**INTERNATIONAL JOURNAL
FOR RESEARCH**

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

www.ijraset.com

Call: ○08813907089  |  E-mail ID: ijraset@gmail.com

# Computer Science and High Dimensional Data Modelling

Adarsh Tiwari[1], Pradeep Kanyal[2], Himanshu Panchal[3], Manjot Kaur Bhatia[4]

[1, 2, 3, 4]*Jagan Institute of Management Studies, Rohini Sector-5, New Delhi*

*Abstract: The need to grasp large database structures is a very important issue in biological and life science. This review paper is aimed toward quantitative medical researchers searching for guidance in modeling large numbers of variables in medical research, how this relates to straightforward linear models and therefore the geometry that underlies their analysis. Issues reviewed include LASSO-related approaches, principal-component based analysis, and problems with model stability and interpretation. Model misspecification issues associated with potential nonlinearities are examined, as is that the Bayesian perspective on these issues.*

## I. INTRODUCTION

As high-dimensional data structures have begun to be available and studied in many areas of medical research, the requirement for intuitive, geometric, and infrequently linear modelbased understanding of such data has grown. Genetic data, imaging data, health outcomes and clinical data, spatial positioning data, internet-based data: all are samples of settings where the flow of knowledge is huge and also the ability to investigate such a flow is usually restricted. a really sizable amount of variables and comparatively few subjects (large p and little n) are often the mark of such data and therefore the goal of the analysis is often to detect various patterns within the dataset. In genomic settings these is also simple mean differences in organic phenomenon levels across treatment groups, comprehensive correlated network clusters, or more detailed epigenetic patterns. Often there's limited theoretical modeling and far of the applied statistical research is empirically driven, falling under the hypothesis or model generating label, with the term "data science" sometimes used. Standard methods of statistical analysis often don't blockage well in such settings. Multiple comparison issues where an outsized number of case-control comparisons are conducted require careful application and interpretation.1 It is true that multiple testing of one-at-a-time mean differences may also be of limited utility for comprehending genomic and epigenetic data structures or gene networks pertinent to certain cell and phenotypic structures. A more complex three-dimensional nonlinear aspect of chromosome structure could also be relevant to such analyses.2 a geometrical perspective is helpful in understanding the properties of estimators and models developed during a linear model or associated analysis of variance (ANOVA) setting. These are often supported orthogonal projections onto a linear plane and therefore the space orthogonal thereto. The squared lengths of those projections will be compared, interpreted, and accustomed develop statistics for estimation and testing.3 When p. n. and standard projections are constrained by constrained dimensions, however, the basic geometric sense of linear models is disrupted. Here, rigorous use of conventional methods is required if the models that emerge are to be understandable. For linear models, correlation and inherent nonlinear interactions are equally problematic. the applying of linear models to correlated structures might not be appropriate, with nonlinear functional relationships potentially going undetected and creating instabilities within the predictive model. Many developmentally related growth factors are nonlinear in pattern.4 Scaling issues, gene clusters, and little embedded networks also affect the applicability of the linear model. Statisticians have developed techniques for restricted or sparse situations, including least-angle regression (LARS) extended via application of the smallest amount absolute shrinkage and selection operator (LASSO),5 and Dantzig6 approaches, which extend older techniques like restricted method of least squares, ridge regression, forward stagewise variable selection techniques, and principal components. Several earlier, more detailed reviews is found in Johnstone and Titterington,7 and in Bickel et al.8 Note that these approaches don't always converge to a fitted model (in which case an all-subsets search is required to search out a best fitting model, often impractical in terms of time)9 or provides a useful predictive model. Some methods have phase-threshold cutoff patterns that give insight into possible convergence.10 a part of application for high-dimensional methods is genetic data structures, which began with Southern blot electrogenesis technology and other fairly simple DNA-related technology and have developed into rather more detailed approaches including: single-nucleotide polymorphisms, copy-number variation, gene splicing, and RNA-related deep sequencing.11 These datasets often reflect specialized bioassays and there remains much to be done regarding standardizing platforms, alignment techniques, etc.1 Recently, the increase of epigenetics, the chemical triggers governing organic phenomenon (chromatin, histone, DNA methylation, for example), have cause one more level of complexity, as have the growing number of detected epigenetically triggered networks or clusters of genes

International Journal for Research in Applied Science & Engineering Technology (IJRASET)
*ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538*
*Volume 10 Issue XII Dec 2022- Available at www.ijraset.com*

that govern protein, cell, and other developmental and maintenance-related activities within the organism.12,13 Outside of genetics, the areas of systems and networkrelated biology, imaging, clinical data repositories, internetbased information, and other "large data" settings are all growing very quickly.14 Data science or big data-based approaches to identifying patterns in these large sets of collected data are quite varied, often reflecting a mixture of methods drawn from engineering, engineering science, and arithmetic.15 These statistically based methods are equally appropriate for certain research areas. From a practical and geometric standpoint, we will cover and analyse a number of related statistical techniques and tools of high-dimensional data analysis in this article. These techniques relate to and expand upon established standard techniques like ANOVA and main components analysis (PCA). the soundness of linear methods is examined and that we investigate the effect of misspecification, especially where this is often associated with nonlinearity. We review practical interpretations of p . n methods using the geometry of method of least squares, restricted statistical method, correlation, and simulated examples, briefly mentioning Bayesian perspectives during this setting.

## II. HIGH-DIMENSIONAL STATISTICAL APPROACHES AND LINEAR MODELS

A standard tool for understanding data and linking a response variable to numerous explanatory variables is that the linear model $y = X\beta + \varepsilon$ where y is additionally a $(n \times 1)$ vector of responses, $X = [x1 , …, xn]$ a $(n \times p)$ matrix of p measured variables, and $\varepsilon$ a $(n \times 1)$ vector of error components. If all variables xi are thought to be relevant, the fi ed method of method model is given by $\hat{y} = Xb$ where $b = (X'X)^{-1} X'y$. Typically within the settings reviewed here, many variables xi are collected and just some are thought to be relevant to predicting the very best result y. In such a setting, placing a particular restriction on the model, as an example, expecting some $\beta i$ values to differ significantly from zero, could even be helpful find the "best" underlying linear model, which is sometimes allotted using stepwise or stagewise methods. Correlations among the explanatory variables and resulting rank deficiency within the X matrix may additionally require the use of modified or restricted linear model-based approaches like ridge regression16 that have a protracted history. Sometimes a sparseness restriction is written as I | m I k t = 1 for arbitrarily tiny values of t and k. Sparseness limitations commonly take m = 1 or 2 into account. this is often most useful in

### A. Ridge Regression

In the case of highly correlated variables within the X design matrix, which affect the soundness and existence of $(X'X)^{-1}$, the older and more commonly used ridge-regression approach is applied and uses m = 2. it's worth examining ridge regression within the case n . p. Assuming centered data, the resulting estimator is given by $bR = (X'X + \lambda I)^{-1} X'y$ for scalar $\lambda$. Even with high correlation within the X design matrix this might exist, with $\lambda$ chosen graphically or via Bayesian posterior calculation. The singular value decomposition (SVD) guides the event. to use SVD normally we write $X = UDV'$, where $U = (u1 , …, up)$ could even be a n by p orthogonal matrix, the uj form an orthonormal basis for the column space of X, and V is additionally a similarly constructed orthogonal matrix for the row space of X. D is a square matrix $(d1 , …, dp)$. Geometrically, ridge regression is like projecting y onto the normalized principal components of X (ie, uj ) where the jth principal component of X is given by dj uj . Specifically, $bR j = [d2 j /d2 j + \lambda] uj ' y$ which could be viewed as weighting the projection of y onto the principal component uj by the relative weights of dj and $\lambda$. The important role played by eigenvalues within the appliance of linear models in restricted settings and dimension reduction generally is further discussed below. PCA LASSO procedures The sparseness restriction itself is applied directly as an additional restriction on the calculation of eigenvalues underlying multivariate techniques like cluster and statistical method,17 and these are observed as PCA LASSO

### B. Procedures

They are closely related to ridge-regression procedures.18 If the matrix A represents the transformation relating the initial data X to the principal components Y = AX then use of a sparseness restriction during this context gives the model Y Ax a t ij j p = < = : | ∑ | 1 where aij are the relevant coefficient elements of the Y = Ax principal components transformation, subject to the quality PCA constraints. As PCA methods are themselves often an initial try to understand or explore the underlying dimensionality or structure of the information, and thus the degree of sparseness itself in large data settings, this might overly restrict an initially explorative approach. In cases where p . n. is present, it does, however, help with convergence and interpretability of the resulting PCA. The enormous p behaviour reported in Hall et al.19 and Ahn et al.20 seems to suggest limited effectiveness of the PCA technique in scenarios without the sparseness limitation, therefore it may be interesting to investigate how the LASSO constraint interacts with the p. n geometry presented below.

## III. LARS

The LARS algorithm underlying mandatory LASSO approach to fitting models in standard n . p linear models is extremely stable21 and obtains a fitted model of size m in m steps. This approach be supported adding new variables during a forward stagewise search approach that uses equiangular bisectors to hunt out the foremost correlated variables within the dataset, adding them sequentially. As we are just projecting y onto a single or a limited number of xi vectors, y = xi (xi ′xi) 1 xi ′y or related residual vectors, the LARS technique is stable in and of itself because linear geometry still holds true. To handle a sparseness restriction in the p. n situation, the LARS method just needs to be significantly modified. Since this type of forward stagewise search and projection mostly avoids the multidimensional geometric elements addressed above, as demonstrated by Efron et al.21, the LASSO sparseness restriction has little impact on it. LARS is a straightforward algorithm. The xj that is best associated with y should be found starting with all bj = 0. Increase bj within the direction of the sign of its correlation with y and acquire the residuals (y - yˆ) stopping when another xk is found such the corr((y - yˆ), xk ) = corr((y - yˆ), xj ). Increase (bj , bk ) in their joint method direction until another predictor xm has the foremost amount correlation with the new residual vector. Continue during this fashion until all useful predictors are within the model. It are often shown that, if when a coefficient hits zero and is faraway from the active set of predictors the joint direction is recomputed, this procedure gives the whole path of LASSO solutions, as t is varied from zero to infinity. The LARS algorithm with LASSO restriction is out there as a package in R. Extensions of the LASSO approach are developed for logistic regression22 and survival analysis23 and lots of other settings. A Bayesian approach to those models could even be developed by assuming a modified Laplace prior distribution24 to elucidate the sparseness restriction directly on the parameter space. The Bayesian approach is interesting because it views the info as fixed and probabilities as directly attached to the parameter space. this means that style of the complexity of the p . n geometry when viewed from this attitude has relevancy only within the utmost amount because the likelihood function and prior elements are affected. That said, the employment of normality within the likelihood links Bayesian and frequentist approaches to method of method geometric considerations. this can be often often further discussed below.

### A. Example

Mouse genetic data a well known example drawn from the genomics literature is given in Ghazalpour et al.25 to convey how of the interaction of eigenvalue structure with p . n restrictions we apply PCA directly here, without LASSO restriction, viewing the results almost like selected p and n values for chromosome 11 where we start with 100 genes and their expression levels for 255 subjects. we start with n . p and randomly remove subjects as shown in Table 1, concluding

**Table I** Mouse data principal component analysis for values of n and p

| n | n* | p | e₁ | e₂ | e₃ | e₄ | e₅ | e₆ | e₇ | e₈ | e₉ | e₁₀ | e₁₁ | e₁₂ | Total variation |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 10 | 7 | 100 | 0.380 | 0.239 | 0.145 | 0.127* | 0.068 | 0.040 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| 20 | 15 | 100 | 0.267 | 0.202 | 0.145 | 0.105 | 0.081* | 0.048 | 0.040 | 0.028 | 0.024 | 0.022 | 0.015 | 0.010 | 0.984 |
| 30 | 20 | 100 | 0.253 | 0.175 | 0.128 | 0.094 | 0.069 | 0.063 | 0.038* | 0.036 | 0.027 | 0.022 | 0.019 | 0.017 | 0.942 |
| 40 | 28 | 100 | 0.251 | 0.156 | 0.115 | 0.085 | 0.068 | 0.057 | 0.037 | 0.034* | 0.030 | 0.025 | 0.019 | 0.018 | 0.896 |
| 50 | 37 | 100 | 0.544 | 0.088 | 0.074 | 0.046 | 0.038 | 0.033* | 0.024 | 0.020 | 0.017 | 0.014 | 0.013 | 0.011 | 0.921 |
| 100 | 80 | 100 | 0.365 | 0.107 | 0.079 | 0.064 | 0.045 | 0.043 | 0.036 | 0.032 | 0.030* | 0.020 | 0.017 | 0.014 | 0.853 |
| 150 | 121 | 100 | 0.353 | 0.100 | 0.070 | 0.067 | 0.042 | 0.039 | 0.037 | 0.033 | 0.027 | 0.022 | 0.017* | 0.016 | 0.824 |
| 200 | 157 | 100 | 0.360 | 0.098 | 0.077 | 0.060 | 0.040 | 0.036 | 0.034 | 0.030 | 0.025 | 0.022 | 0.017 | 0.016 | 0.815 |
| 254 | 200 | 100 | 0.338 | 0.101 | 0.084 | 0.056 | 0.039 | 0.036 | 0.033 | 0.030 | 0.030 | 0.025 | 0.022 | 0.018 | 0.800 |

**Notes:** Proportion of total variation shown; *denotes 80% of variation explained. The e are the ordered principal components, n is the sample size, and p the number of variables. Results for the first four principal components with p > n are highlighted in bold.

PCA for every set of subjects and variables, stepping into the p . n context. The 12 largest eigenvalues for every analysis are reported. Note that when p . n There are exactly n non-zero eigenvalues that can exist. For the initial PCA with n = 255 (n* with missing values) 12 PCA variables account for 90% of the total variation implying potential structure within the set of gene expressions. As is common in most PCA analysis, the primary eigenvalue is seen as an overall average. Of greater interest in relevance the results discussed here is that the overall structure of the remaining eigenvalues as p/n increases. As this increases, there are fewer relative sources of variation or degrees of freedom and also the subsequent level of total variation explained. As this is often often often real data, and p is large but not within the realm of enormous asymptotic values, the expected similarity of eigenvalues is seen as slowly occurring, especially beyond the primary or largest PCA, subject to random error. Further restricting our view here to the primary four eigenvalues, we see a fragile increase in their values and similarity as p/n increases when p . n, while the other numbers remain constant or decrease toward zero. As noted within the results of Hall et al19 and Ahn et al20 above, as p . n the knowledge provided by the eigenvalues is also a smaller amount useful with relevance identifying clusters within the data. The results seem to point a growing convergence to a smaller subgroup of eigenvalues.

## REFERENCES

[1]  GALE EBOOKS https://go.gale.com/ps/i.do?id=GALE%7CA411198218

[2]  High-dimensional data and linear models by Brimacombe M

[3]  https://www.sciencedirect.com/topics/computer-science/high-dimensional-data

[4]  https://www.dovepress.com/high-dimensional-data-and-linear-models-a-review-peer-reviewed-fulltext-article-OAMS

[5]  https://www.researchgate.net/publication/330380054_High-Dimensional_LASSO-Based_Computational_Regression_Models_Regularization_Shrinkage_and_Selection

[6]  https://www.softwaretestinghelp.com/dimensional-data-model-in-data-warehouse/

# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)