

## DOCTORAL THESIS

### Variable selection and dimension reduction in high-dimensional regression

Wang, Tao

*Date of Award:*  
2013

[Link to publication](#)

#### General rights

Copyright and intellectual property rights for the publications made accessible in HKBU Scholars are retained by the authors and/or other copyright owners. In addition to the restrictions prescribed by the Copyright Ordinance of Hong Kong, all users and readers must also observe the following terms of use:

- Users may download and print one copy of any publication from HKBU Scholars for the purpose of private study or research
- Users cannot further distribute the material or use it for any profit-making activity or commercial gain
- To share publications in HKBU Scholars with others, users are welcome to freely distribute the permanent URL assigned to the publication

# Variable Selection and Dimension Reduction in High-dimensional Regression

WANG Tao

A thesis submitted in partial fulfillment of the requirements  
for the degree of  
Doctor of Philosophy

Principal Supervisor: Professor ZHU Lixing

Hong Kong Baptist University

August 2013

# Abstract

In this thesis we develop some new methods for variable selection and dimension reduction in regression with high-dimensional data.

The first part of the thesis introduces a linear least squares framework for parameter estimation and variable selection in high-dimensional semi-parametric models. Firstly, we consider a general class of models with a single-index structure. We show that the coefficient vector from the least squares fit of any transformation of the response variable on the predictor vector can be used to identify the single-index parameter. Building upon the least squares formulation, we propose the non-convex penalized least squares estimation. We prove the oracle property of SCAD and the minimax concave penalty estimator when the number of predictors grows at some polynomial rate of the sample size. Secondly, we consider an additive multiple-index model in which each component function has a single-index structure, and of which the partially linear single-index model is a special case. Somewhat surprisingly, we show that all index vectors can be recovered through a single least squares coefficient vector, extending our previous result on the identifiability of single-index models. We propose the SCAD-penalized least squares estimation and establish the oracle property in sparse high-dimensional settings. As an application, for partially linear single-index models we develop a new two-stage estimation procedure that is iterative-free and easily implemented.

In the second part of the thesis, we establish a connection, from the viewpoint of dimension reduction, between classification and regression. Linear discriminant analysis (LDA) is a standard tool both for classification and dimension reduction. It is well-known that LDA can be re-formulated as a regression problem via optimal scoring. We show that sliced inverse regression (SIR), an innovative and effective method for dimension reduction in regression, can also be recast as an optimal scoring problem. Motivated by this regression interpretation of SIR, we propose a sparse and thus interpretable sufficient dimension reduction method that is applicable to high-dimensional settings. We also propose an alternating minimization algorithm for

efficiently solving the optimization problem.

In the third part, we consider the problem of semi-parametric variable selection for multiple-index regression models. We propose a knowledge-based method that takes into account predictor group information by combining a group-wise dimension reduction method and LASSO. Interestingly, the estimation via a LASSO penalty behaves like an adaptive LASSO estimation. We show that the new estimator is consistent in variable selection while retaining the root- $n$  consistency. We use a Bayesian information criterion (BIC) type criterion to select the optimal tuning parameter. Moreover, we derive the consistency of the resulting BIC-type selector.

**Keywords:** Dimension reduction; High dimensionality; LASSO; Minimum average variance estimation; Multiple-index models; SCAD; Single-index models; Sliced inverse regression; Variable selection.

# Table of Contents

<b>Declaration</b>	<b>i</b>
<b>Abstract</b>	<b>ii</b>
<b>Acknowledgements</b>	<b>iv</b>
<b>Table of Contents</b>	<b>v</b>
<b>List of Tables</b>	<b>viii</b>
<b>List of Figures</b>	<b>x</b>
<b>Chapter 1 Introduction</b>	<b>1</b>
1.1 Variable selection in regression . . . . .	2
1.1.1 Parametric regression models . . . . .	2
1.1.2 Semi-parametric regression models . . . . .	5
1.2 Sufficient dimension reduction in regression . . . . .	7
1.2.1 Sliced inverse regression . . . . .	8
1.2.2 Minimum average variance estimation . . . . .	9
1.2.3 Model-free variable selection . . . . .	11
1.3 Outline of the thesis . . . . .	12
<b>Chapter 2 Non-convex Penalized Estimation in High-dimensional Models With a Single-index Structure</b>	<b>15</b>

2.1	Introduction . . . . .	15
2.2	Methodology and main results . . . . .	17
2.2.1	Index estimation and asymptotics . . . . .	17
2.2.2	A distribution transformation . . . . .	20
2.2.3	Model-free variable selection in high dimensions . . . . .	21
2.3	Numerical studies . . . . .	23
2.4	Appendix . . . . .	31

**Chapter 3 Estimation of a Group-wise Additive Multiple-index Model  
and Its Applications 40**

3.1	Introduction . . . . .	40
3.2	Identifiability and estimation . . . . .	42
3.3	Applications . . . . .	45
3.3.1	Partially linear single-index models: A two-stage estimation procedure . . . . .	45
3.3.2	Predictor selection for large- $d$ -small- $n$ problems . . . . .	49
3.4	Simulation study . . . . .	51
3.5	Appendix . . . . .	54

**Chapter 4 Sparse Sufficient Dimension Reduction Using Optimal Scor-  
ing 68**

4.1	Introduction . . . . .	68
4.2	Methodology and main results . . . . .	70
4.2.1	Linear discriminant analysis by optimal scoring . . . . .	70
4.2.2	Sliced inverse regression by optimal scoring . . . . .	71
4.2.3	Sparse dimension reduction via the elastic net approach . . . . .	74
4.3	Simulation studies . . . . .	75
4.3.1	Computation and tuning parameter selection . . . . .	75
4.3.2	Numerical results . . . . .	76

4.3.3	Leukemia data . . . . .	79
4.4	Discussion . . . . .	84
<b>Chapter 5 Shrinkage Estimation and Variable Selection for Multiple-</b>		
	<b>index Models</b>	<b>85</b>
5.1	Introduction . . . . .	85
5.2	Methodology . . . . .	90
5.2.1	A short review . . . . .	90
5.2.2	Shrinkage group-wise minimum average variance estimation . . . . .	92
5.2.3	Asymptotic theory . . . . .	94
5.3	Numerical studies . . . . .	98
5.3.1	Simulation studies . . . . .	98
5.3.2	Baseball hitters' salary data . . . . .	103
5.4	Appendix . . . . .	109
	<b>Bibliography</b>	<b>119</b>
	<b>Curriculum Vitae</b>	<b>131</b>