

Computational Tricks

Simon Rückinger

Seminar:
Modeling, Simulation and Inference of Complex
Biological Systems
(Dr. Korbinian Strimmer)
SoSe 06

9. Juni 2006

Objective

Computationally efficient algorithms for quadratic regularisation

Starting point

- A set of gene expression arrays \mathbf{X}
→ large p (columns), small n (rows).
- Output value \mathbf{y} we want to predict from the expression values.

Idea

Decompose \mathbf{X} in a tricky way and reduce computational effort (Matrix Inversion).

1 / 31

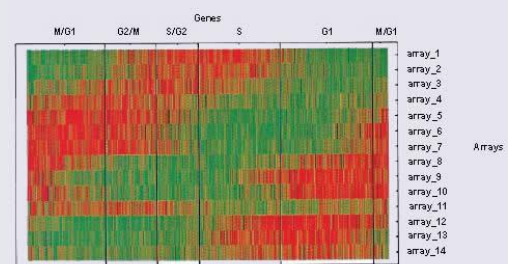
Table of contents

- 1 Motivation
- 2 Singular value decomposition
- 3 Examples
- 4 Summary

Data

Microarray data (p genes, n arrays)

$$\mathbf{X} = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix}$$



2 / 31

3 / 31

What can we do with gene expression data?

Output value

Usually there exists an additional measure \mathbf{y} for each array (individual, time,...). This could be

- cancer class
- biological species
- survival time
- any other quantitative/qualitative measure

Making predictions

A typical goal for a statistician would be to find a model which

- shows the association between gene expression and the output value \mathbf{y}
- is able to make further predictions for \mathbf{y} .

Predicting \mathbf{y} via linear regression

Ordinary least squares

Linear regression could be the model of choice for predicting \mathbf{y} as a continuous measure. To get the model parameters we simply minimize the sum of squared differences between observed (y_i) and predicted values ($\mathbf{x}_i^T \beta$).

$$\min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2$$

Solution for classic linear regression

This leads to the well-known solution

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

4 / 31

5 / 31

Linear regression does not work for $p > n$

Normal equations

When $p > n$, the so-called normal equations

$$\mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T \mathbf{y}$$

which lead to

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

do not have a unique solution for β . Instead they provide infinitely many solutions which

- all fit the training data perfectly
- usually are not able to make reliable predictions.

In other words: $(\mathbf{X}^T \mathbf{X})$ is not invertible.

Quadratic Penalization

Ridge Regression (Hoerl *et al.*, 1970)

A classic attempt for the $p > n$ situation is to add a quadratic penalty to the OLS-criterion.

$$\min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \beta^T \beta$$

which leads to the solution

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Benefits

$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})$ is an invertible matrix. We get a unique solution for β .

6 / 31

7 / 31

8 / 31

Matrix Inversion

$(X^T X + \lambda \mathbb{I})$ is a $p \times p$ matrix. Thus the computational cost of inverting this matrix is $O(p^3)$ operations.

p is large

Typical expression arrays include between 1,000 and 20,000 genes. For 10,000 genes it means we have to invert a $10,000 \times 10,000$ -matrix. This involves $O(10,000^3)$ operations.

How long does it take?

Even on a modern PC this is a matter of hours or even days.

1 Motivation

2 Singular value decomposition

3 Examples

4 Summary

9 / 31

Getting 'tricky'

Singular value decomposition (SVD)

Every Matrix $X \in \mathbb{R}^{p \times n}$ can be decomposed in the following way

$$X = UDV^T = RV^T$$

$$\underbrace{\begin{bmatrix} X \end{bmatrix}}_{n \times p} = \underbrace{\begin{bmatrix} U \end{bmatrix}}_{n \times n} \underbrace{\begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_n \end{bmatrix}}_{n \times n} \underbrace{\begin{bmatrix} V \end{bmatrix}}_{p \times n}^T$$

11 / 31

Applying SVD to ridge regression

Another solution for ridge regression

Plugging $X = RV^T$ into

$$\hat{\beta} = \underbrace{(X^T X + \lambda \mathbb{I})}_{p \times p}^{-1} X^T y$$

one gets

$$\hat{\beta} = V \underbrace{(R^T R + \lambda \mathbb{I})}_{n \times n}^{-1} R^T y$$

Conclusion

We can get the parameters for ridge regression by

- applying SVD to X
- using R as a 'data-matrix'
- multiplying β by V

12 / 31

Applying SVD to other models I

Linear predictors

Many models 'connect' the variables with the outcome through a **linear predictor**. E.g.

- logistic regression
- linear discriminant analysis
- support-vector machines

Loss functions

The parameters are estimated by minimizing a **loss function**

$\sum_{i=1}^n L(y_i, \beta_0 + x_i^T \beta)$ which can be the

- squared error
- negative log-likelihood
- etc...

14 / 31

Getting 'tricky'

Singular value decomposition (SVD)

Every Matrix $X \in \mathbb{R}^{p \times n}$ can be decomposed in the following way

$$X = UDV^T = RV^T$$

$$\underbrace{\begin{bmatrix} X \end{bmatrix}}_{n \times p} = \underbrace{\begin{bmatrix} U \end{bmatrix}}_{n \times n} \underbrace{\begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_n \end{bmatrix}}_{n \times n}^R \underbrace{\begin{bmatrix} V \end{bmatrix}}_{p \times n}^T$$

11 / 31

About Singular Value Decomposition

Computational costs

Getting the parameters for ridges regression via the SVD trick involves

- reducing p variables to n variables by SVD $\rightarrow O(pn^2)$
- solving n dimensional ridge regression $\rightarrow O(n^3)$
- transforming back to p dimensions $\rightarrow O(np)$

Thus it saves inverting a $p \times p$ matrix with $O(p^3)$ operations.

Software

SVD is a standard linear algebra tool. It is implemented in most mathematical languages. In R

- `svd()`
- `fast.svd()` \rightarrow faster for small n , large p
(`library(corpcor)`)

13 / 31

Applying SVD to other models II

Quadratic regularization

Like linear regression, all models using a linear predictor and some kind of loss function don't work properly when $p \gg n$.

This can be fixed by adding a **quadratic penalty**

$$\min_{\beta_0, \beta} \sum_{i=1}^n L(y_i, \beta_0 + x_i^T \beta) + \lambda \beta^T \beta.$$

The good news is

- The SVD trick can be used with all models in the above mentioned form.
- All aspects of model evaluation can be performed in this reduced space.

15 / 31

Theorem 1

Let $X = RV^T$ be the SVD of X , and denote by r_i the i th row of R , a vector of n predictor values for the i th observation. Consider the pair of optimization problems:

$$(\hat{\beta}_0, \hat{\beta}) = \operatorname{argmin}_{\beta_0, \beta \in \mathbb{R}^p} \sum_{i=1}^n L(y_i, \beta_0 + x_i^T \beta) + \lambda \beta^T \beta$$

$$(\hat{\theta}_0, \hat{\theta}) = \operatorname{argmin}_{\theta_0, \theta \in \mathbb{R}^n} \sum_{i=1}^n L(y_i, \theta_0 + r_i^T \theta) + \lambda \theta^T \theta$$

Then $\hat{\beta}_0 = \hat{\theta}_0$, and $\hat{\beta} = V\hat{\theta}$.

Proof

see Hastie *et al.*, 2004

Where to get λ from?

The regularization parameter λ is often selected by k -fold cross-validation.

How does cross-validation work?

- divide the training data into k groups of size $\frac{n}{k}$
- fit the model to $\frac{k-1}{k}$ and test it on $\frac{1}{k}$ of the data
- repeat this k separate times
- average the results

This is done for a series of values for λ , and a preferred value is chosen.

16 / 31

17 / 31

Cross-validation II

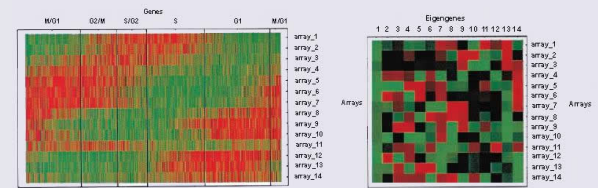
Illustration / Eigengenes

Corollary 1

The entire model-selection process via cross-validation can be performed using a single reduced data set R . Hence, when we perform cross-validation, we simply sample from the rows of R .

Proof

Cross-validation relies on predictions $x^T \beta$, which are equivariant under orthogonal rotations. \square

 $X \rightarrow R$ 

Eigengenes

If the columns of X are centered, the columns or R

- are the principal components of X
- also called **eigengenes**.

18 / 31

19 / 31

Derivatives I

Derivatives II

1st/2nd derivate

If the loss function is based on a log-likelihood we can:

- perform score tests with its first derivative
- gain variances of the parameters from its second derivative (information matrix / hessian matrix)

Expensive computation

Differentiation in the p -dimensional space is a waste of time.

SVD once again

SVD helps us to obtain the p -dimensional derivatives from the n -dimensional versions.

Corollary 2

Define

$$L(\beta) = \sum_{i=1}^n L(y_i, \beta_0 + x_i^T \beta), \quad L(\theta) = \sum_{i=1}^n L(y_i, \beta_0 + r_i^T \theta).$$

Then with $\theta = V^T \beta \Rightarrow L(\beta) = L(\theta)$.

If L is differentiable, then

$$\frac{\partial L(\beta)}{\partial \beta} = V \frac{\partial L(\theta)}{\partial \theta};$$

$$\frac{\partial^2 L(\beta)}{\partial \beta \partial \beta^T} = V \frac{\partial^2 L(\theta)}{\partial \theta \partial \theta^T} V^T,$$

with the partial derivatives in the right-hand side evaluated at $\theta = V^T \beta$

20 / 31

21 / 31

Derivatives III

Parenthesis: Woodbury matrix identity

Proof of Corollary 2

The equivalence

$$L(\beta) = L(\theta)$$

follows immediately from the identity

$$X = RV^T,$$

and the fact that x_i^T and r_i^T are the i th rows of X and R . The derivatives are simple applications of the chain rule to $L(\beta) = L(\theta)$. \square

Matrix inversion lemma

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

where

- A is $p \times p$
- U is $p \times n$
- C is $n \times n$
- V is $n \times p$

Link to SVD

For

- $A = \lambda I$
- $UCV = X^T X = (VDU^T)(UDV^T) = VD^2V^T$

we see the coherency with SVD.

22 / 31

23 / 31

- 1 Motivation
- 2 Singular value decomposition
- 3 Examples
- 4 Summary

Classes of models

SVD is applicable to all models in the mentioned form, of which there are many.

Examples

- Generalized linear models
- The Cox proportional hazards model
- Multiple logistic regression
- Regularized linear discriminant analysis
- Neural networks

24 / 31

25 / 31

Example 1

Ramaswamy *et al.*, 2001

Multiple logistic regression in order to predict tumor class:

- $n = 144$ training tumor samples
- 14 tumor classes
- $p = 16063$ genes for each sample

Amount of time for parameter estimation

- without SVD \rightarrow 8 days
- using SVD \rightarrow 0.4 seconds

Example 2

My own decent example

Simulation of a

- gene expression data set X : $p = 2000$ genes, $n = 3$ arrays (`rnorm()`).
- continuous output value y for each array (again `rnorm()`)

Fitting a ridge regression model

- conventional \rightarrow approx. 1 minute
- with SVD \rightarrow less than a second

see 'svd.R' on course homepage

26 / 31

27 / 31

Table of contents

Summary I

- 1 Motivation
- 2 Singular value decomposition
- 3 Examples
- 4 Summary

Why?

Estimating parameters in models for microarray data often involves inverting large ($p \times p$) matrices. This comes along with high computational costs.

How?

$$X = RV^T \quad (\text{SVD})$$

R can be used as a new data matrix with dimension $n \times n$. The p -dimensional parameter is then calculated in the following way

$$\hat{\beta} = V\hat{\theta} \quad \beta \in \mathbb{R}^p, \theta \in \mathbb{R}^n$$

28 / 31

29 / 31

Summary II

Literature

When?

SVD works with all models involving

- a linear predictor
- some loss function
- a quadratic penalty

And...

Important aspects of model evaluation, such as

- cross-validation
- testing parameters

can be performed with the reduced data set R .

Does it afford?

The time saving can be enormous!

- Hastie, T. and Tibshirani, R. 2004. Efficient quadratic regularization for expression arrays. *Biostatistics* **5**:329-340.
- Alter, O., Brown, P. and Botstein, D. 2000. Singular value decomposition for genome-wide expression data processing and modelling. *Proceedings of the National Academy of Sciences, USA* **97**, 10101-10106
- Ramaswamy, S., Tamayo, P., Rifkin, R., Mukherjee, S., Yeang, C., Angelo, M., Ladd, C., Reich, M., Latulippe, E., Mesirov, J. *et al.*, (2001). Multiclass cancer diagnosis using tumor gene expression signature. *Proceedings of the National Academy of Sciences, USA* **98**, 15149-15154
- Hoerl, A. E., Kennard, R. W. 1970. Ridge Regression: Biased Estimation for Non-Orthogonal Problems, *Technometrics*, **12**, 55-67.