

Sparse CCA using Lasso

Anastasia Lykou & Joe Whittaker

Department of Mathematics and Statistics,
Lancaster University

July 23, 2008

Outline

- 1 Introduction**
 - Motivation
- 2 CCA**
 - Definition
 - CCA as least squares problem
- 3 Lasso**
 - Definition
 - Lasso algorithms
 - The Lasso algorithms contrasted
- 4 Sparse CCA**
 - SCCA
 - Algorithm for SCCA
 - Example
- 5 Summary**

Motivation

- **SCCA**

 - improve the interpretation of CCA

 - sparse principal component analysis (SCoTLASS by Jolliffe et al. (2003) and SPCA by Zou et al. (2004))

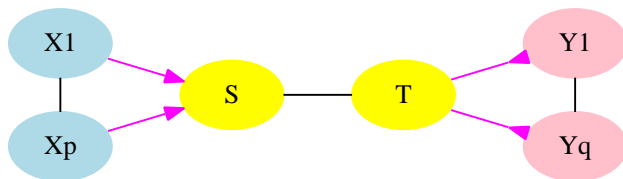
 - interesting data sets (market basket analysis)

- **Sparsity**

 - shrinkage and model selection simultaneously

 - (may reduce the prediction error, can be extended to high-dimensional data sets)

Canonical Correlation Analysis



- seek linear combinations $S = \alpha^T X$ and $T = \beta^T Y$ such that $\rho = \max_{\alpha, \beta} \text{corr}(S, T)$
- S, T are the canonical variates
- α, β are called conical loadings
- Standard solution through eigen decomposition.

1st dimension

Theorem 1

Let α, β be p, q dimensional vectors, respectively.

$$(\hat{\alpha}, \hat{\beta}) = \operatorname{argmin}_{\alpha, \beta} \left\{ \operatorname{var}(\alpha^T \mathbf{X} - \beta^T \mathbf{Y}) \right\},$$

$$\text{subject to } \alpha^T \operatorname{var}(\mathbf{X}) \alpha = \beta^T \operatorname{var}(\mathbf{Y}) \beta = 1.$$

Then $\hat{\alpha}, \hat{\beta}$ are proportional to the first dimensional ordinary canonical loadings.

2nd dimension

Theorem2

Let α, β be p, q dimensional vectors.

$$(\hat{\alpha}, \hat{\beta}) = \underset{\alpha, \beta}{\operatorname{argmin}} \left\{ \operatorname{var}(\alpha^T \mathbf{X} - \beta^T \mathbf{Y}) \right\},$$

$$\text{st } \alpha^T \operatorname{var}(\mathbf{X}) \alpha = \beta^T \operatorname{var}(\mathbf{Y}) \beta = 1 \text{ and}$$

$$\alpha_1^T \operatorname{var}(\mathbf{X}) \alpha = \beta_1^T \operatorname{var}(\mathbf{Y}) \beta = 0$$

where α_1, β_1 are the first canonical loadings.

Then, $\hat{\alpha}, \hat{\beta}$ are proportional to the second dimensional ordinary canonical loadings.

The theorems establish an **Alternating Least Squares algorithm for CCA**.

2nd dimension

Theorem2

Let α, β be p, q dimensional vectors.

$$(\hat{\alpha}, \hat{\beta}) = \underset{\alpha, \beta}{\operatorname{argmin}} \left\{ \operatorname{var}(\alpha^T \mathbf{X} - \beta^T \mathbf{Y}) \right\},$$

$$\text{st } \alpha^T \operatorname{var}(\mathbf{X}) \alpha = \beta^T \operatorname{var}(\mathbf{Y}) \beta = 1 \text{ and}$$

$$\alpha_1^T \operatorname{var}(\mathbf{X}) \alpha = \beta_1^T \operatorname{var}(\mathbf{Y}) \beta = 0$$

where α_1, β_1 are the first canonical loadings.

Then, $\hat{\alpha}, \hat{\beta}$ are proportional to the second dimensional ordinary canonical loadings.

The theorems establish an **Alternating Least Squares algorithm for CCA**.

ALS for CCA

Let the objective function be

$$Q(\alpha, \beta) = \text{var}(\alpha^T \mathbf{X} - \beta^T \mathbf{Y})$$

subject to $\alpha^T \text{var}(\mathbf{X})\alpha = \beta^T \text{var}(\mathbf{Y})\beta = 1.$

Q is continuous with closed and bounded domain $\Rightarrow Q$ attains its infimum

ALS algorithm

- Given $\hat{\alpha}$
- $\hat{\beta} = \arg \min_{\beta} Q(\hat{\alpha}, \beta)$ subject to $\text{var}(\beta^T \mathbf{Y}) = 1$
- Given $\hat{\beta}$
- $\hat{\alpha} = \arg \min_{\alpha} Q(\alpha, \hat{\beta})$ subject to $\text{var}(\alpha^T \mathbf{X}) = 1$

Q decreases over the iterations and is bounded from below $\Rightarrow Q$ converges.

ALS for CCA

Let the objective function be

$$Q(\alpha, \beta) = \text{var}(\alpha^T \mathbf{X} - \beta^T \mathbf{Y})$$

subject to $\alpha^T \text{var}(\mathbf{X})\alpha = \beta^T \text{var}(\mathbf{Y})\beta = 1.$

Q is continuous with closed and bounded domain $\Rightarrow Q$ attains its infimum

ALS algorithm

- Given $\hat{\alpha}$
- $\hat{\beta} = \arg \min_{\beta} Q(\hat{\alpha}, \beta)$ subject to $\text{var}(\beta^T \mathbf{Y}) = 1$
- Given $\hat{\beta}$
- $\hat{\alpha} = \arg \min_{\alpha} Q(\alpha, \hat{\beta})$ subject to $\text{var}(\alpha^T \mathbf{X}) = 1$

Q decreases over the iterations and is bounded from below $\Rightarrow Q$ converges.

ALS for CCA

Let the objective function be

$$Q(\alpha, \beta) = \text{var}(\alpha^T \mathbf{X} - \beta^T \mathbf{Y})$$

subject to $\alpha^T \text{var}(\mathbf{X})\alpha = \beta^T \text{var}(\mathbf{Y})\beta = 1.$

Q is continuous with closed and bounded domain $\Rightarrow Q$ attains its infimum

ALS algorithm

- Given $\hat{\alpha}$
- $\hat{\beta} = \arg \min_{\beta} Q(\hat{\alpha}, \beta)$ subject to $\text{var}(\beta^T \mathbf{Y}) = 1$
- Given $\hat{\beta}$
- $\hat{\alpha} = \arg \min_{\alpha} Q(\alpha, \hat{\beta})$ subject to $\text{var}(\alpha^T \mathbf{X}) = 1$

Q decreases over the iterations and is bounded from below $\Rightarrow Q$ converges.

Lasso (least absolute shrinkage and selection operator)

- Introduced by Tibshirani (1996)
- Imposes the L_1 norm on the linear regression coefficients.

Lasso

$$\hat{\beta}_{lasso} = \operatorname{argmin}_{\beta} \left\{ \operatorname{var}(\mathbf{Y} - \beta^T \mathbf{X}) \right\}$$

$$\text{subject to } \sum_{j=1}^p |\beta_j| \leq t$$

- The L_1 norm properties shrink the coefficients towards zero and exactly to zero if t is small enough.

Lasso algorithms available in the literature

- **Lasso by Tibshirani**

Expresses the problem as a least squares problem with 2^p inequality constraints

Adapts the NNLS algorithm

- **Lars-Lasso**

A modified version of Lars algorithm introduced by Efron et al. (2004)

Lasso estimates are calculated such that the angle between the active covariates and the residuals is always equal.

Proposed algorithm

Lasso with positivity constraints

Suppose that the sign of the coefficients does not change during shrinkage of the coefficients

Positivity Lasso

$$\hat{\beta}_{lasso} = \underset{\beta}{\operatorname{argmin}} \left\{ \operatorname{var}(\mathbf{Y} - \beta^T \mathbf{X}) \right\}$$

subject to $s_0^t \beta \leq t$ and $s_{0j} \beta_j \geq 0$ for $i = 1 \dots, p$

where s_0 is the sign of the OLS estimate.

- simple algorithm, but quite general
- restricted version of Lasso algorithms, since the sign of the coefficients cannot change
- up to $p + 1$ constraints imposed, $\ll 2^p$ constraints of Tibshirani's Lasso

Numerical solution

The solution is given through quadratic programming methods,

Positivity Lasso solution

$$\hat{\beta} = b_0 - \lambda \text{var}(X)^{-1} s_0 + \text{var}(X)^{-1} \text{diag}(s_0) \mu$$

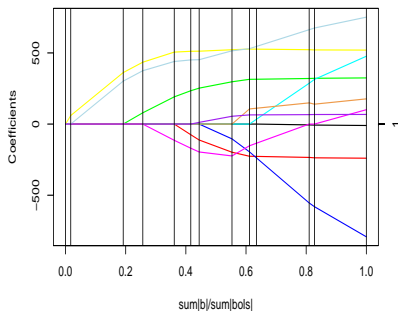
- b_0 is the OLS estimate.
- λ is the shrinkage parameter and there is a one to one correspondence between the λ and t
- μ is zero for active and positive for nonactive coefficients
- parameters λ and μ are calculated satisfying the KKT conditions under the positivity constraints

The Lasso algorithms contrasted

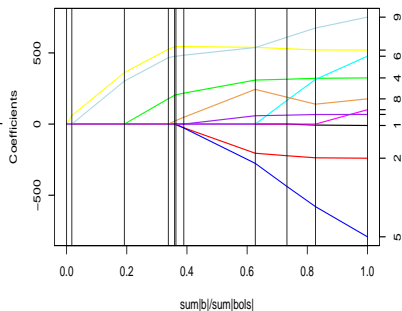
Diabetes data set

- 442 observations
- age, sex, body mass index, average blood pressure and six blood serum measurements
- disease progression one year after baseline

Lars - Lasso



Positivity Lasso



Simulation studies

We simulate 200 data sets consisting of 100 observations each from the following model,

$$\mathbf{Y} = \beta^T \mathbf{X} + \sigma \epsilon, \quad \text{corr}(X_i, X_j) = \rho^{|i-j|}$$

Dataset	n	p	β	σ	ρ
1	100	8	$(3, 1.5, 0, 0, 2, 0, 0, 0)^T$	3	0.50
2	100	8	$(3, 1.5, 0, 0, 2, 0, 0, 0)^T$	3	0.90
3	100	8	$0.85^{\forall j}$	3	0.50
4	100	8	$(5, 0, 0, 0, 0, 0, 0, 0)^T$	2	0.50

Table: Proportions of the cases the correct model selected.

Dataset	Tibs-Lasso	Lars-Lasso	Pos-Lasso
1	0.06	0.13	0.14
2	0.02	0.04	0.04
3	0.84	0.89	0.87
4	0.09	0.19	0.19

Table: Proportions of agreement between Pos-Lasso and

Dataset	Tibs-Lasso	Lars-Lasso
1	0.76	0.83
2	0.63	0.65
3	0.95	0.98
4	0.77	0.78

ALS for CCA and Lasso

First dimension

Given the canonical variate $T = \beta^T \mathbf{Y}$,

$$\begin{aligned} \hat{\alpha} &= \arg \min_{\alpha} \{ \text{var}(T - \alpha^T \mathbf{X}) \} \\ \text{st } \text{var}(\alpha^T \mathbf{X}) &= 1 \text{ and } \|\alpha\|_1 \leq t \end{aligned}$$

We seek an algorithm solving this optimization problem

or

Modify the Lasso algorithm in order to incorporate the equality constraint.

ALS for CCA and Lasso

First dimension

Given the canonical variate $T = \beta^T \mathbf{Y}$,

$$\begin{aligned} \hat{\alpha} &= \arg \min_{\alpha} \{ \text{var}(T - \alpha^T \mathbf{X}) \} \\ \text{st } \text{var}(\alpha^T \mathbf{X}) &= 1 \text{ and } \|\alpha\|_1 \leq t \end{aligned}$$

We seek an algorithm solving this optimization problem

or

Modify the Lasso algorithm in order to incorporate the equality constraint.

ALS for CCA and Lasso

First dimension

Given the canonical variate $T = \beta^T \mathbf{Y}$,

$$\begin{aligned} \hat{\alpha} &= \arg \min_{\alpha} \{ \text{var}(T - \alpha^T \mathbf{X}) \} \\ \text{st } \text{var}(\alpha^T \mathbf{X}) &= 1 \text{ and } \|\alpha\|_1 \leq t \end{aligned}$$

We seek an algorithm solving this optimization problem

or

Modify the Lasso algorithm in order to incorporate the equality constraint.

ALS for CCA and Lasso

- **Tibshirani's Lasso**
NNLS algorithm cannot incorporate the equality constraint
- **Lars Lasso**
the equality constraint violates the equiangular condition
- **Positivity Lasso**
by additionally imposing positivity constraints the above optimization problem can be solved.

ALS for CCA and Lasso

- **Tibshirani's Lasso**
NNLS algorithm cannot incorporate the equality constraint
- **Lars Lasso**
the equality constraint violates the equiangular condition
- **Positivity Lasso**
by additionally imposing positivity constraints the above optimization problem can be solved.

ALS for CCA and Lasso

- **Tibshirani's Lasso**
NNLS algorithm cannot incorporate the equality constraint
- **Lars Lasso**
the equality constraint violates the equiangular condition
- **Positivity Lasso**
by additionally imposing positivity constraints the above optimization problem can be solved.

SCCA with positivity

First dimension

$$\min_{\alpha} \left\{ \text{var}(T - \alpha^T \mathbf{X}) \right\} \text{ st } \alpha^T \text{var}(\mathbf{X})\alpha = 1,$$

and $s_0^T \alpha \leq t, \quad s_{0j} \alpha_j \geq 0 \quad \text{for } j = 1, \dots, p$

- The entire Lasso path is derived by considering KKT conditions.
- Cross-validation methods select the shrinkage level applied.
- α_{sp} and β_{sp} for each set of variables are derived alternately until the $\text{corr}(S_{sp}, T_{sp})$ converges to its maximum.

SCCA with positivity

First dimension

$$\min_{\alpha} \left\{ \text{var}(T - \alpha^T \mathbf{X}) \right\} \text{ st } \alpha^T \text{var}(\mathbf{X})\alpha = 1,$$

and $\mathbf{s}_0^T \alpha \leq t, \quad \mathbf{s}_{0j} \alpha_j \geq 0 \quad \text{for } j = 1, \dots, p$

- The entire Lasso path is derived by considering KKT conditions.
- Cross-validation methods select the shrinkage level applied.
- α_{sp} and β_{sp} for each set of variables are derived alternately until the $\text{corr}(S_{sp}, T_{sp})$ converges to its maximum.

SCCA with positivity

First dimension

$$\min_{\alpha} \left\{ \text{var}(T - \alpha^T \mathbf{X}) \right\} \text{ st } \alpha^T \text{var}(\mathbf{X})\alpha = 1,$$

and $s_0^T \alpha \leq t, \quad s_{0j} \alpha_j \geq 0 \quad \text{for } j = 1, \dots, p$

- The entire Lasso path is derived by considering KKT conditions.
- Cross-validation methods select the shrinkage level applied.
- α_{sp} and β_{sp} for each set of variables are derived alternately until the $\text{corr}(S_{sp}, T_{sp})$ converges to its maximum.

SCCA with positivity

Second dimension

$$\min_{\alpha} \left\{ \text{var}(T - \alpha^T \mathbf{X}) \right\} \quad \text{st} \quad \alpha^T \text{var}(\mathbf{X}) \alpha = 1, \quad \alpha_1^T \text{var}(\mathbf{X}) \alpha = 0,$$

$$\text{and} \quad s_0^T \alpha \leq t, \quad s_{0j} \alpha_j \geq 0 \quad \text{for} \quad j = 1, \dots, p$$

where α_1 is the first dimensional loading.

- Cross-validation methods select the shrinkage level.
- Again alternating algorithm derives the second dimensional canonical loadings

SCCA with positivity

Second dimension

$$\min_{\alpha} \left\{ \text{var}(T - \alpha^T \mathbf{X}) \right\} \quad \text{st} \quad \alpha^T \text{var}(\mathbf{X}) \alpha = 1, \quad \alpha_1^T \text{var}(\mathbf{X}) \alpha = 0,$$
$$\text{and} \quad s_0^T \alpha \leq t, \quad s_{0j} \alpha_j \geq 0 \quad \text{for} \quad j = 1, \dots, p$$

where α_1 is the first dimensional loading.

- Cross-validation methods select the shrinkage level.
- Again alternating algorithm derives the second dimensional canonical loadings

SCCA with positivity

Second dimension

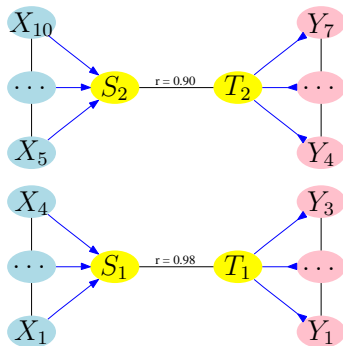
$$\min_{\alpha} \left\{ \text{var}(T - \alpha^T \mathbf{X}) \right\} \quad \text{st} \quad \alpha^T \text{var}(\mathbf{X}) \alpha = 1, \quad \alpha_1^T \text{var}(\mathbf{X}) \alpha = 0,$$
$$\text{and} \quad s_0^T \alpha \leq t, \quad s_{0j} \alpha_j \geq 0 \quad \text{for} \quad j = 1, \dots, p$$

where α_1 is the first dimensional loading.

- Cross-validation methods select the shrinkage level.
- Again alternating algorithm derives the second dimensional canonical loadings

Simulations

We simulate 300 observations of the following model.



Example

Simulations

Variable	1st dim		2nd dim	
	CCA	SCCA	CCA	SCCA
X_1	0.229	0.248	0.122	0.056
X_2	0.350	0.366	-0.052	0
X_3	0.337	0.341	0.027	0
X_4	0.304	0.298	0.114	0
X_5	0.135	0.014	0.198	0.208
X_6	-0.037	0	0.381	0.472
X_7	-0.052	0	0.212	0.183
X_8	-0.052	0	0.205	0.266
X_9	-0.111	0	0.166	0.177
X_{10}	-0.019	0	0.168	0
Y_1	0.402	0.419	0.112	0.014
Y_2	0.460	0.444	-0.018	0
Y_3	0.309	0.325	0.085	0
Y_4	-0.018	0	0.279	0.421
Y_5	0.032	0	0.183	0.008
Y_6	-0.113	-0.028	0.395	0.361
Y_7	-0.089	-0.025	0.384	0.427
ρ	0.745	0.737	0.654	0.638
RdX (%)	14.2	13.9	13.4	12.6
RdY (%)	16.6	16.4	15	14
Var.Ext of X (%)	25.7	25.5	31.4	30.9
Var.Ext of Y (%)	30	30.1	35	33.8

Summary

Extra work

- Sparse CCA without positivity constraints
using Lars-Lasso algorithm

Further work

- Compare the performance of SCCA with and without positivity constraints
- Bayesian model selection
 - Imposing different Lasso penalties
 - Using GVS, Dellaportas et al. (2002)
 - Bayesian version of the SCCA

Literature

- Dellaportas, P., Forster, J., and Ntzoufras, I. (2002). On bayesian model and variable selection using mcmc. *Statistics and Computing*, 12:27–36.
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004). Least angle regression. *Annals of Statistics*, 32:407–499.
- Jolliffe, I., Trendafilov, N., and Uddin, M. (2003). A modified principal component technique based on the lasso. *Journal of Computational and Graphical Statistics*, 12(3):531–547.
- Lawson, C. and Hanson, R. (1974). *Solving Least Square Problems*. Prentice Hall, Englewood Cliffs, NJ.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *J. Royal. Statist. Soc. B.*, 58:267–288.
- Zou, H., Hastie, T., and Tibshirani, T. (2004). Sparse principal component analysis. *to appear, JCGS*.

SCCA

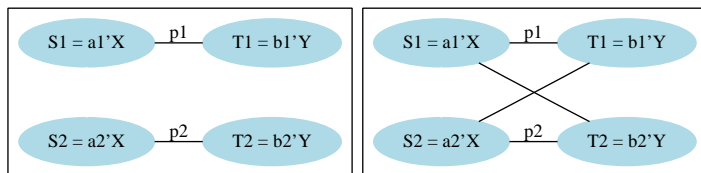


Figure: CCA and SCCA with positivity

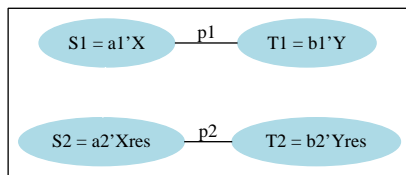


Figure: SCCA without positivity

NNLS

Lawson and Hanson (1974) define the following problems,

LSI problem

LSI problem: $\min_{\beta} \|\mathbf{Y} - \beta^T \mathbf{X}\|$ subject to $\mathbf{G}\beta \geq \mathbf{h}$

NNLS problem: $\min_{\beta} \|\mathbf{Y} - \beta^T \mathbf{X}\|$ subject to $\beta \geq 0$

LDP problem: $\min_{\beta} \|\beta^T\|$ subject to $\mathbf{G}\beta \geq \mathbf{h}$

LSI is equivalent to Lasso \rightarrow LDP \rightarrow NNLS