

ON PRINCIPAL COMPONENTS AND REGRESSION: A STATISTICAL EXPLANATION OF A NATURAL PHENOMENON

Andreas Artemiou and Bing Li

Pennsylvania State University

Abstract: In this note we give a probabilistic explanation of a phenomenon that is frequently observed but whose reason is not well understood. That is, in a regression setting, the response (Y) is often highly correlated with the leading principal components of the predictor (\mathbf{X}) even though there seems no logical reason for this connection. This phenomenon has long been noticed and discussed in the literature, and has received renewed interest recently because of the need for regressing Y on \mathbf{X} of very high dimension, often with comparatively few sampling units, in which case it seems natural to regress on the first few principal components of \mathbf{X} . This work stems from a discussion of a recent paper by Cook (2007) which, along with other developments, described a historical debate surrounding, and current interest in, this phenomenon.

Key words and phrases: Dimension reduction, orientationally uniform distribution, principal components, random covariance matrices, regression, stochastic ordering.

1. Introduction

Cook (2007) described an intriguing historical debate surrounding the relation between the regression of a scalar response Y on a random vector \mathbf{X} and the principal components of \mathbf{X} . The debate arose from the practice of regressing Y on the first few principal components of \mathbf{X} , as suggested and advocated in Kendall (1957, p.75), Hocking (1976), Mosteller and Tukey (1977, p.307), Scott (1992). Some authors, however, question the logic behind this practice, on the basis that in the computation of principal components of \mathbf{X} , the response Y is never used in any direct or indirect way. For example, Cox (1968, p.272) writes:

A difficulty seems to be that there is no logical reason why the dependent variable should not be closely tied to the least important principal component.

This view is shared by authors such as Hotelling (1957) and Hawkins and Fatti (1984). See Cook (2007) for a more detailed description.

This question has received renewed interest recently due to the need for handling regression problems with very high-dimensional predictors but relatively few observation units, as one encounters when analyzing microarray data, so that the sample covariance matrix of \mathbf{X} is singular and the usual regression techniques cannot be directly applied. Under these circumstances, regressing Y on the first few principal components is a practical solution and often gives reasonable results. See Alter, Brown and Botstein (2000), Chiaromonte and Martinelli (2002), Bura and Pfeiffer (2003), and Li and Li (2004). See also Cook, Li and Chiaromonte (2007) for a different approach to this problem.

In his comments on Cook (2007), Li (2007) made a conjecture in an attempt to explain probabilistically why the response should be related to the leading principal components of the predictors, which was stated roughly as follows:

If nature arbitrarily selects a covariance matrix Σ for \mathbf{X} and coefficients β for the regression of Y on \mathbf{X} , then the principal components of \mathbf{X} of higher ranks tend to have stronger correlations with Y than do those of lower ranks.

Li (2007) argued intuitively that if \mathbf{X} is concentrated on a single direction, then the only way for Y to be correlated with \mathbf{X} at all is to be correlated with its first principal component. Likewise if \mathbf{X} has an elongated distribution the \mathbf{X} components in the longer axes should on average bear stronger correlations with Y . Now if Σ is selected arbitrarily then \mathbf{X} would have a large probability of having an elongated distribution, and would therefore effect the similar probabilistic ordering of correlations, even if the relation between Y and \mathbf{X} is independent of the shape of the distribution of \mathbf{X} . He supported this conjecture with several simulation studies, that affirmed it.

In this paper we give a precise formulation and a rigorous proof of the conjecture. This provides at least a partial justification for regressing Y on the principal components of \mathbf{X} , and gives fresh insights into a debate of historical interest and of importance in contemporary data analysis.

In Section 2, we will demonstrate that the conjectured ordering of correlations does occur naturally in practice by analyzing 33 data sets chosen arbitrarily from a database. The conjecture is then formulated and proved in Section 3.

2. The Phenomenon as Seen Through 33 Data Sets

In this section we analyze a collection of data sets in *Arc* software database, which can be found at <http://www.stat.umn.edu/arc/software.html>. From this collection we select 33 suitable data sets. The excluded data sets either have

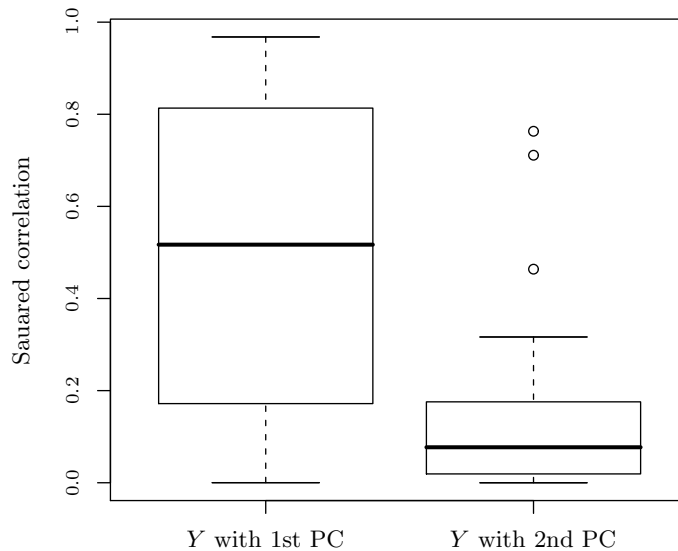


Figure 1. Boxplots of the squared correlations between Y and the first principal component of \mathbf{X} (left) and those between Y and the second principal component of \mathbf{X} (right) for 33 data sets.

only one predictor, or have categorical predictors or responses, or are artificially constructed.

The selected data sets contain from 2 to 12 predictors. For each of them we calculated all principal components of the predictor and their squared correlation with the response. Among these 33, in 24 cases the first principal components have the highest correlation with the responses, in 6 cases the second principal components have the highest squared correlation, in 2 cases the third principal components have the highest squared correlation, and in 1 case the fifth principal component has the highest square correlation. In Figure 1 we present the boxplots of the squared correlation coefficients between the responses and the first principal components of the predictors (left) and between the responses and the second principal components of the predictors (right) for the 33 data sets. This does indicate the tendency for the response to have higher squared correlation with the first principal component of the predictor.

3. Formulation and Proof of the Conjecture

Recall that p random elements, say W_1, \dots, W_p , are exchangeable if, for any permutation (i_1, \dots, i_p) of $(1, \dots, p)$, we have $(W_{i_1}, \dots, W_{i_p}) \stackrel{\mathcal{D}}{=} (W_1, \dots, W_p)$ where $\stackrel{\mathcal{D}}{=}$ indicates two random elements having the same distribution. To give the conjecture a rigorous formulation we need a precise definition of a random

covariance matrix that has equal probability of any orientation.

Definition 3.1. We say that a $p \times p$ positive semidefinite random matrix Σ has an orientationally uniform distribution if $\Sigma = \sigma_1^2 \mathbf{v}_1 \mathbf{v}_1^T + \dots + \sigma_p^2 \mathbf{v}_p \mathbf{v}_p^T$, where each $(\sigma_i^2, \mathbf{v}_i)$ is a pair of random elements in which σ_i^2 is a positive random variable and \mathbf{v}_i is a p -dimensional random vector, such that $(\sigma_1^2, \dots, \sigma_p^2)$ are exchangeable with distribution dominated by Lebesgue measure, $(\mathbf{v}_1, \dots, \mathbf{v}_p)$ are exchangeable and orthonormal, $(\sigma_1^2, \dots, \sigma_p^2)$ and $(\mathbf{v}_1, \dots, \mathbf{v}_p)$ are independent.

Thus if Σ is orientationally uniform, the random ellipsoid $\{\mathbf{x} : \mathbf{x}^T \Sigma \mathbf{x} \leq c\}$ has equal probability to have any orientation. Or, from a different perspective, suppose that \mathbf{X} is a p -dimensional random vector satisfying $E(\mathbf{X}|\Sigma) = \mathbf{0}$ and $\text{Var}(\mathbf{X}|\Sigma) = \Sigma$, then any random variable among $\mathbf{v}_1^T \mathbf{X}, \dots, \mathbf{v}_p^T \mathbf{X}$ is equally likely to be the 1st, 2nd, ..., or p th principal component of \mathbf{X} .

In the following, the symbol $\perp\!\!\!\perp$ indicates independence between random elements. We adopt the convention that if U is a random variable, then any constant α that satisfies $P(U < \alpha) \leq 1/2 \leq P(U \leq \alpha)$ is a median of U .

Lemma 3.1. Suppose β and $\mathbf{v}_1, \mathbf{v}_2$ are p -dimensional random vectors such that 1. $\beta \perp\!\!\!\perp (\mathbf{v}_1, \mathbf{v}_2)$; 2. $P(\beta \in G) > 0$ for any nonempty open set G ; 3. \mathbf{v}_1 and \mathbf{v}_2 are linearly independent and exchangeable. Then $(\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2$ has a unique median of 1.

Proof. First, we show that 1 is a median of $(\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2$; that is,

$$P\left(\frac{(\beta^T \mathbf{v}_2)^2}{(\beta^T \mathbf{v}_1)^2} < 1\right) \leq \frac{1}{2} \leq P\left(\frac{(\beta^T \mathbf{v}_2)^2}{(\beta^T \mathbf{v}_1)^2} \leq 1\right). \tag{3.1}$$

Because $(\mathbf{v}_1, \mathbf{v}_2)$ are exchangeable and $\beta \perp\!\!\!\perp (\mathbf{v}_1, \mathbf{v}_2)$, the random variables $(\beta^T \mathbf{v}_1)^2$ and $(\beta^T \mathbf{v}_2)^2$ are exchangeable. Hence $P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 \leq 1) = P((\beta^T \mathbf{v}_1)^2 / (\beta^T \mathbf{v}_2)^2 \leq 1) = 1 - P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 < 1)$. It follows that $P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 < 1) \leq 1 - P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 < 1)$ and $P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 \leq 1) \geq 1 - P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 \leq 1)$, which implies (3.1).

Now we show that 1 is the only number that satisfies (3.1). In other words, for any $0 < c_1 < 1$ and $c_2 > 1$ we have $P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 \leq c_1) < 1/2$ and $P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 < c_2) > 1/2$. We only show the first inequality; the second can be shown similarly. Since $P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 \leq c_1) = E[P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 \leq c_1 | \mathbf{v}_1, \mathbf{v}_2)]$, it suffices to show that for any nonrandom, linearly independent $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$, we have $P((\beta^T \mathbf{v}_2)^2 / (\beta^T \mathbf{v}_1)^2 \leq c_1 | (\mathbf{v}_1, \mathbf{v}_2) = (\mathbf{a}, \mathbf{b})) < 1/2$. However, because $(\mathbf{v}_1, \mathbf{v}_2) \perp\!\!\!\perp \beta$, the above inequality is equivalent to

$$P\left(\frac{(\beta^T \mathbf{b})^2}{(\beta^T \mathbf{a})^2} \leq c_1\right) < \frac{1}{2}. \tag{3.2}$$

Let $c_3 \in (c_1, 1)$. Since (\mathbf{a}, \mathbf{b}) has full column rank, the system of equations

$$\begin{cases} \boldsymbol{\beta}^T \mathbf{b} = \sqrt{c_3} \\ \boldsymbol{\beta}^T \mathbf{a} = 1 \end{cases}$$

has a solution, say $\boldsymbol{\beta}_0$. Note that $(\boldsymbol{\beta}_0^T \mathbf{b})^2 / (\boldsymbol{\beta}_0^T \mathbf{a})^2 = c_3 \in (c_1, 1)$. Because $\boldsymbol{\beta} \mapsto (\boldsymbol{\beta}^T \mathbf{b})^2 / (\boldsymbol{\beta}^T \mathbf{a})^2$ is continuous there is a neighborhood of $\boldsymbol{\beta}_0$, say G , such that $\boldsymbol{\beta} \in G \Rightarrow (\boldsymbol{\beta}^T \mathbf{b})^2 / (\boldsymbol{\beta}^T \mathbf{a})^2 \in (c_1, 1)$. By assumption, $P(\boldsymbol{\beta} \in G) > 0$. Therefore $P((\boldsymbol{\beta}^T \mathbf{b})^2 / (\boldsymbol{\beta}^T \mathbf{a})^2 \in (c_1, 1)) > 0$ which, combined with (3.1), implies (3.2).

We are now ready to establish the main result of the paper.

Theorem 3.1. *Suppose*

1. $\boldsymbol{\Sigma}$ is a $p \times p$ orientationally uniform random matrix,
2. \mathbf{X} is a p -dimensional random vector with $E(\mathbf{X}|\boldsymbol{\Sigma}) = 0$ and $\text{Var}(\mathbf{X}|\boldsymbol{\Sigma}) = \boldsymbol{\Sigma}$,
3. $Y = \boldsymbol{\beta}^T \mathbf{X} + \delta$, where $\boldsymbol{\beta}$ is a p -dimensional random vector and δ is a random variable such that $\boldsymbol{\beta} \perp (\mathbf{X}, \boldsymbol{\Sigma})$, $\delta \perp (\mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\Sigma})$, $E(\delta) = 0$ and $\text{Var}(\delta) < \infty$.
4. $P(\boldsymbol{\beta} \in G) > 0$ for any nonempty open set $G \in \mathbb{R}^p$.

Let w_1, \dots, w_p be the 1st, ..., p th principal components of \mathbf{X} , and let $\rho_i = \rho_i(\boldsymbol{\beta}, \boldsymbol{\Sigma}) = \text{Corr}^2(Y, w_i | \boldsymbol{\beta}, \boldsymbol{\Sigma})$. Then, whenever $i < j$, $P(\rho_i \geq \rho_j) > 1/2$.

Proof. Let τ^2 denote $\text{Var}(\delta)$. Let $(\sigma_{(1)}^2, \mathbf{v}_{(1)}), \dots, (\sigma_{(p)}^2, \mathbf{v}_{(p)})$ be the reordered $(\sigma_1^2, \mathbf{v}_1), \dots, (\sigma_p^2, \mathbf{v}_p)$ such that $\sigma_{(1)}^2 \geq \dots \geq \sigma_{(p)}^2$. We derive an explicit expression for ρ_i . Note that

$$\begin{aligned} \text{Cov}(Y, \mathbf{v}_{(i)}^T \mathbf{X} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) &= \text{Cov}(\boldsymbol{\beta}^T \mathbf{X} + \delta, \mathbf{v}_{(i)}^T \mathbf{X} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) \\ &= \boldsymbol{\beta}^T \boldsymbol{\Sigma} \mathbf{v}_{(i)} + \text{Cov}(\delta, \mathbf{v}_{(i)}^T \mathbf{X} | \boldsymbol{\beta}, \boldsymbol{\Sigma}). \end{aligned} \tag{3.3}$$

Because $\delta \perp (\boldsymbol{\Sigma}, \mathbf{X}, \boldsymbol{\beta})$, we have $\delta \perp (\mathbf{v}_{(i)}^T \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\Sigma})$. This implies $\delta \perp \mathbf{v}_{(i)}^T \mathbf{X} | (\boldsymbol{\beta}, \boldsymbol{\Sigma})$, and hence that the second term in (3.3) is zero. Because $(\sigma_{(i)}^2, \mathbf{v}_{(i)})$ is an eigen pair of $\boldsymbol{\Sigma}$, we have $\boldsymbol{\Sigma} \mathbf{v}_{(i)} = \sigma_{(i)}^2 \mathbf{v}_{(i)}$. Hence

$$\text{Cov}^2(Y, \mathbf{v}_{(i)}^T \mathbf{X} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \sigma_{(i)}^4 (\boldsymbol{\beta}^T \mathbf{v}_{(i)})^2. \tag{3.4}$$

In the meantime, $\text{Var}(Y | \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \text{Var}(\boldsymbol{\beta}^T \mathbf{X} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) + 2 \text{Cov}(\boldsymbol{\beta}^T \mathbf{X}, \delta | \boldsymbol{\beta}, \boldsymbol{\Sigma}) + \text{Var}(\delta | \boldsymbol{\beta}, \boldsymbol{\Sigma})$. Because $\delta \perp (\boldsymbol{\beta}, \boldsymbol{\Sigma})$, the last term on the right is simply τ^2 . Because $\delta \perp (\boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{X})$, we have $\delta \perp \boldsymbol{\beta}^T \mathbf{X} | (\boldsymbol{\beta}, \boldsymbol{\Sigma})$. So the second term on the right is 0. Hence

$$\text{Var}(Y | \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \boldsymbol{\beta}^T \boldsymbol{\Sigma} \boldsymbol{\beta} + \tau^2. \tag{3.5}$$

Moreover,

$$\text{Var}(\mathbf{v}_{(i)}^T \mathbf{X} | \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \mathbf{v}_{(i)}^T \boldsymbol{\Sigma} \mathbf{v}_{(i)} = \sigma_{(i)}^2. \tag{3.6}$$

Now combine (3.4), (3.5), and (3.6) to obtain

$$\rho_i = \text{Corr}(Y, \mathbf{v}_{(i)}^T \mathbf{X}) = \frac{\sigma_{(i)}^2 (\boldsymbol{\beta}^T \mathbf{v}_{(i)})^2}{\boldsymbol{\beta}^T \boldsymbol{\Sigma} \boldsymbol{\beta} + \tau^2}. \tag{3.7}$$

Cook (2007), in his rejoinder to Li (2007), gave a special case of (3.7).

Let $i < j$. Then, using (3.7), we deduce

$$P(\rho_i \geq \rho_j) = P\left(\frac{\sigma_{(i)}^2 (\boldsymbol{\beta}^T \mathbf{v}_{(i)})^2}{\boldsymbol{\beta}^T \boldsymbol{\Sigma} \boldsymbol{\beta} + \tau^2} \geq \frac{\sigma_{(j)}^2 (\boldsymbol{\beta}^T \mathbf{v}_{(j)})^2}{\boldsymbol{\beta}^T \boldsymbol{\Sigma} \boldsymbol{\beta} + \tau^2}\right) = P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_{(i)})^2}{(\boldsymbol{\beta}^T \mathbf{v}_{(j)})^2} \geq \frac{\sigma_{(j)}^2}{\sigma_{(i)}^2}\right).$$

The right hand side can be written as

$$\sum_{k \neq \ell} P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_{(i)})^2}{(\boldsymbol{\beta}^T \mathbf{v}_{(j)})^2} \geq \frac{\sigma_{(j)}^2}{\sigma_{(i)}^2} \middle| \sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2\right) P(\sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2).$$

Because $\sigma_1^2, \dots, \sigma_p^2$ are exchangeable, $(\sigma_{(i)}^2, \sigma_{(j)}^2)$ has equal probability to be $(\sigma_k^2, \sigma_\ell^2)$ for any $k \neq \ell$, and that probability is $\binom{p}{2}^{-1}$. Hence the above reduces to

$$\begin{aligned} & \binom{p}{2}^{-1} \sum_{k \neq \ell} P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_{(i)})^2}{(\boldsymbol{\beta}^T \mathbf{v}_{(j)})^2} \geq \frac{\sigma_{(j)}^2}{\sigma_{(i)}^2} \middle| \sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2\right) \\ &= \binom{p}{2}^{-1} \sum_{k \neq \ell} P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_k)^2}{(\boldsymbol{\beta}^T \mathbf{v}_\ell)^2} \geq \frac{\sigma_\ell^2}{\sigma_k^2} \middle| \sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2\right), \end{aligned} \tag{3.8}$$

where the equality follows from the fact that, conditioning on the event $(\sigma_{(i)}^2, \sigma_{(j)}^2) = (\sigma_k^2, \sigma_\ell^2)$, one has $(\mathbf{v}_{(i)}^2, \mathbf{v}_{(j)}^2) = (\mathbf{v}_k^2, \mathbf{v}_\ell^2)$.

Reexpress each term in the summation in (3.8) as

$$E\left[P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_k)^2}{(\boldsymbol{\beta}^T \mathbf{v}_\ell)^2} \geq \frac{\sigma_\ell^2}{\sigma_k^2} \middle| \sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2, \sigma_k^2, \sigma_\ell^2\right) \middle| \sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2\right]. \tag{3.9}$$

From Definition 3.1 we have

$$\begin{aligned} (\mathbf{v}_k, \mathbf{v}_\ell) \perp\!\!\!\perp (\sigma_1^2, \dots, \sigma_p^2) &\Rightarrow (\mathbf{v}_k, \mathbf{v}_\ell) \perp\!\!\!\perp (\sigma_1^2, \dots, \sigma_p^2; \sigma_{(1)}^2, \dots, \sigma_{(p)}^2) \\ &\Rightarrow (\mathbf{v}_k, \mathbf{v}_\ell) \perp\!\!\!\perp (\sigma_k^2, \sigma_\ell^2, \sigma_{(i)}^2, \sigma_{(j)}^2) \\ &\Rightarrow (\mathbf{v}_k, \mathbf{v}_\ell) \perp\!\!\!\perp (\sigma_{(i)}^2, \sigma_{(j)}^2) | (\sigma_k^2, \sigma_\ell^2). \end{aligned}$$

Thus the event $\{\sigma_k^2 = \sigma_{(i)}^2, \sigma_\ell^2 = \sigma_{(j)}^2\}$ can be removed from the conditional probability inside the conditional expectation (3.9), which then reduces to

$$E\left[P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_k)^2}{(\boldsymbol{\beta}^T \mathbf{v}_\ell)^2} \geq \frac{\sigma_\ell^2}{\sigma_k^2} \middle| \sigma_k^2, \sigma_\ell^2\right) \middle| \sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2\right]. \tag{3.10}$$

Because $(\boldsymbol{\beta}, \mathbf{v}_k, \mathbf{v}_\ell) \perp\!\!\!\perp (\sigma_k^2, \sigma_\ell^2)$, for each fixed $0 < s < t$,

$$P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_k)^2}{(\boldsymbol{\beta}^T \mathbf{v}_\ell)^2} \geq \frac{s}{t} \mid \sigma_k^2 = t, \sigma_\ell^2 = s\right) = P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_k)^2}{(\boldsymbol{\beta}^T \mathbf{v}_\ell)^2} \geq \frac{s}{t}\right) > \frac{1}{2},$$

where the inequality follows from Lemma 3.1. By Definition 3.1, the event $\{\sigma_k^2 = \sigma_\ell^2\}$ has probability 0. It follows that

$$P\left(\frac{(\boldsymbol{\beta}^T \mathbf{v}_k)^2}{(\boldsymbol{\beta}^T \mathbf{v}_\ell)^2} \geq \frac{\sigma_\ell^2}{\sigma_k^2} \mid \sigma_k^2, \sigma_\ell^2\right) > \frac{1}{2}$$

almost surely on the event $\{\sigma_{(i)}^2 = \sigma_k^2, \sigma_{(j)}^2 = \sigma_\ell^2\}$. Therefore (3.10), and hence (3.8), are strictly greater than 1/2.

The inequality $P(\rho_i \geq \rho_j) > 1/2$ in Theorem 3.1 is equivalent to

$$P(\rho_i \geq \rho_j) > P(\rho_i < \rho_j). \tag{3.11}$$

In other words ρ_i has a larger probability to be larger than ρ_j than to be smaller than ρ_j . We recall the definition of stochastic ordering.

Definition 3.2. Let U_1 and U_2 be two random variables whose distributions are dominated by a common measure μ . We say that U_1 is stochastically no greater than U_2 if, for each $c \in \mathbb{R}$, $P(U_1 \leq c) \geq P(U_2 \leq c)$. In this case we write $U_1 \stackrel{D}{\leq} U_2$. If in addition, $\mu(\{c : P(U_1 \leq c) > P(U_2 \leq c)\}) > 0$ then we say that U_1 stochastically (strictly) less than U_2 , and write $U_1 \stackrel{D}{<} U_2$.

This version of the definition of stochastic ordering is used, for example, in Li, Zha and Chiaromonte (2005). In general, inequality (3.11) is neither stronger nor weaker than stochastic ordering. However, in a special case it is weaker than stochastic ordering, as we show below.

For $i = 1, 2$, let F_i be the distribution of U_i and f_i be the density of U_i with respect to μ . We say that U_1 and U_2 have a common support if $\{f_1 > 0\} = \{f_2 > 0\}$. It is easy to see that if $U_1 \stackrel{D}{<} U_2$ and U_1 and U_2 have a common support, then

$$F_1(\{c : F_1(c) > F_2(c)\}) > 0, \quad F_2(\{c : F_1(c) > F_2(c)\}) > 0. \tag{3.12}$$

The following proposition gives a sufficient condition for $U_1 \stackrel{D}{<} U_2$ to imply $P(U_1 \leq U_2) > 1/2$.

Proposition 3.1. *Suppose U_1 and U_2 are random variables whose distributions are dominated by a common measure μ ; $U_1 \stackrel{D}{<} U_2$; $U_1 \perp\!\!\!\perp U_2$; and U_1 and U_2 have a common support. Then $P(U_1 \leq U_2) > 1/2$.*

Proof. By independence of U_1 and U_2 and by Fubini's Theorem,

$$\begin{aligned} P(U_1 \leq U_2) &= \int_{\mathbb{R}} \left[\int_{u_1 \leq u_2} f_1(u_1) \mu(du_1) \right] f_2(u_2) \mu(du_2) \\ &= \int_{\mathbb{R}} F_1(u_2) f_2(u_2) \mu(du_2) = \int_{\mathbb{R}} F_1(u_2) dF_2(u_2). \end{aligned}$$

By the second inequality in (3.12) the right hand side above is (strictly) greater than

$$\int_{\mathbb{R}} F_2(u_2) dF_2(u_2) = \left[\frac{F_2^2(u_2)}{2} \right]_{-\infty}^{\infty} = \frac{1}{2},$$

which completes the proof.

We must point out that the natural tendency described in this paper is neither definite nor particularly strong, and there is much room for improvement by sufficient dimension reduction, which reduces the dimension of \mathbf{X} in reference to Y . See Cook (2007).

Acknowledgements

We are grateful to a referee, an associate editor, and the Editor for their very helpful and prompt reviews. The research is supported in part by a National Science Foundation grant (DMS-0704621) to Bing Li.

References

- Alter, O., Brown, P. and Botstein, D. (2000). Singular value decomposition for gene-wide expression data processing and modelling. *Proc. Nat. Acad. Sci.* **97**, 10101-10106.
- Bura, E. and Pfeiffer, R. M. (2003). Graphical methods for class prediction using dimension reduction techniques on DNA microarray data. *Bioinformatics* **19**, 1252-1258.
- Chiaromonte, F. and Martinelli, J. (2002). Dimension reduction strategies for analyzing global gene expression data with a response. *Math. Biosci.* **176**, 123-144.
- Cook, R. D. (2007). Fisher lecture: Dimension reduction in regression. *Statist. Sci.* **22**, 1-40.
- Cook, R. D., Li, B. and Chiaromonte, F. (2007). Dimension reduction without matrix inversion. *Biometrika* **94**, 569-584.
- Cox, D. R. (1968). Notes on some aspects of regression analysis. *J. Roy. Statist. Soc. Ser. A* **131**, 265-279.
- Hawkins, D. M. and Fatti, L. P. (1984). Exploring multivariate data using the minor principal components. *The Statistician* **33**, 325-338.
- Hocking, R. R. (1976). The analysis and selection of variables in linear regression. *Biometrics* **32**, 1-49.
- Hotelling, H. (1957). The relationship of the newer multivariate statistical methods to factor analysis. *British J. Math. Statist. Psych.* **10**, 69-79.
- Kendall, M. G. (1957). *A course in Multivariate Analysis*. Griffin, London.

- Li, B. (2007). Comment: Fisher lecture: Dimension reduction in regression. *Statist. Sci.* **22**, 32-35.
- Li, B., Zha, H. and Chiaromonte, F. (2005). Contour regression: a general approach to dimension reduction. *Ann. Statist.* **33**, 1580-1616.
- Li, L. and Li, H. (2004). Dimension reduction methods for microarrays with application to censored survival data. *Bioinformatics* **20**, 3406-3412.
- Mosteller, F. and Tukey, J. W. (1977). *Data Analysis and Regression*. Addison-Wesley, Reading, Massachusetts.
- Scott, D. (1992). *Multivariate Density Estimation*. Wiley, New York.

Department of Statistics, Pennsylvania State University, 325 Thomas Building, University Park, PA 16802, U.S.A.

E-mail: aaa195@stat.psu.edu

Department of Statistics, Pennsylvania State University, 410 Thomas Building, University Park, PA 16802, U.S.A.

E-mail: bing@stat.psu.edu

(Received February 2008; accepted April 2008)