

Copyright
by
Ali Jalali
2012

The Dissertation Committee for Ali Jalali
certifies that this is the approved version of the following dissertation:

Dirty Statistical Models

Committee:

Sujay Sanghavi, Supervisor

Constantine Caramanis

Inderjit Dhillon

Joydeep Ghosh

Pradeep Ravikumar

Dirty Statistical Models

by

Ali Jalali, B.S.; M.S.

DISSERTATION

Presented to the Faculty of the Graduate School of
The University of Texas at Austin
in Partial Fulfillment
of the Requirements
for the Degree of

DOCTOR OF PHILOSOPHY

THE UNIVERSITY OF TEXAS AT AUSTIN

May 2012

Dedicated to Soltän.

Acknowledgments

I wish to thank all professors, students and staff of The Wireless Networking and Communication Networking Group at UT Austin for their kind helps throughout my PhD. In particular, I would like to thank my adviser, Prof. Sujay Sanghavi, whose constant support was always encouraging me to learn more and without his dedication, this thesis was impossible. I would also like to express my gratitude to Prof. Pradeep Ravikumar for insightful discussions that helped me find my research path. My special thanks to my parents and little sister for their unconditional love and support during my hardest times. I am running out of words when it comes to thanking my friends, my family by choice, for being always with me in difficulties.

Dirty Statistical Models

Publication No. _____

Ali Jalali, Ph.D.

The University of Texas at Austin, 2012

Supervisor: Sujay Sanghavi

In fields across science and engineering, we are increasingly faced with problems where the number of variables or features we need to estimate is much larger than the number of observations. Under such high-dimensional scaling, for any hope of statistically consistent estimation, it becomes vital to leverage any potential structure in the problem such as sparsity, low-rank structure or block sparsity. However, data may deviate significantly from any one such statistical model. The motivation of this thesis is: can we simultaneously leverage more than one such statistical structural model, to obtain consistency in a larger number of problems, and with fewer samples, than can be obtained by single models? Our approach involves combining via simple linear superposition, a technique we term *dirty models*. The idea is very simple: while any one structure might not capture the data, a superposition of structural classes might. Dirty models thus searches for a parameter that can be *decomposed* into a number of simpler structures such as (a) sparse plus block-sparse, (b) sparse plus low-rank and (c) low-rank plus block-sparse. In this thesis, we propose dirty model based algorithms for different problems such as multi-task learning, graph clustering and time-series analysis with latent factors. We analyze these algorithms in terms of the number of observations we need to estimate the variables. These algorithms are based on convex optimization and sometimes they are relatively slow. We provide a class of low-complexity greedy algorithms that not only can solve these optimizations faster, but also guarantee the solution. Other than theoretical results, in each case, we provide experimental results to illustrate the power of dirty models.

Table of Contents

Acknowledgments	v
Abstract	vi
List of Tables	xi
List of Figures	xii
Chapter 1. Introduction	1
1.1 Structure Recovery Techniques	3
1.1.1 Convex Optimization	3
1.1.2 Greedy Algorithms	5
1.2 Studied Dirty Models	5
1.2.1 Sparse + Block Sparse	6
1.2.2 Sparse + Low-Rank	7
1.2.2.1 Graph Clustering	7
1.2.2.2 Time-Series Analysis	8
Chapter 2. A Dirty Model for Multiple Sparse Regression	10
2.1 Introduction: Motivation and Setup	10
2.2 Problem Set-up and Our Method	13
2.2.1 Our Method	14
2.3 Main Results and Their Consequences	14
2.3.1 Sufficient Conditions for Deterministic Designs	16
2.3.2 General Gaussian Designs	18
2.3.3 Quantifying the gain for 2-Task Gaussian Designs	20
2.4 Simulation Results	21
2.4.1 Synthetic Data Simulation	22
2.4.2 Handwritten Digits Dataset	25

2.5	Proof Outline	27
2.5.1	Definitions and Setup	28
2.5.1.1	Towards Identifying Optimal Solution	28
2.5.1.2	Sparse Matrix Setup	29
2.5.1.3	Row-Sparse Matrix Setup	29
2.5.2	Proof Overview	30
2.6	Proofs	34
2.6.1	Proof of Theorem 1	34
2.6.2	Proof of Theorem 2	39
2.6.3	Proof of Theorem 3	48
2.6.4	Proof of Theorem 4	51
2.7	Deterministic Necessary Optimality Conditions	62
2.7.1	Sub-differential of ℓ_1/ℓ_∞ and ℓ_1/ℓ_1 Norms	62
2.7.2	Necessary Conditions	63
2.8	Coordinate Descent Algorithm	68
2.8.1	Correctness of Algorithms	68
Chapter 3. Clustering Partially Observed Graphs		72
3.1	Introduction	72
3.1.1	Related Work	74
3.2	Main Contributions	76
3.3	Proofs	82
3.3.1	Proof of Theorem 5	82
3.3.2	Proof Outline for Theorem 6 and 7	83
3.3.2.1	Preliminaries	83
3.4	Worst Case Analysis	84
3.5	Average Case Analysis	86
3.6	Experimental Results	88
3.7	Additional Notations	90
3.8	Proof of Theorem 6	93
3.8.1	Auxiliary Lemmas	95
3.9	Proof of Theorem 7	96
3.9.1	Auxiliary Lemmas	106

Chapter 4. Graph Clustering using Max-norm Optimization	113
4.1 Introduction	113
4.1.1 Relationship to the Goemans Willimason SDP Relaxation	114
4.1.2 Other Clustering Approaches	115
4.2 Problem Setup	116
4.3 Max-Norm Relaxation	117
4.3.1 Theoretical Guarantee	118
4.3.2 Comparison to Trace-Norm Constrained Clustering . . .	120
4.4 Max-norm + ℓ_1 -norm Optimization	122
4.4.1 Semi-Definite Programming Method	122
4.4.2 Factorization Method	123
4.4.3 Loss Function Method	123
4.4.4 Dual Decomposition Method	124
4.4.5 Numerical Comparison	125
4.5 Tighter Relaxations	126
4.5.1 Single-linkage Post Processing	127
4.5.2 Comparison with Other Algorithms	128
4.6 Proof of Lemma 27	129
4.7 Proof of Lemma 26	129
4.8 Proof of Theorem 8	130
4.8.1 Notation	131
4.8.1.1 Residual Matrix Notations	131
4.8.1.2 Clustering Matrix Notations	131
4.8.2 Sufficient Optimality Conditions	133
4.8.3 Dual Variable Construction	134
Chapter 5. Learning the Dependence Graph of Time Series with Latent Factors	137
5.1 Introduction	137
5.2 Problem Setting and Main Idea	139
5.2.1 Main Idea	140
5.2.2 Identifiability	141
5.2.3 Algorithm	142

5.2.4	High-dimensional setting	143
5.2.5	Related Work	143
5.3	Main Results	144
5.4	Proof of the Theorem	147
5.4.1	Proof Technique	148
5.5	Experimental Results	150
5.5.1	Synthetic Data	150
5.5.2	Stock Market Data	152
5.6	Proof of Lemma 32	153
5.7	Illustrative Example	154
5.8	Proof of Lemma 33	154
5.9	Auxiliary Optimality Lemmas	157
5.10	Concentration Results	161
5.11	Proof of the Continuous Time Theorem	170
Chapter 6.	Greedy Dirty Models	171
6.1	Introduction	171
6.2	Greedy Algorithm for Dirty Model	174
6.3	Theoretical Guarantee	177
6.4	Experimental Results	180
6.4.1	Synthetic Data	180
6.4.2	Handwritten Digits Dataset	182
6.5	Auxiliary Lemmas for Theorem 10	184
6.6	Lemmas on the Stopping Size	188
Chapter 7.	Conclusion	194
7.1	Flexible and Robust High-dimensional Statistics	195
7.2	Graphical Data Modeling in High-dimensions	196
7.3	Information Theoretic Limits of High-dimensional Statistics	197
Bibliography		198
Index		212
Vita		213

List of Tables

2.1	Six different classes of features provided in the dataset. The dynamic ranges are approximate not exact. The dynamic range of different morphological features are completely different. For those 6 morphological features, we provide their different dynamic ranges separately.	25
2.2	Simulation Results for our model, ℓ_1/ℓ_∞ and LASSO.	27
6.1	Handwriting Classification Results for greedy algorithm, dirty model, group LASSO and LASSO. The greedy method provides much better classification errors with simpler models. The greedy model selection is more consistent as the number of samples increases.	184

List of Figures

2.1	Probability of success in recovering the true signed support using dirty model, Lasso and ℓ_1/ℓ_∞ regularizer. For a 2-task problem, the probability of success for different values of feature-overlap fraction α is plotted. As we can see in the regimes that Lasso is better than, as good as and worse than ℓ_1/ℓ_∞ regularizer ((a), (b) and (c) respectively), the dirty model outperforms both of the methods, i.e., it requires less number of observations for successful recovery of the true signed support compared to Lasso and ℓ_1/ℓ_∞ regularizer. Here $s = \lfloor \frac{p}{10} \rfloor$ always.	16
2.2	Verification of the result of the Theorem 3 on the behavior of phase transition threshold by changing the parameter α in a 2-task (n, p, s, α) problem for our method, LASSO and ℓ_1/ℓ_∞ regularizer. The y -axis is $\frac{n}{s \log(p - (2 - \alpha)s)}$, where n is the number of samples at which threshold was observed. Here $s = \lfloor \frac{p}{10} \rfloor$. Our method shows a gain in sample complexity over the entire range of sharing α . The pre-constant in Theorem 3 is also validated.	24
2.3	An instance of images of the ten digits extracted from the dataset	26
3.1	The adjacency matrix of a graph before (a) and after (b) proper reordering (i.e. clustering) of the nodes. The figure (b) is indicative of the matrix as a superposition of a sparse matrix and a low-rank one.	77
3.2	Simulation results for fully observed 1000-node graph with all clusters of the same size. For different cluster sizes K_{\min} and different number of disagreements per node b , we plot the probability of success.	85
3.3	Simulation results for fully observed 1000-node graph with cluster of non-uniform sizes. The graph has clusters of at least size k . For different minimum cluster size K_{\min} and number of disagreement per node b , we plot the probability of success.	86
3.4	Simulation results for partially observed 400-node network with minimum cluster size fixed at $K_{\min} = 60$. Disagreements are placed on each (potential) edge with probability τ , and each edge is observed with probability p_0 . The figure shows the probability of success in recovering the ideal cluster under different τ and p_0 . Brighter colors show higher success.	88

3.5	Simulation results for partially observed 400-node network with fixed probability $\tau = 0.04$ of placing a disagreement, and different K_{\min} and p_0	89
4.1	Theorem 8 guarantee region of the noise level D_{\max} vs the unbalanceness parameter $\frac{1}{k^*} \sum_i \left(\frac{ C_i^* }{ C_{\min} } \right)^2$	119
4.2	Probability of exact clustering recovery for max-norm and trace-norm constrained algorithms under absolute $\ A - K\ _1$ and linear $\sum_{i,j} K_{ij}(1 - 2A_{ij})$ objectives. There are 4 clusters of size 25 for the balanced case and three clusters of size 30 + one cluster of size 10 for the unbalanced case. We consider two cases for each graph; where the affinity matrix is binary and when it is not. We both show the results for simple max-norm relaxation (basic algorithm) and tighter relaxations presented in Section 4.5 (enhanced algorithm). The result shows that max-norm constrained optimization recovers the exact clustering matrix under higher noise regimes better than trace-norm and single-linkage algorithm. Also, the linear objective seems to be performing better than the absolute objective for the clustering problem in most cases.	121
4.3	Comparison of the proposed numerical optimization methods in terms of the sparsity of the solution they provide and the ℓ_1 error of the estimation.	126
4.4	Summary of possible convex relaxations of the set of valid clustering matrices and their relations. Here, $\ \cdot\ _*$ represents the trace (nuclear) norm, $\ \cdot\ _{\infty,2}$ represents the maximum ℓ_2 norm of the rows, “ \geq ” is used for element-wise positiveness and “ \succeq ” is used for positive semi-definiteness. Each double-ended arrow represents the equivalence of two sets. Each single-ended arrow in this figure represents a <i>strict</i> sub-set relation between two sets.	127
4.5	Comparison of our <i>best</i> proposed method which is the linear objective over tight relaxation (followed by a single-linkage algorithm) with trace-norm counterpart, single-linkage algorithm and spectral clustering. Here, we plot the entropy-based distance of the recovered clustering with the underlying true clustering.	128
4.6	Illustration of two alternative clusterings on the same graph with $D_{\max} = \gamma$. Each gray cloud of points is a clique. Each link between two clouds of points connects every points on one cloud to every points on the other cloud.	130

5.1	Probability of success in recovering the true signed support of A^* versus the control parameter Θ (rescaled ηn) with $p = 200$, $r = 10$ and $s = 20$ for different values of η (left), and, with $p = 200$, $s = 20$ and $\eta = 0.01$ for different number of latent time series r (middle), and, with $p = 200$, $r = 10$ and fixed $\eta = 0.01$ for different sparsity sizes s (right).	150
5.2	Comparison of the stock dependencies recovered by Pure LASSO [15] and our algorithm.	152
5.3	Prediction error and model sparsity versus the ratio of the training/testing sample sizes for prediction of the stock price. Prediction error is measured using mean squared error and the model sparsity is the number of non-zero entries divided by the size of \hat{A}	153
6.1	Probability of success in recovering the exact sign support using greedy algorithm, dirty model, Lasso and group LASSO (ℓ_1/ℓ_∞). For a 2-task problem, the probability of success for different values of feature-overlap fraction κ is plotted. Here, we let $s = p/10$ and the values of the parameter and design matrices are i.i.d standard Gaussians. Also, the noise variance is set to be $\sigma = 0.1$. As we can see, greedy method outperforms all methods in the minimum number of samples required for sign support recovery.	181
6.2	Behavior of phase transition threshold versus the parameter κ in a 2-task problem for greedy algorithm, dirty model, LASSO and group LASSO (ℓ_1/ℓ_∞ regularizer). The y-axis is $\Theta = \frac{n}{s \log(p - (2 - \kappa)s)}$, where n is the number of samples at which threshold was observed. Here, we let $s = p/10$ and the values of the parameter and design matrices are i.i.d standard Gaussians. Also, the noise variance is set to be $\sigma = 0.1$. The greedy algorithm shows substantial improvement in terms of the minimum number of samples required for exact sign support recovery over the other methods.	183

Chapter 1

Introduction

In many applications such as pattern recognition, machine vision, bioinformatics, data mining, financial engineering, etc, we face parameter estimation problems where the number of observations is much less than the number of the dimension of the parameter. In such a high-dimensional regime, there is no hope for consistent estimation unless we restrict the parameter to have a certain structure model. In other words, the true parameter lives in a high-dimensional space, but the observed data has much lower intrinsic dimensions. The hope is that few observations suffice to recover the parameter in the low-dimensional sub-space. Common, and recently popular, structure models are

- **Sparsity:** This model suggests that most entries of the unknown parameter are zero and there are only few non-zero elements. This structure has been studied in many areas such as compressed sensing [13] and LASSO [131].
- **Block Sparse:** This model partitions the entries of the unknown parameter into a number of subsets (blocks) and suggests that within each block, either all entries are zero or most of them, maybe all, are non-zero [26, 111].
- **Low-Rank:** This model suggests that the parameter matrix has a low-dimensional row/column space and hence has much lower rank compare to its size [102, 116].
- **Sparse Markov Random Field:** This model suggests that the graphical model associated with a set of random variables is sparse [112, 113]; that is the joint probability distribution of a set of random variables can be factorized such that each factor is the joint probability distribution of a subset of those random variables.

In each of these models, if we knew apriori the exact instance of the overall dimensional structure, inference is easy. The challenge however is finding the correct instance of the structure from a very large number of possibilities. For example, in sparse linear regression, if we know the location of non-zero entries, then we can set the other entries to zero and estimate only non-zero entries, say by solving a convex optimization problem. Similarly, in block sparse model, if we know which blocks are non-zero, or in low-rank model, if we know the singular vectors of the matrix, then we can estimate non-zero entries with very few observations. Hence, recovering the structure of a model is the critical task in high-dimensional setting. Once the structure of the model is recovered, estimation is a fairly easy job.

There is a hidden assumption behind choosing a certain model for the unknown parameter. In fact, we assume that the underlying true data structure falls into one of these models completely, i.e., the data is *clean*. Notice that when we use the term clean data, we do not necessarily mean that there is no noise. For example, if x is a vector, then $Q = xx^T + w$, with additive noise w , can be considered to have low-rank model (which is a clean model). However, in many real applications, we face more complex structures, where the underlying data structure does not fit any model solely. Instead, the data structure model can be expressed as the superposition of a number of simpler models, i.e., the data is *dirty*. For example, we can consider the aforementioned matrix Q to be a superposition of low-rank and sparse models, assuming that the noise has only few non-zero elements. We call such a superposition of simple models a *dirty model*.

Our main idea in this thesis is to analyze dirty models for dirty data and try to recover the structure of each of these models. We have analyzed recovery of different dirty models via convex optimization and greedy algorithms in terms of their sample complexity, structure recovery guarantees and comparison to other clean models. As discussed before, once the structure is revealed, the estimation is easy because the unknown parameter will be restricted to a low-dimensional subspace.

1.1 Structure Recovery Techniques

In this section, we discuss two main structure recovery techniques and their properties. Instances of each of these techniques as well as relevant applications are presented in the consequent chapters of this thesis.

1.1.1 Convex Optimization

The use of convex optimization for structure recovery is very broad and popular. At a high level, given a set of observation D , we would like to estimate a target parameter θ^* . There are two ingredients for a consistent convex optimization based structure recovery algorithm as follows:

- **Loss Function:** We need a convex loss function $\mathcal{L}(\theta; D)$ that given our observation, it assigns a *penalty* to each parameter θ . This loss function must have θ^* as the asymptotically optimal point, i.e., $\mathbb{E}_D [\nabla \mathcal{L}(\theta^*; D)] = 0$. Moreover, the more curvature the loss function around θ^* has, the easier the estimation will be.
- **Regularizer:** In a high-dimensional setting, to get a consistent estimation, we assume that $\theta^* \in \mathcal{C}$ for some structure set \mathcal{C} . Examples of such a set are the set of 1-sparse vectors or the set of 1-rank matrices. We then construct a function $R_{\mathcal{C}}(\theta)$ to be the function that encourages to pick elements from the set \mathcal{C} as opposed to an element from outside, i.e., $R_{\mathcal{C}}(\theta^*) \leq R_{\mathcal{C}}(\theta)$ for any $\theta \notin \mathcal{C}$. This function is often referred to as *regularizer*. Since the set \mathcal{C} and consequently the regularizer $R_{\mathcal{C}}$ are often not-convex, we pick the convex envelope $\bar{R}_{\mathcal{C}}$ to ensure tractability.

Using these two ingredients, we solve the following convex program to get an estimate of θ^* .

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; D) + \lambda \bar{R}_{\mathcal{C}}(\theta).$$

The parameter λ can be chosen via techniques like cross-validation or bootstrapping. This class of algorithms are known to have two main issues:

- **Robustness:** A single observation can arbitrarily change $\hat{\theta}$ and hence, the method is not robust with respect to outliers.
- **Flexibility:** If θ^* represents a multi-modal data, then $\underline{\theta}$ potentially represents an *average* mode which is not a representative of neither modes of the data.

Using our dirty model, we can fight against both issues. Imagine $\theta^* = \theta_1^* + \theta_2^*$ where θ_1^* and θ_2^* come from two structures \mathcal{C}_1 and \mathcal{C}_2 , respectively. To get robustness, we choose \mathcal{C}_1 to be the structure of our data and \mathcal{C}_2 to be the structure of the outliers. Similarly, to get flexibility, we choose \mathcal{C}_1 and \mathcal{C}_2 to represent different modes of the data. Subsequently, we solve the following convex optimization problem:

$$(\hat{\theta}_1, \hat{\theta}_2) = \arg \min_{\theta_1, \theta_2} \mathcal{L}(\theta_1 + \theta_2; D) + \lambda_1 \bar{R}_{\mathcal{C}_1}(\theta_1) + \lambda_2 \bar{R}_{\mathcal{C}_2}(\theta_2).$$

In the robust case, we output $\hat{\theta} = \hat{\theta}_1$ and in the flexible case, we output $\hat{\theta} = \hat{\theta}_1 + \hat{\theta}_2$. In this thesis, we investigate some instances of this problem for both robustness (Chapters 3,4,5) and flexibility (Chapter 2) and show the advantages of this dirty model based algorithm over simple model counterparts.

One of the difficulties with robust dirty model analysis is to characterize different structures so that without any ambiguity they can be distinguished. Recall the example of low-rank and sparse models superposition and notice that a low-rank matrix can be also sparse or conversely, a sparse matrix can be low-rank. Thus, we need to not only reduce the size of the problem by imposing the structure, but also further restricting each structure to be consistently incoherent from each other to get robustness. Without this separation, the problem is not well-defined and characterization of the models is impossible. Motivated by this, we normally impose some conditions that uniquely partitions the space of parameters so that each non-zero element of the space belongs to only one structure.

1.1.2 Greedy Algorithms

Considering the original non-convex regularizer $R_{\mathcal{C}}(\cdot)$, the *quality* of the convex optimization method depends on the tightness of the convex relaxation $\bar{R}_{\mathcal{C}}(\cdot)$. Often to ensure tightness, convex optimization problems require strong assumptions such as irrepresentable condition that imposes strong restrictions on the loss function. Besides, convex optimization methods, although polynomial, are still computationally expensive. In this thesis, we show that greedy algorithms are good candidates to both reduce the computational complexity and relax strong assumptions while maintaining the same guarantees as convex optimization.

The main idea of greedy algorithms is basis pursuit and in particular orthogonal matching pursuit [135, 154] – that is iteratively find the next best “coordinate” in \mathcal{C} and add it to the estimated structure. This simple forward greedy algorithm has been shown to perform as good as the convex optimization. However, it still requires the same strong assumptions. Like many forward algorithms, the beginning steps are far more important than the others and if the first steps are mistaken, then the algorithm performs poorly.

We introduce a forward-backward greedy algorithm to solve sparse + block-sparse dirty model based on the algorithms introduced in [68, 152]. We show this algorithm not only lowers the sample complexity, but also requires only mild assumption of restricted strong convexity [100]. This algorithm greedily picks the best coordinates in the forward step and then removes “bad” coordinates in the backward step. We show both theoretically and empirically outperforms the corresponding convex optimization method.

1.2 Studied Dirty Models

We studied different instances of two major classes of dirty models introduced in the next two sections. These dirty models include sparse+block-sparse and sparse+low-rank models.

1.2.1 Sparse + Block Sparse

Suppose we have multiple linear regression problems, i.e., multi-task learning problem [26]. Here, multiple tasks share some common structure such as sparsity, and estimating these tasks jointly by leveraging this common structure could be more statistically efficient. We have $r > 1$ response variables (tasks), and a common set of p features or covariates. The setting we focus on is where the response variables have *simultaneously sparse* structure: the index set of relevant features for each task is sparse; and there is a large overlap of these relevant features across the different regression problems. Such “simultaneous sparsity” arises in a variety of contexts [133]; indeed, most applications of sparse signal recovery in contexts ranging from graphical model learning, kernel learning, and function estimation have natural extensions to the simultaneous-sparse setting [7, 109, 111].

It is useful to represent the multiple regression parameters via a matrix, where each column corresponds to a task, and each row to a feature. Having simultaneous sparse structure then corresponds to the matrix being largely “block-sparse” – where each row is either all zero or mostly non-zero, and the number of non-zero rows is small. A lot of recent research in this setting has focused on ℓ_1/ℓ_q norm regularizations, for $q > 1$, that encourage the parameter matrix to have such block-sparse structure. Particular examples include results using the ℓ_1/ℓ_∞ norm [101, 136, 151], and the ℓ_1/ℓ_2 norm [89, 104].

Our method searches for a parameter matrix that can be *decomposed* into a row-sparse matrix (corresponding to the overlapping or shared features) and an elementwise sparse matrix (corresponding to the non-shared features). As we show both theoretically and empirically, with this simple fix we are able to leverage any extent of shared features, while allowing disparities in support and values of the parameters, so that we are *always* better than both the Lasso or block-sparse regularizers (at times remarkably so). In Chapter 2, we provide a convex optimization based algorithm and in Chapter 6, we provide a forward-backward greedy algorithm for the same problem.

1.2.2 Sparse + Low-Rank

In many applications such as principal component analysis, often times we want to recover a low-rank matrix and a sparse matrix from their sum. Additive large-magnitude noise potentially can change the rank of the matrix by much. In such a heavily noisy regime, PCA based approaches has poor performance. Thus, recovering the underlying low-rank structure in presence of the noise is a challenging job. We model this problem using dirty models and cast it as a convex optimization problem. We consider graph clustering and time-series analysis as applications of using this dirty model.

1.2.2.1 Graph Clustering

The problem we look at is: given an unweighted graph, partition the nodes, i.e., cluster, so as to minimize the sum of (a) number of edges between endpoints that are in different partitions, and (b) the number of missing edges between endpoints in the same partition. This is one particular non statistical application of sparse and low rank matrix separation. Given a clustering, we call edges of type (a) or (b) “disagreements”; we are thus interested in optimal clusterings – those that minimize the number of disagreements. This formulation has the advantage of not requiring an external parametric input of how many clusters there should be in the final solution; this is fully determined by the data at hand.

This problem, as formulated above, is exactly the same as the problem of correlation clustering, first proposed by Bansal, Blum and Chawla [10]. Specifically, edges that [10] labels as “+” are those that are “present in the graph” for us, and those labelled as “-” in [10] are those that are “missing”. They consider the problem with +/- labels on the complete graph. [10] showed that finding the exact optimum is NP-complete; they then proceed to provide a constant-factor approximation algorithm for minimizing the number of disagreements.

We take an alternative approach to the problem; instead of looking for an approximation that holds for all problem inputs, we provide algorithms that either (a) yields the optimum (i.e. disagreement minimizing) clustering, or (b) generates a FAILURE flag, yielding no clustering. Our algorithm is

based on using recently developed matrix uncertainty principles [22, 31] - that (certain classes of) matrices cannot be simultaneously sparse and low-rank. For these matrices, we can *exactly* recover the sparse matrix (\mathbf{B}) and low-rank matrix (\mathbf{K}) given only their sum ($\mathbf{B}+\mathbf{K}$), using convex optimization. In the graph context, the low-rank matrix corresponds to the cliques that would be an ideal input corresponding to the optimal clustering, the composite matrix represents the given data, and the sparse matrix represents the disagreements. Existing results on sparse and low-rank matrix decompositions provide weak guarantees for the graph clustering case; we consider two types of relaxation based on nuclear norm in Chapter 3 and max-norm in Chapter 4 and provide much stronger guarantees.

1.2.2.2 Time-Series Analysis

Suppose we have a non-stationary random vector whose mean and variance evolves over the time and we observe some entries of this vector for a finite time. The evolution of the mean and variance of each entry depends on both observed and unobserved (latent) entries. The problem we would like to investigate is that whether or not it is possible to learn the evolution process of the mean and variance of the observed entries. We consider linear stochastic first-order evolution model $\dot{x}(t) = Ax(t) + \dot{B}(t)$, where x is the random vector, A is the evolution coefficients matrix and $B(t)$ is the standard Brownian motion noise, and learn the matrix A ; the problem that is investigated for the case when there is no latent variable in [15].

We formulate this problem as a sparse plus low-rank dirty model of $A = S + L$, where, S captures the effect of observed variables on each other and L captures the effect of latent variables on observed entries. Here, we assume that few observed entries affect each observed entry and also each observed entry affects only few observed entries, and hence, S is sparse. Moreover, we show that if the number of latent variables is less than observed entries, then the matrix L is low-rank. Unlike most cases of studies in sparse and low-rank decomposition, here, the focus is to learn the structure of the sparse matrix S rather than the low-rank matrix L .

There are similarities between this problem and learning Gaussian graphical models with latent variables [28]. However, at least in [28] the focus has

been to recover the number of latent variables, which is the rank of L , and more importantly, in graphical model learning, the assumption is that the samples are independent and law of large numbers can be applied easily. We need to provide much more subtle analysis since because of the time evolution, our samples are highly correlated and hence, the usual concentration results (of independent variables) cannot be applied. In Chapter 5, we investigate this problem and provide guarantees of success for our method. Further, we apply our method to stock market prediction problem and observe significant improvements.