are the true possible subspaces for the data and noise free.

Given the noise-free subspaces and the possibly noisy data, clustering repetitively should eventually promote the true clustering for the data. This can be compared to the Central Limit Theorem. Imagine the true subspaces as the true mean. Every time you randomly cluster, you receive a sample mean. Eventually, after enough random clusters, your sample means will converge to the true mean – the correct underlying subspace. A variation of EKSS, EKSS-0, performs well empirically when compared to Thresholded Subspace Clustering.

---

#### Algorithm 2 Subspace Clustering using Ensembles of $K$ -Subspaces

---

**Initialize:**  $\mathcal{X}$  data,  $\bar{K}$  candidate subspaces,  $\bar{d}$  dimensions for each  $\bar{K}$  subspace,  $K$  number of clusters, threshold  $q$ , Number of base clusterings  $B$ ,  $T$  number of iterations  
**for**  $b=1, \dots, B$  **do**  
    Generate Random Bases,  $\mathcal{U} = \{U_1, \dots, U_{\bar{K}}\}$ , each  $\bar{d}$  in dimension for KSS Clustering  
    Form cluster,  $c_k$ , by projecting data onto generated clusters, Eq. 7  
    **for**  $t=1, \dots, T$  **do**  
        Use PCA to estimate subspace bases  
        Refined cluster,  $c_k$ , using estimated subspace bases  
    **end for**  
    form,  $C^b$ , a matrix of clusters. Save up all of these for each  $b$ .  
**end for**  
From all the runs, count up number of times a cluster was used. Store in similarity matrix  $A$   
Make  $A$  sparse by thresholding at  $q$ , producing  $\bar{A}$  **State Output:** Spectral Clusters,  $C$ , by spectral clustering  $\bar{A}$  with  $K$  clusters

---

### 1.2.2 Theoretical Results

Given two dot products,  $x^T y$  and  $x^T z$ , and one is significantly larger than the other,  $x^T y \gg x^T z$ , one could say  $y$  is more similar to  $x$  than  $z$ . Thresholding dot products produced a thresholding subspace clustering approach. Dot products can be viewed from a different perspective, however. In Eq. 8, dot products are computed as the magnitude of each vector with the angle,  $\theta$  between them. Using this approach, one could say similar vectors – vectors that have a larger dot product – have a small angular separation. This holds true for high dimensional vector spaces. In EKSS, angular separation is used to calculate the angles from one point in a subspace to another point. This is used to quantify the distance between clusters. If the data has a high, positive angular separation, EKSS proved an upper bound to guarantee correct clustering. Specifically, EKSS' probability of clustering points,  $x_i$  and  $x_j$ , correctly is monotonically increasing with the absolute value of the inner product, 9, [2]. Assuming EKSS preserves angular separation, it has a theoretical "recovery guarantee" that is historically unproved with KSS and thresholding subspace clustering. EKSS, being a repetition of KSS, converges to a local minimum.

$$x^T y = |x| |y| \cos \theta \quad (8)$$

$$P(\theta) = \mathbf{P}\{Qx_i, Qx_j \text{ both assigned to } U_1\} \quad (9)$$

### 1.2.3 Relation to Course Material

There are a few core components that relate to the course material. The idea of residuals and projections is familiar, however, clustering them is new. Instead of clustering by the closest subspace, we stated that PCA

told us to cluster using the subspaces that contain the most variance, which we used to unknowingly cluster in Homework 6. Similarly, KSS was generally discussed in the course but not explicitly. The final step of EKSS – spectral clustering – was also covered in the course.

There were a few bits of the paper that I did not catch. First, I did not catch that the subspaces needed to be the true, noise-free subspaces. I ended up reading [your dissertation](#) and [this talk by Dr. Laura Balzano](#). I think a student from this course could understand concepts from these resources except angular separation theory – which seems to be a key component in the recovery guarantees. I still do not understand it.

### 1.3 Paper 3

**Paper Title:** Adaptive Online k-Subspaces with Cooperative Realization

**Student Name:** Will Mijangos

This process attempts to solve the rather typical problem of classification. The potential problems lie in large data sets where alternative solutions scale much faster than the proposed algorithm. Not only that, but it may also be able to solve some uncertainties of the true number of subspaces and their dimensions where these are unknown or variable. Applications of this include any where classification is a problem, image processing, signal processing, data mining, data compression and recovery. The typical form fits a collection of  $k$  subspaces that are a collection of data points and attempts to separate them with a dependance on similarity. This is standard when you have a known number of subspaces and a finite dataset. As for the algorithm, is an online algorithm, which means that it can process the data one point at a time, in contrast to other methods, which require the entire data set to be available before they can be run. It is a greedy algorithm, which means that it iteratively refines the set of  $k$  subspaces by greedily re-initializing each subspace with the remaining data. This ensures that each subspace is updated to reflect the current state of the data set, and that the final set of subspaces is a good fit for the data. It has been shown to be effective on a variety of data sets, and it has been shown to be more robust to noise and outliers than other methods.

The algorithm is initialized with  $k$  randomly chosen points from the data set. In each iteration, each subspace is re-initialized with the remaining data. There are regularization terms that are functional. The first is a penalty scaled relative to cluster size, another is fixed and this allows convergence for large data sets or unknown number of subspaces. The points are then assigned to the subspace that they are closest to. This process is repeated until the algorithm converges.

---

#### Algorithm 3 CoRe Algorithm

---

```

Initialize:  $k$ -Subspaces,  $S$ , with randomly chosen points from the data set,  $\epsilon$ 
error =  $\epsilon + 1$ 
while error >  $\epsilon$  do
    Re-initialize each subspace  $S$  with the remaining data.
    Assign each point to the subspace  $S$  that it is closest to.
end while

```

---

The CoRe algorithm can be initialized with  $k$  randomly chosen points from the data set. However, it has been shown that the algorithm can be improved by initializing the subspaces with the results of a  $k$ -means algorithm.

#### 1.3.1 Complexity

The complexity of the CoRe algorithm is  $\mathcal{O}(k^2 \log(k))$  where  $k$  is the number of bases.

The CoRe algorithm has two tuning parameters:

- $k$ : The number of subspaces.
- $\epsilon$ : The convergence tolerance.

The value of  $k$  should be chosen such that the data set is well-represented by the subspaces. The value of  $\epsilon$  should be chosen such that the algorithm converges quickly.

The main theoretical result of the paper is that CoRe is able to find  $k$  subspaces that are close to the true subspaces, even in the presence of noise and outliers. The paper also shows that CoRe is able to achieve better performance than state-of-the-art algorithms, such as  $k$ -means and  $k$ -svd, on both synthetic and real-world data sets. The experiments show that CoRe is able to find  $k$  subspaces that are closer to the true subspaces, even in the presence of noise and outliers using image data. The conclusion includes a discussion on the implications of the results for the field of machine learning. The results suggest that CoRe could be used to improve the performance of machine learning algorithms that require  $k$  subspaces, such as image clustering and data compression. The main takeaways are that this scales very well and it may adapt to unknown  $k$  (number of subspaces) and  $d$  (dimensions) of those subspaces.

This algorithm directly employs Stochastic Gradient Descent (SGD) which we have covered in class, albeit with a different regularization term and many other factors but that was immediately discernible. The paper also compares results to  $k$ -means and methods using the SVD which were also introduced. Subspace clustering and classification is a recurring problem in this subject field.

## 2 Comparison of Algorithms

### 2.1 Interpretability

The Sparse Subspace Clustering was the easiest to understand and implement. It felt like two broad, open to interpretation steps: get some sparse coefficients, spectral cluster them. I felt like this paper alone could birth a wide range of spin-off topics and experiments. It also had the most theory that felt similar and understandable over one pass. EKSS was the easiest to visualize using Dr. Balzano, however, the unfamiliar theory made it difficult to engage with.

I think SSC was easier to understand because it was an optimization problem as well. EKSS, being geometric, is indirectly related to the course while SSC is directly related. Similarly, CoRe is a happy medium that introduces directly related topics such as stochastic gradient descent, least squares, and general programming operations. It also included a few tangential stretch topics like self-expression, Hessian matrices, and batch computations.

### 2.2 Theoretical Guarantees

Sparse Subspace Clustering offers theoretical guaranteed success if two conditions are met, Eq. 4 and Eq. 6. EKSS has a theoretical guarantee as long as the subspaces are not too close in every direction and the angles are preserved. The probability of success is related to the magnitude of the dot products. While SSC promotes guarantees most related to this course, one could argue that a proper theoretical comparison of SSC and EKSS is related to the trade-offs between SSC Eq. 6 and 4 and EKSS' probability of co-clustering, Eq. 9. This, of course, is comparing an optimization problem to a geometric one. Comparing SSC to CoRe, CoRe seems to be a variation of SSC that uses the Hessian and stochastic gradient descent instead of Eq. 5 and incurs a weaker theoretical guarantee as a result.

### 2.3 Empirical Guarantees

SSC is guaranteed to recover a sparse solution even after the convex relaxation, Eq. 5 [1]. For CoRe, there are guaranteed intersections as the number of subspace dimensions increases with a fixed number of clusters, and it is guaranteed to cluster all points. EKSS is guaranteed to cluster all points as long as the subspaces are separated. From reading papers, SSC seems to be accepted as a generally good baseline. EKSS is demonstrated to outperform SSC with subspace angles  $\theta \in [0.2, 0.6]$ . With other angles, they perform identically. CoRe to be a different version of SSC. Instead of trying to compete for performance, it focuses on computational complexity and scalability.

## 2.4 Computational Complexity

The computational complexity of Sparse Subspace Clustering grows exponential from the subspace’s dimensions – which is also the sparsity level. EKSS boasts a complexity of  $\mathcal{O}(\bar{K}D\bar{d}(D+N))$  where  $\bar{K}$  is the number of randomly generated subspaces with  $\bar{d}$  dimensions, and  $N$  samples of data  $D$  in dimension. CoRe’s most computationally expensive step is  $\mathcal{O}(R\bar{k}\bar{d}^2(D+n_{bs}))$  and is overall  $\mathcal{O}(k^2 \log(k))$  with  $k$  being the number of subspaces. EKSS could easily become the most computationally efficient – especially with low dimensional data. Given the choice, I would select SSC if there number of subspaces is low.

## 3 Algorithm Implementation & Testing

We implementing the Sparse Subspace Clustering (SSC) algorithm on a synthetic data set and two digits from the MNIST handwritten digit data set. It seemed the most related to the Spectral Clustering demo and the recently discussed alternative direction method of multipliers. It also had the most citations which means it is the coolest. There are also a few changes we made to the algorithm. The paper states that, unless the data has previously proven subspaces, its best to perform SSC on the whole data set. This, however, is computationally  $\mathcal{O}(e^d)$  depending on the number of dimensions to define a subspace,  $d$ . For the MNIST handwritten digit dataset,  $d$  could be around 784. In their real world experiments, they briefly discuss how to use a Scree plot to select the number of PCA components by identifying the "knee" from the Scree. The number of components is the sparsity level to span the subspace – which is the number of required coefficients. The alternative method is to recover a complete coefficient matrix and make it sparse using a threshold. For us, this provided worse empirical results.

The link to the Google Colab file can be found [here](#). You must be signed into PSU’s Google Account to access.

### 3.1 Results on Synthetic Data

The synthetic data composed of two orthogonal subspaces: one with basis vector  $[1, 1, 0]$  and the other with basis vector  $[1, 0, 1]$ . A total of  $N$  many samples were randomly scaled with an even split of samples in both subspaces. Additive noise in  $\mathbb{R}^3$  from a normal distribution was added to each sample. The results of the Sparse Subspace Clustering, Figure 1, proved to be unaffected by noisy data.

Increasing the number of samples without changing the number of possible clusters allowed for more unique coefficients – which should improve the performance of clustering. Our results confirmed this intuition, Figure 2.

There are two different parameters worth noting: the weight of the regularizer,  $\lambda$ , and the number of dimensions to consider,  $P$ . Their influence is demonstrated in Figure 3 and 4 respectively. If the number of dimensions is unknown,  $P$  would change to a tolerance,  $\tau$ , where coefficients less than  $\tau$  are set to zero.

The best results came with  $\lambda = 1e - 6$  and  $P = 5$  coefficients, which produced a 99.5% accuracy, Figure 4. With  $N = 400$  and two clusters, the run time was 153.08 seconds.

The run time of the algorithm empirically proved to be at least polynomial in terms of the number of samples,  $N$ , Figure 7.

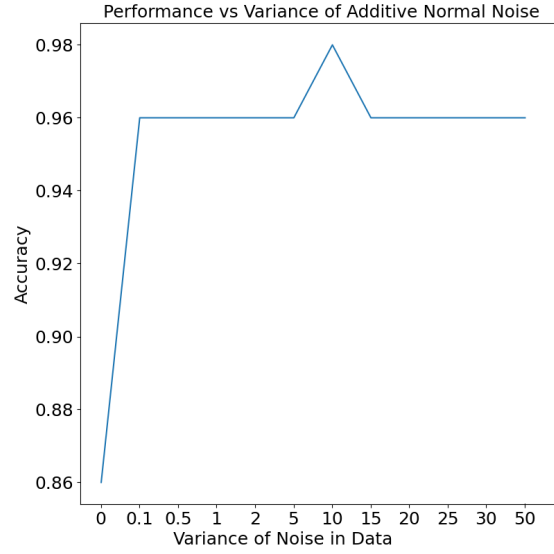


Figure 1: Additive white noise did not affect the performance of SSC.

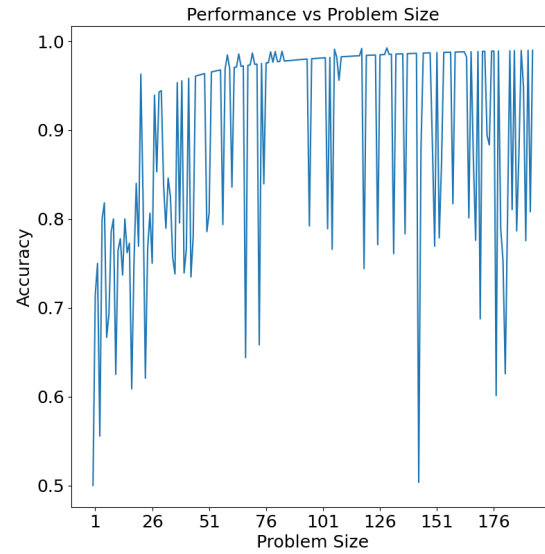


Figure 2: When the number of samples increased and the number of clusters stayed the same, the self-expressiveness of data points improved. There are more possible coefficients to promote tailored linearly dependent combinations to recover difficult data.



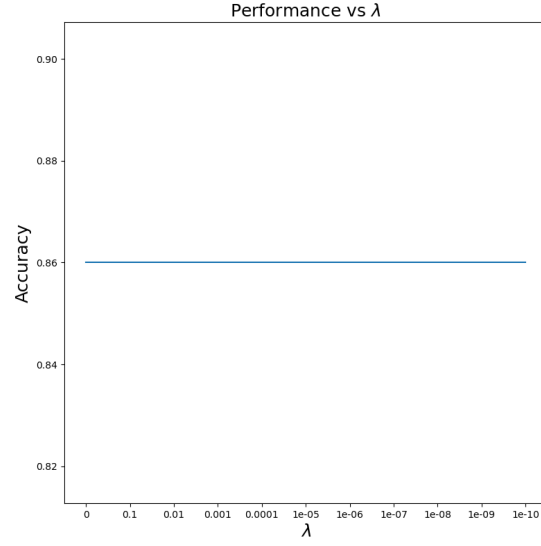


Figure 3: For our synthetic data, the weight of the constraint was uncorrelated with the performance of the SSC. It is assumed that the synthetic data always had some linear combination that perfectly reconstructed all data points. This could indicate that the minimum number of basis vectors is lower than our estimation.

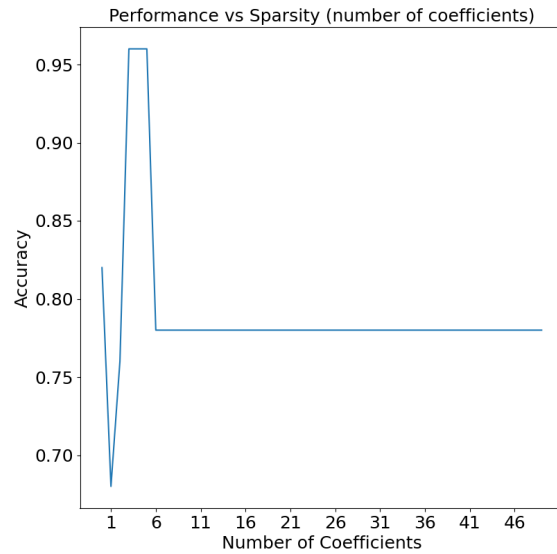


Figure 4: With 50 samples of data, the number of components increased as it approached the true number of basis vectors to span the subspace. As  $P$  advanced beyond, the performance dipped as the resulting basis vectors begin to bleed into each other.

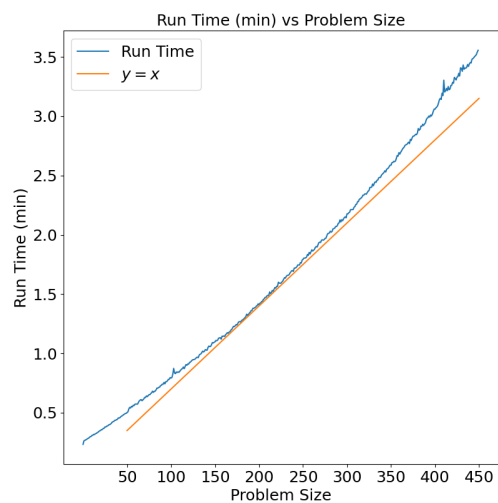


Figure 5: As the number of samples increased, the run time appeared to increase polynomially (blue). A linear line (orange) with a slope determined from the first and final run time is shown for comparison.

### 3.2 Results on Benchmark Data

SSC was implemented using 200 sample from both digit 1 and digit 5 from MNIST's handwritten digits data set [3]. A sweep across  $\lambda$  and  $P$  revealed a local minima for  $\lambda$  while increasing  $P$  decreased performance. As  $P$  increased, the proportion of number of samples compared to the number of coefficients lowered, producing a harder problem to solve. The performance did not match what was expected. SSC was able to differentiate between many faces while our implementation failed to completely differentiate two numbers. In [1], it was known that the faces has subspaces in  $\mathbb{R}^9$ . For our implementation, we empirically concluded to subspaces in  $\mathbb{R}^5$ . Perhaps, we were wrong and missing information that led to incorrect clustering.

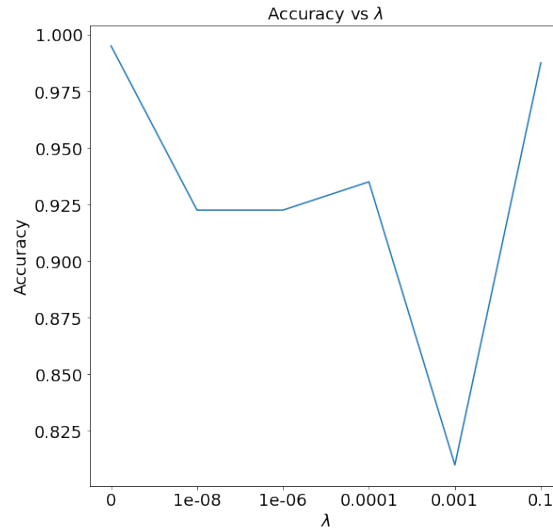


Figure 6: Increasing the weight of the regularizer generally worsened the performance of SSC.

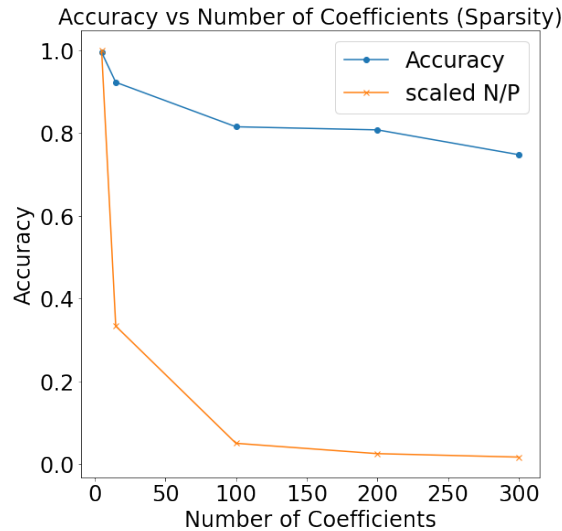


Figure 7: Increasing the number of coefficients per sample shows a decrease in the performance (blue). As the number of coefficients increases, the ratio of coefficients to samples decreases (orange), justifying the worsening accuracy.

## References

- [1] E. Elhamifar and R. Vidal, “Sparse subspace clustering: Algorithm, theory, and applications,” *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 35, pp. 2765–2781, Nov. 2013.
- [2] J. Lipor, D. Hong, Y. S. Tan, and L. Balzano, “Subspace clustering using ensembles of k-subspaces,” *arXiv preprint arXiv:1709.04744*, 2017.
- [3] L. Deng, “The mnist database of handwritten digit images for machine learning research,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.

## A Appendix

This is an appendix where you may include your code.