

A Collection of Numerical Recipes Useful For Building Scalable Psychometric Applications

Harold Doran

hdoran@humrro.org

Human Resources Research Organization

Working Technical Paper

December 5, 2021

Abstract

Next generation psychometrics will inevitably leverage advancements made available with modern computing and massive data to create psychometric innovations. Realizing this future requires that psychometricians adopt better approaches for modern computing known to support scalability. Unfortunately, well-known numerical recipes for large-scale computing have not been fully embraced in psychometrics and this limits innovation. Commonly, textbook representations of mathematical solutions are confused to be the same computational strategy for implementation or trivial approximations are often used when faced with complex scenarios. Psychometricians will benefit from focused conversations around a core skill set needed to build scalable applications including numerical analysis, computing, and software development. This work sets out to focus on only one of those components and a collection of numerical recipes is provided.

Keywords: computational psychometrics; gaussian quadrature; item response theory; linear models; numerical analysis; psychometric data scientist

Introduction

Psychometrics is evolving and increasingly becoming a multidisciplinary field with a more direct integration of modern computing with psychometric methods (ACTNext, 2019). This presents a future where next generation psychometricians will blend expertise in psychological measurement with their computational acumen to generate novel psychometric methods. However, a new challenge facing psychometricians is that massive data are now typical and new methods of analysis, such as machine learning (ML), are taking greater foothold in psychometrics. As a result, many

of the traditional methods used in psychometric computing need reconsideration and the changing dynamics of large data and computationally demanding models must be confronted in new ways.

The National Research Council published an important report on this topic directly addressing how science is changing (National Research Council, 2013). The NRC report made clear that in the world of massive data, statisticians should be concerned specifically with scalability and algorithms and adapt to the new types of challenging computational demands. One approach tacitly assumed as a solution for tackling scalability is to engage larger environments for psychometric computing. For instance, distributed computing in cloud-based infrastructures can involve shared computing power to solve big data problems. This is a helpful approach, but it is critically more important to first consider the smallest and most manageable unit—the computational methods used in the underlying code to build the psychometric application. Then, good computational methods can more readily take advantage of larger computing environments than expecting poor methods to be remedied by more powerful computing.

As a first priority, it is important to reconsider the methods of numerical analysis used to estimate parameters and build models. These are the core at which scientific ideas are implemented and computational psychometricians require a portfolio of efficient, scalable techniques at the ready so that new ideas can unfold. Unfortunately, sound numerical methods are not widely embraced within the field of psychometrics and, perhaps, this is rooted in how we learn. One general observation is that textbooks commonly represent psychometric solutions in mathematical ways that conflate computational implementation with the notational expression. Notational representations should be understood as valuable for developing a conceptual understanding of the psychometric problem and to communicate the idea of the model or its parameterization. However, those same representations should not be interpreted as the same way to build a computational implementation, although it appears they often are.

Bates (2004) noted that the algebraic textbook solution for the least squares model is commonly used within R code and provided details for better computational implementation. Genz and Kass (1997) report that statisticians often ignore superior methods for high dimensional integration, a problem also noted by Antal and Oranje (2007) in psychometric software applications finding that trivial rectangular rules are common even in the National Assessment of Educational Progress (Allen, Donoghue, & Schoeps, 2001) and commercially available psychometric software.

There are, however, many interesting conversations among psychometric innovators regarding the challenges in building scalable applications. Cai, Yang, and Hansen (2011) openly discuss the specific challenges between approaches for dealing with high dimensional integration problems as do the developers of the R package Dire (Doran, Bailey, Buehler, & joo Lee, 2021) for estimating the marginal likelihood regression (Mislevy, 1984). Rosseel (2021) provides an impressive example that translates a complex problem into a feasible computational solution. This work makes good use of a standard inversion lemma to reduce the overhead associated with inverting large matrices and takes advantage of certain special structures and reduces the problem to a yield a simpler solution and is an exemplar of psychometric computing.

Notably, these conversations are typically focused on a specific application or problem. In contrast, the work presented here is intentionally problem agnostic with the goal of organizing a general class of useful recipes to be used as templates for creating new solutions. For this reason, this work emphasizes application and is therefore an organization of tested numerical approaches with the hope that they become more familiar to psychometricians. This work does not set out to formally prove theorems or engage in deriving estimators and proving their properties other than

one simple proof in the appendix. Rather, it is an exploration of computational psychometrics considered through the lens of integrating psychological measurement with templates intending to springboard new psychometric models that are computationally demanding.

It’s important to acknowledge that numerical analysis techniques are situational and implemented as preference by the innovator suitable for a given condition. For instance, some of the methods here scale better than others or combinations of the approaches might be used. Also, this work makes use of a certain decomposition; however, others may have legitimate reasons for others. Hence, there is no motivation to oversell these specific concepts as the most desirable for any scenario. Instead, the underlying motivation is to encourage staring at a problem and finding a recipe through exploration and testing that yields a stable and reduceable structure to implement a complex procedure at scale. In other words, the idea proposed is that psychometricians need a toolkit allowing them to make a big problem small and the ideas presented here are generally useful, but not intended as the most ideal for every scenario.

The intent of the recipes provided here is to organize what is hoped to be a useful resource for those interested in computational psychometrics. The recipes cover a core set of numerically stable and scalable methods needed to support psychometric work. They are introduced in the framework of familiar and simple problems and the ideas are advanced sequentially to more complex scenarios. The work ends with three examples of methods not traditionally found in standard software simply to illustrate implementation of computational ideas and not to center attention on those specific methods.

Notation and Terminology

Standard notation is generally used throughout and the terms and notation are defined within each topic. For instance, matrices assume the general notation of uppercase bold, \mathbf{A} , transpose \mathbf{A}' , determinants, $|\cdot|$, vectors are generally lowercase bold, \mathbf{x} , elements within sets $\{x_1, x_2, \dots, x_n\}$, $\mathbb{E}(\cdot)$ is used for expected values, and functions of variables with a general form $f(\cdot)$. One convention to clarify is the distinction between $g(\theta; \beta)$ and $g(\theta|\beta)$. The former is used to mean the density of the random variable at the value of θ with the parameters β and the latter is used to mean the conditional distribution of the random variable given the parameters.

Distinctions in terminology are made throughout using “algebraic” representation and “computational” representation to mean very different things. Specifically, “algebraic solution” is used to mean the way a solution to a problem is written in order to convey a concept. Instead, a “computational solution” is a numerical recipe that provides a stable way to implement a solution in software to estimate parameters that may not resemble the algebraic solution. The distinction between the two is a non-trivial point.

Defining Scalable Psychometric Applications

The term “scale” is used in this work but requires clarification on its meaning and purpose in psychometrics. The term scale can be used to describe algorithms, networks, databases, and other complex systems. To remove ambiguity, scale in this work refers to the ability of an algorithm (i.e., the computational implementation) and the system making use of the algorithm (i.e., the software application) to handle big inputs and grow in relationship to the size of the problem

(Teng, 2016). The algorithm is how the idea is implemented and the application is how that idea becomes available to others.

In psychometric practice, there are three primary ways input sizes tend to grow. These include 1) the magnitude of the data used for analysis, 2) the complexity of the computational problem¹, and 3) the large number of simultaneous users accessing a common system. A scalable application is one that can flexibly handle all three components efficiently.

Traditional psychometric techniques, such as item calibrations in large populations, are commonly managed with sampling techniques as a data reduction method. This has been convenient for decades with simple, low-dimensional psychometric models that require relatively small computations with software installed on a local machine. In essence, this statement captures how scale is managed in traditional psychometrics—make big data small, use simple computational methods, and rely on single user applications installed on a local machine. However, the future of psychometrics will bring new models to assess learning, cognition, and human behavior in ways that will break the traditional psychometric paradigm. These new approaches will require the richness associated massive data, the benefits of complex psychometric models in high-dimensional space, and the modern nature of computing in a cloud-based infrastructure. Hence, the new psychometric paradigm requires a focused conversation on scale.

Introductory Concepts For Dealing With Correlated Data

Whitening Transformations for Correlated Variables

A concept first introduced and used throughout almost all subsequent recipes is referred to as the whitening transformation. Whitening is a process of transforming a collection of correlated vectors into a collection of uncorrelated vectors with the aid of a suitable decomposition (Kessy, Lewin, & Strimmer, 2018). The process is critical to the solution of linear systems and for fast quadrature routines and is generally helpful for reducing complex problems into more computationally feasible solutions.

The concept is introduced as follows. Let \mathbf{Y} be an $n \times p$ matrix of correlated variables from a p -dimensional multivariate normal, $\mathbf{Y} \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where $\boldsymbol{\Sigma}$ is symmetric and positive definite. Now find a decomposition that satisfies $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}'$ where \mathbf{L} is lower triangular known as the Cholesky factor, and $\boldsymbol{\Sigma}^{-1} = \mathbf{L}'^{-1}\mathbf{L}^{-1}$. This provides that the transformation $\mathbf{X} = \mathbf{Y}\mathbf{L}'^{-1}$ can be used so that the variables in the matrix \mathbf{X} are “whitened” or “de-correlated”. Notice this can also be used to create correlated vectors when the operation is reversed, $\mathbf{Y} = \mathbf{X}\mathbf{L}'$, which may be useful for data generation, simulation, or applications such as high dimensional integration.

There are many useful decompositions in numerical analysis (Searle, 1982; Zhang, 1999) and there are reasons to choose others depending on the computational objective. Cholesky is just one of many that is helpful for decorrelation and the appendix shows how the eigendecomposition can be used to support inverting a matrix. However, the Cholesky will be a focus in this work as it relates to the computational objectives described here and some work has demonstrated its potential advantage over others, such as the QR when used for solving normal equations (Lira, Iyer,

¹The computational complexity is considered as an input here because many complex problems require more inputs. For instance, high dimensional integrals require a large number of quadrature points.

Trindade, & Howle, 2016). Some applications have fully adopted a sparse Cholesky decomposition (Davis & Hager, 2005) for linear mixed models (Bates, 2004b) with tremendous success.

Multivariate Densities of Whitened Data

The multivariate normal density function has an important role in psychometric practice. It is possible to reduce this from a complex problem into separable parts and show how the multivariate density can be computed as a product of transformed univariate densities. This change of variable transformation has broader applications for psychometric computing in maximization methods or dealing with high-dimensional integration problems.

The general form of the p -dimensional multivariate normal for a non-diagonal covariance matrix Σ is

$$f_y(y_1, \dots, y_p) = \frac{\exp(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})' \Sigma^{-1}(\mathbf{y} - \boldsymbol{\mu}))}{\sqrt{(2\pi)^p |\Sigma|}} \quad (1)$$

If the mean vector is $\boldsymbol{\mu} = \mathbf{0}$ and $\Sigma = \text{diag}(1, \dots, 1)$, then the density of the multivariate normal is equal to the product of the marginal densities. Let φ denote the standard normal probability density function (p.d.f.), then

$$f_y(y_1, \dots, y_p) = \prod_{i=1}^p \varphi(y_i) \quad (2)$$

Note that (1) and (2) are related only under the special case of uncorrelated variables with unit variance. With correlated variables, rewrite (1) as

$$f_y(y_1, \dots, y_p) = \frac{\exp(-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})' \mathbf{L}'^{-1} \mathbf{L}^{-1}(\mathbf{y} - \boldsymbol{\mu}))}{\sqrt{(2\pi)^p |\mathbf{L} \mathbf{L}'|}} \quad (3)$$

$$= \frac{\exp(-\frac{1}{2} \mathbf{y}'^* \mathbf{y}^*)}{\sqrt{(2\pi)^p |\mathbf{L}^2|}} \quad (4)$$

$$= |\mathbf{L}|^{-1} \left(\prod_{i=1}^p \varphi(y_i^*) \right) \quad (5)$$

where $\mathbf{y}^* = \mathbf{L}^{-1} \mathbf{y}$ with elements $\mathbf{y}^* = \{y_1^*, \dots, y_p^*\}'$. This transform expresses (5) into the form of (2), thus showing how the multivariate density can be computed as a product of the (transformed) marginal densities after whitening.

Computational Concepts for Linear Systems

This section explores efficient computational recipes for linear models to create simple, reduceable structures to build scalable psychometric algorithms. It's best to begin with a simple, common problem—ordinary least squares (OLS), that will serve as a foundation and subsequently show how to frame more complex models with the same OLS solution. While the examples are provided within the context of linear regressions, the same concepts are used widely across psychometric

practice. For example, these same methods are often used in the expectation-maximization (EM) algorithm during the M-step, in natural language processing problems, latent semantic analysis, solutions to factor analysis problems, in highly parameterized item calibrations problems using the newton-raphson algorithm, and within iteratively reweighted least solutions for generalized linear models. These are just a subset of the possible applications.

Ordinary Least Squares as a Unifying Framework

The standard representation of the linear least squares problem is via the following normal equations (McCulloch & Searle, 2001)

$$\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{X}'\mathbf{y} \quad (6)$$

where \mathbf{X} is an $n \times p$ design matrix of full column rank, $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown parameters to be estimated, and \mathbf{y} is an $n \times 1$ vector of observed outcomes. Provisionally assume in the following examples $p \ll n$, however many machine learning applications are trained on data when $n < p$, a problem referred to as high dimensional regression. The algebraic representation in standard notation is shown having an invertible $\mathbf{X}'\mathbf{X}$ with a solution written as

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}. \quad (7)$$

It is common for texts to convey strong statements around the representation shown in Equation (7) with claims insinuating that this representation is also its computational implementation. In fact, this is a poor implementation and it is unnecessary to invert $\mathbf{X}'\mathbf{X}$. Instead a more stable approach is to find the decomposition, $\mathbf{X}'\mathbf{X} = \mathbf{L}\mathbf{L}'$, and then rewrite (6) as $\mathbf{L}\mathbf{L}'\boldsymbol{\beta} = \mathbf{X}'\mathbf{y}$. This expression provides a simple solution via two triangular systems requiring to first forwardsolve for \mathbf{z} and then backsolve for $\boldsymbol{\beta}$ as

$$\mathbf{L}\mathbf{z} = \mathbf{X}'\mathbf{y}, \quad \mathbf{L}'\boldsymbol{\beta} = \mathbf{z}. \quad (8)$$

Obtaining least squares estimates is very simple and (8) shows estimates for $\boldsymbol{\beta}$ are obtained without inverting any matrices along the way. However, it serves as a baseline and a general unifying framework for computing linear systems as subsequently shown when more complex models can be rewritten into the form of (6). In such cases, other models can be shown to inherit the solution shown in (8) even when they appear to be more complex, hence it unifies approaches to computing linear systems.

Generalized Least Squares in OLS Form

The generalized least squares (GLS) problem is an incrementally more advanced model than its OLS counterpart. The standard textbook normal equations for the GLS model is

$$\mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{X}\boldsymbol{\beta} = \mathbf{X}'\boldsymbol{\Sigma}^{-1}\mathbf{y} \quad (9)$$

where $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}'$ is a symmetric $(n \times n)$ positive definite covariance matrix and all others matrices are as previously noted. The algebraic solution for the GLS approach is commonly shown to involve matrix inversions as follows

$$\hat{\beta} = (X' \Sigma^{-1} X)^{-1} X' \Sigma^{-1} y. \quad (10)$$

The algebraic representation in (10) implies inverting the matrices Σ and $(X' \Sigma^{-1} X)$ which can be entirely avoided. The dimensions of Σ may very large as it is equal to the number of observations in X and the GLS problem often involves an iterative solution with elements of Σ changing over the iterations, so inversions would be costly if not impossible. The concept of whitening can again be used to rewrite the normal equations in (9) such that the data in X are “de-correlated” as

$$X^{*'} X^* \beta = X^{*'} y^* \quad (11)$$

where $X^* = L^{-1} X$ and $y^* = L^{-1} y$. This whitening transform writes (9) into the form (6) so that it can inherit the same solution show in (8).

If (11) were used in an iterative process and elements of Σ were updated at each step, then recomputing its Cholesky factor repeatedly would layer additional expense. This is where numerical recipes become situational and we encourage staring at Σ to determine if it can be separated and reduced into a simple structure. One such example is provided in the section on linear mixed models where a very large matrix is repeatedly updated in an iterative algorithm and finding its inverse at each iteration is costly. Instead, a standard inversion lemma is used and it is extended to exploit some special characteristics of that matrix allowing a complicated problem to be reduced into a very trivial computation implementation.

Weighted Least Squares in OLS Form

Weighted least squares (WLS) is a special case of GLS when the covariance matrix Σ is diagonal. In this case, it is entirely unnecessary to actually create and store the entire matrix Σ and instead only its diagonal elements, v_{jj} , are needed. In this case there is a special matrix property providing that $X' \Sigma^{-1} X = X' X^*$ where is arrived at X^* via the elementwise calculation of $1/v_{jj}$ over each column of the matrix X .

Performing this elementwise calculation to the matrix X on the left and right hand side of Equation 9 when Σ is diagonal yields $X' X^* \beta = X^{*'} y$, which is now written in the form of (6) and inherits the same computational solution.

Stochastic Gradient Descent for Linear Systems

The previously discussed methods show that linear models in the form of (6) inherit the general solution of (8). However, when the dimensions of the model matrix X become extremely large (generally n but also sometimes p), then computing by (8) may prove challenging. Stochastic gradient descent (SGD) is a powerful approach for building highly scalable applications that reduces the computational burden associated with large data by simplifying it to an iterative solution passing over the gradient one sample at a time (Cizek & Cizkova, 2004; Shamir & Zhang, 2013; Tran, Toulis, & Airolidi, 2015).

While SGD is a general purpose optimization technique useful for minimization of any differentiable objective function, $\mathcal{J}(\theta)$, it is presented within the section on linear models as it provides a fast and clear pathway for many common psychometric solutions, such as with machine learning

applications. It's conceivable that some training data sets used for ML could have hundreds of thousands or even millions of observations, in which case computing by (8) may be unrealistic.

Formally, let $\nabla \mathcal{J}_i(\boldsymbol{\theta})$ denote the gradient vector of the objective (loss) function with respect to the parameters to be minimized for an i th sample in the data. SGD performs the optimization by drawing samples from the observed data and repeatedly updating the gradient and the parameters at each iteration as

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha \nabla \mathcal{J}_i(\boldsymbol{\theta}) \quad (12)$$

where $\alpha > 0$ is a learning rate parameter, $\boldsymbol{\theta}_t$ is the provisional value of the parameters at iteration t , and $\nabla \mathcal{J}_i(\boldsymbol{\theta})$ is the value of the gradient computed using sample i .

With respect to models assuming the form of (6), the gradient is $\nabla \mathcal{J}(\boldsymbol{\theta}) = 2n^{-1}(\mathbf{X}'\mathbf{X}\boldsymbol{\theta} - \mathbf{X}'\mathbf{y})$. However, very large dimensions in \mathbf{X} make general optimization difficult with either general gradient descent or using the traditional closed form solution in (8). Instead, the iterative process outlined in Algorithm (1) allows for a large problem to be reduced simply.

Algorithm 1 Stochastic Gradient Descent

Input: Choose initial values for $\boldsymbol{\theta}$, $\nabla \mathcal{J}_i(\boldsymbol{\theta})$, and select a value for α .

- 1: Set $\boldsymbol{\theta}_{t+1} := \boldsymbol{\theta}_t - \alpha \nabla \mathcal{J}_i(\boldsymbol{\theta})$ at iteration t .
 - 2: Sample \mathbf{X}_i , in \mathbf{X} and its corresponding value, \mathbf{y}_i .
 - 3: Compute $\nabla \mathcal{J}_i(\boldsymbol{\theta})$ using $\{\mathbf{X}_i, \mathbf{y}_i\}$.
 - 4: Iterate between steps (1) and (3) until criterion for a stopping rule has been satisfied.
-

When $i = 1$, then \mathbf{X}_i is the i th row of the model matrix and this process is termed stochastic gradient descent. When $1 < i < n$, the same iterative process can be used in small batches termed mini-batch gradient descent and \mathbf{X}_i is a collection of rows in the model matrix. When $i = n$ with n being the total rows in \mathbf{X} , this is termed batch gradient descent.

SGD, in contrast to traditional least squares algorithms, works only with samples of data and not with the full model matrix \mathbf{X} at any point. This hints towards a different way to manage the amount of data needed in storage and the concept of streaming in only subsets of data needed for computing instead of requiring the complete data for memory management can be used with SGD (Tran et al., 2015). Consequently, SGD not only uses smaller computations, it also provides a means for using smaller subsets of data as it computes and the combination of the two is pivotal for computing with modern data that may contain many millions of observations.

A concept only noted in passing is that α is a learning rate and selecting values for this scalar are not fully explored here. Pragmatically speaking, choosing values for the learning rate parameter is a harder problem than it may seem at the surface and there is a good deal of research to further explore. A “large” value for learning is generally desired for initial evaluations of the gradient allowing the algorithm to take large steps towards the optimal value. However, using that same large value when approaching the optimum may cause for the algorithm to jump over the optimal value and not converge to the expected result. Hence, the learning rate “schedule” is a concept that adjusts the learning rate to smaller values as the algorithm proceeds. In addition, the example here assumes $p \ll n$, in which case the loss function associated with standard linear regression can be used. In fact, many machine learning applications involve scenarios where $n < p$, in which case

other approaches such as Ridge or Lasso regressions may prove useful, but can adopt the general SGD framework described in Algorithm (1).

As an aside, iterative algorithms need a stopping rule, often referred to as convergence criteria. In many respects, it's an arbitrary rule constructed for a given condition and there are multiple ways to achieve this outcome. For example, one might compute changes in the parameter estimates from iteration t to $t + 1$ and take the difference between them. Stopping might occur if no parameter differs by more than ϵ , a term set to be very small indicating the parameters are not changing enough from iteration to iteration any longer. Or, one might stop when differences in the log-likelihood differ by less than ϵ over iterations. In the SGD scenario, one might stop after reaching a fixed number of iterations or when the sum of the gradients squared is less than ϵ . No particular rule can be highly endorsed, although some scenarios are susceptible to being stuck in a saddle point.

Linear Mixed Models with Henderson's Method

The linear mixed model is generally written with the form $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$ (Laird & Ware, 1982) and is often presented as having the normal equations (McLean, Sanders, & Stroup, 1991)

$$\begin{bmatrix} \mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X} & \mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{Z} \\ \mathbf{Z}'\boldsymbol{\Omega}^{-1}\mathbf{X} & \mathbf{Z}'\boldsymbol{\Omega}^{-1}\mathbf{Z} + \mathbf{G}^{-1} \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta} \\ \mathbf{u} \end{bmatrix} = \begin{bmatrix} \mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{y} \\ \mathbf{Z}'\boldsymbol{\Omega}^{-1}\mathbf{y} \end{bmatrix} \quad (13)$$

where \mathbf{y} is an $n \times 1$ vector of outcomes, \mathbf{X} is an $n \times p$ model matrix for the fixed effects, $\boldsymbol{\beta}$ is a $p \times 1$ vector of estimates for the fixed effects, $\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_Q]$ is an $n \times Q$ model matrix for the random effects, $\mathbf{u}' = [\mathbf{u}'_1, \mathbf{u}'_2, \dots, \mathbf{u}'_Q]$ is the vector of random effects, and \mathbf{e} is the residual error term. The matrices $\boldsymbol{\Omega}$ and \mathbf{G} can have many forms depending on the structure of the model and are developed below for one variant of the linear mixed model.

Algorithm 2 Henderson Mixed Model Sketch

- Input:** Create starting values for σ_ϵ^2 and $\sigma_q^2 \forall q$.
- 1: Construct $\boldsymbol{\Omega} = \sigma_\epsilon^2 \mathbf{I}_n$ and $\mathbf{G} = \text{diag}(\{\sigma_1^2, \dots, \sigma_1^2\}, \{\sigma_2^2, \dots, \sigma_2^2\}, \dots, \{\sigma_Q^2, \dots, \sigma_Q^2\})$. The length of the diagonal element for block q is equal to the number of columns in its corresponding matrix, \mathbf{Z}_q .
 - 2: Solve the linear system for $\boldsymbol{\beta}$ and \mathbf{u} .
 - 3: Update the values of the variances of the random effects including σ_ϵ^2 and σ_q^2 .
 - 4: Iterate between steps (1) and (3) until criterion for a stopping rule has been satisfied.
-

Given the representation of (13) there is a temptation to assume a naive solution to the linear system by inversion of the leftmost matrix. This is in fact how the SAS software documents the solution to the mixed model in its technical manual (SAS Documentation 14.2, n.d.). This may be perhaps just the algebraic representation and not actually the computational implementation. Instead, the computational solution can be written to resemble the solutions to linear systems previously described and a sketch of the iterative process is provided in Algorithm (2). Begin by writing

$$\mathbf{L}\mathbf{L}' = \begin{bmatrix} \mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X} & \mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{Z} \\ \mathbf{Z}'\boldsymbol{\Omega}^{-1}\mathbf{X} & \mathbf{Z}'\boldsymbol{\Omega}^{-1}\mathbf{Z} + \mathbf{G}^{-1} \end{bmatrix}, \quad \mathbf{y}^* = \begin{bmatrix} \mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{y} \\ \mathbf{Z}'\boldsymbol{\Omega}^{-1}\mathbf{y} \end{bmatrix} \quad (14)$$

where elementwise calculation can be used for the diagonal matrices (e.g., $\Omega^{-1}X$) as presented in the WLS section and $\Theta' = [\beta' u']$. The model can now be written as $LL'\Theta' = y^*$ sharing the representation of the OLS form and also inherits the solution using (8) through the two triangular systems, $Lz = y^*$ for z and then $L'\Theta = z$ for Θ . Of course, this is only one iteration and solutions for the mixed model iterate with updated values for the variance components. Note, that when all elements of the matrix $G = 0$, this reduces to the partitioned fixed effects model.

Variance Estimation with Woodbury Identity and Decompositions for Matrix Inversions

Achieving scalability with iterative algorithms often requires finding less expensive ways to manage computations that are repeated. Linear mixed models are one such example where estimation depends on iterating between steps to compute variance components and then use those to estimate fixed parameters. It is one use case to further explore how better scalability can be achieved by reducing the computational overhead by rewriting a problem to take advantage of an inversion lemma, a decomposition, and recycling certain components. In the linear mixed model, obtaining updated values for the variance components at each iteration of the algorithm can be estimated by first computing (McCulloch & Searle, 2001)

$$T = \left[I + \left(Z'Z - Z'X(X'X)^{-1}X'Z \right) D \right]^{-1} \quad (15)$$

where I is an identity matrix with dimensions equal to the model matrix for the random effects and $D = \sigma_e^{-2}G$ and all other matrices are as above. The variances at each iteration $t - 1$ are computed (removing subscripts for iteration)

$$\sigma_e^2 = \frac{y'e}{N - p}, \quad \sigma_q^2 = \frac{\tilde{u}_q' \tilde{u}_q}{m_q - \text{tr}(T_q)} \quad (16)$$

where σ_e^2 is the residual variance, σ_q^2 is the marginal variance level q , \tilde{u}_q are the predictions of the random effects at level q , m_q are the number of units in level q , and $\text{tr}(T_q)$ is the trace of the matrix T_q , which is the block of the matrix T corresponding to the q th block.

The challenge here is that the square matrix T can be very large. As a result, it is extremely expensive to find its inverse inside an iterative problem, thus rendering this particular approach computationally prohibitive unless a reduceable structure can be identified. In some cases, there is an alternative method for taking the inverse of a matrix with this structure via the Woodbury Matrix Identity (Woodbury, 1950) that provides a very efficient route for computation. First, Equation (15) is rewritten as

$$T = \left(I + Z'ZD - Z'X(X'X)^{-1}X'ZD \right)^{-1} \quad (17)$$

where D is a diagonal matrix in the case of random intercepts and within each block q is constant along the diagonal. The Woodbury Identity provides that T can be computed as

$$T = A^{-1} + A^{-1}B(P - B'DA^{-1}B)^{-1}B'DA^{-1} \quad (18)$$

where $A = I + Z'ZD$, $B = Z'X$, and $P = X'X$. Now using (18) T can be computed by separately finding the inverse of A and then also of $P - B'DA^{-1}B$. The latter of the two has

very small dimensions, only $p \times p$ where p is the number of fixed effects in the model. The portion $\mathbf{I} + \mathbf{Z}'\mathbf{Z}\mathbf{D}$ is a bit more cumbersome because the dimensions of the model matrix for the random effects is $Q \times Q$, where Q is the number of columns in the model matrix for the random effects which can be extremely large.

The components \mathbf{I} and $\mathbf{Z}'\mathbf{Z}$ are fixed and never change over the iterative process, only \mathbf{D} is changing at each iteration. For nested random effects, $\mathbf{Z}'\mathbf{Z}$ is block diagonal but it has no special structure when the random effects are fully or partially crossed. This fact is motivation to explore if the special structures can be exploited in order to more easily find \mathbf{A}^{-1} given that it will be computed many times. Of course, it is a well-known result that the sum of inverses is not equal to the inverse of the sum. However, \mathbf{I} and \mathbf{D} are both diagonal within a block of q for nested random effects and have constant values along the diagonal. These structures make it possible to find \mathbf{A}^{-1} in a trivial way for the linear mixed model when dealing with nested random effects. Because $(\mathbf{Z}'\mathbf{Z})^{-1}$ is block-diagonal in this case, its inverse is found one block at a time.

It's pointed out here that the following works in the special case of nested random effects for random intercepts and is not globally generalizable (see appendix). However, the emphasis in this work is to identify special structures for a given scenario and exploit them to reduce computational burden. The inverse of the sum provided here may also have broader applications in other types of problems, such as when repeatedly finding the inverse of a Hessian matrix in iterative maximization problems, and so it is a motivating use case to document the condition.

For convenience, the subscript on \mathbf{Z} is dropped, but note that this performed on each block, or perhaps on large \mathbf{Z} , when it too has a special structure. First, precompute and store the eigendecomposition for $\mathbf{Z}'\mathbf{Z}$, which is the most expensive part of the calculation. Then, all subsequent iterations treat the eigenvectors as fixed and are reused and the eigenvalues are updated repeatedly with virtually no computational overhead. The inverse of the sum is then simply the inverse of the eigenvalues as

$$(\mathbf{Z}'\mathbf{Z}\mathbf{D} + \mathbf{I})^{-1} = \mathbf{A}^{-1} = \mathbf{Q}_1\boldsymbol{\lambda}^{*-1}\mathbf{Q}_1' \quad (19)$$

where the details and a proof for (19) are provided in the appendix. Now \mathbf{A}^{-1} is obtained in a much less expensive way by reusing, \mathbf{Q}_1 , and simply updating $\boldsymbol{\lambda}^*$ at each iteration, which is a nominal task given that it is diagonal. So, pulling together the Woodbury Identity and the decomposition derived for Equation (19), the inverse of the very big matrix \mathbf{T} is easily computed over and over again with little overhead.

Scalable Methods for Integration

Psychometricians routinely encounter challenging marginal likelihood problems or posterior distributions that require numerical integration methods. These methods approximate difficult integrals by summing over a finite number of points. Integrals with a gaussian factor are rather common in psychometric applications and, as such, the Gauss Hermite rule (GHR) may be used for integrals of the following form over the domain $(-\infty, \infty)$

$$\int f(\theta)e^{-\theta^2}d\theta \simeq \sum_{q=1}^Q f(\theta_q)w_q \quad (20)$$

where summation is over $\mathcal{P} = \{\theta_q, w_q\}_{q=1}^Q$, the so called nodes, θ_q , and weights, w_q , for $q = 1, 2, \dots, Q$. The choice of nodes and weights in \mathcal{P} is optimal under GHR and will generally outperform other methods (Quarteroni, Saleri, & Gervasio, 2010). Unfortunately, psychometric practice seems to have codified the basic rectangular rule with fixed, equally spaced points from the domain of $(-4, 4)$ (Bock & Mislevy, 1982) despite its serious limitations in most cases for at least two reasons. First, the equally spaced points may not be sampled in the best region of the integrand and second it lacks the computational advantage of precision until an extremely large number of fixed points are used. Building scalable applications requires approaches with minimal computational overhead that achieve precision rapidly.

One-Dimensional Gaussian Quadrature

Conveniently, psychometricians encounter scenarios where the weight function is a standard normal density

$$\int f(\theta) \frac{e^{-\theta^2/2}}{\sqrt{2\pi}} d\theta \simeq \pi^{-1/2} \sum_{q=1}^Q f(\sqrt{2}\mathbf{L}\theta_q + \tilde{\theta})w_q \quad (21)$$

where $\tilde{\theta}$ is a location parameter for centering (Liu & Pierce, 1994)². The approximation in (21) implements a transformation on θ from (20) and adapts to the best region of the integrand—hence it is the adaptive GHR. Scalability and accuracy in evaluating integrals depends not only on the number of quadrature points, but also on the location of those points and the nature of the function surface being evaluated (Lesaffre & Spiessens, 2001).

Multivariate Gaussian Quadrature

The one-dimensional GHR can be extended to the case for integrals of the general form in multiple dimensions

$$\int \dots \int f(\boldsymbol{\theta})\psi(\boldsymbol{\theta})d\boldsymbol{\theta}. \quad (22)$$

where $\boldsymbol{\theta} = \{\theta_1, \dots, \theta_p\}'$ and continue with the assumption $\psi(\boldsymbol{\theta})$ is multivariate normal. The integral in (22) is often intractable and requires approximation. This may generally be approximated with nodes rotated to reflect a covariance structure and possibly shifted over a mean. That rotation can be implemented with the Cholesky factor of a covariance matrix as (Chowdhary, Salloum, Debusschere, & Larson, 2015; Jäckel, 2005; Judd, Maliar, & Maliar, 2011; Stringer, 2021)

$$\int \dots \int f(\boldsymbol{\theta})\psi(\boldsymbol{\theta})d\boldsymbol{\theta} \simeq \pi^{-P/2} \sum_{q_1=1}^Q \dots \sum_{q_p=1}^Q f(\sqrt{2}\mathbf{L}(\theta_{q_1}, \dots, \theta_{q_p})' + (\tilde{\theta}_{q_1}, \dots, \tilde{\theta}_{q_p})')w_q^* \quad (23)$$

where θ_{qp} is node $q = 1, 2, \dots, Q$ in dimension $p = 1, 2, \dots, P$, w_{qp} is the corresponding weight, $w_q^* = \prod_{p=1}^P w_{qp}$, and $\tilde{\theta}_{qp}$ is used to mean some points around which the nodes are centered. The nodes are expansions of the points in the one-dimensional set, \mathcal{P} , via a cartesian product rule to form a multivariate grid, \mathbf{G} , with dimensions $Q^P \times P$. Hence, the q th row of \mathbf{G} is $\mathbf{g}_q =$

²The \mathbf{L} in (21) is actually a scalar, σ , but is left as a matrix so it appears as a special case of (23).

$(\theta_{q1}, \dots, \theta_{qp})$. Note that (21) is simply a special case of (23), hence no new concepts are needed for implementation beyond what has been described for the one-dimensional scenario.

The extremely large number of quadrature points is what renders this approach challenging as, hueristically, this is effectively a two step process. Step 1 requires evaluating the function $f(\cdot)$ over all Q^P nodes and then step 2 is the summation of (23) making use of those values over all rows of the grid G . Even with a modest selection for Q , the exponential growth in the function evaluations quickly becomes large with increasing P . Hence, we might explore ways in which we evaluate $f(\cdot)$ as few times as possible and whether fewer function evaluations than $Q^P \times P$ can be used. These options are explored in the applications section.

General Purpose Optimization Tools

Perhaps the most pervasive need for psychometric application building is a collection of tools for optimizing likelihoods. The likelihood theory of inference is well-developed (King, 1998) and found across all corners of psychometric practice. The common challenge is that likelihood functions for psychometric problems rarely have analytic forms and so iterative numerical approaches are used for maximum likelihood estimation (MLE). Two methods are explored here with a reminder that the stochastic gradient descent method is also a general purpose tool, but was presented within the section on linear models.

Maximization with Newton-Raphson

Optimization problems in psychometrics generally involve a multi-parameter likelihood function intending to simultaneously maximize all parameters. In this instance, let $\mathcal{L}(\theta)$ denote a twice differentiable likelihood function with respect to the parameters θ with gradient $\nabla \mathcal{L}(\theta)$ and Hessian, $\mathcal{H}(\theta)$. The Newton-Raphson procedure maximizes $\mathcal{L}(\theta)$ by choosing initial starting values for θ and then iteratively updating the parameters via

$$\theta_{t+1} = \theta_t - \mathcal{H}_t(\theta)^{-1} \nabla \mathcal{L}_t(\theta) \quad (24)$$

where the subscript t denotes the values of the gradient and Hessian at iteration t and continues until a convergence criterion is reached. The notation, $\mathcal{H}_t(\theta)^{-1}$, implies inversion of this matrix at each iteration. However, as previously discussed, the operation should be performed using a numerically stable decomposition. Additionally, because $\mathcal{H}_t(\theta)^{-1}$ is found iteratively and may be expensive, it may be feasible to find a trivial way to perform this operation similar to the eigen-decomposition method in Equation (35) provided in the appendix that takes advantage of precomputing and reusing certain components. While the method in the appendix is unique to the mixed model scenario, it is intended to be a concept that can be generalized to other situations where a matrix inverse computed iteratively can be made less expensive. Quasi-Newton methods, such as the BFGS algorithm (Fletcher, 1987), are other less expensive options that are quite simple to implement and avoid directly computing $\mathcal{H}_t(\theta)^{-1}$, which may be difficult in some scenarios.

It can be pointed out that the SGD approach in (12) bears similarity to Newton's method of (24) where the learning rate of SGD is replaced by a Hessian in Newton's method. Conceptually, Newton's method will more directly head towards the optimal value and SGD will essentially wiggle its way there. However, the cost of Newton's approach is much greater than that of the SGD.

Some work combining the faster convergence of Newton’s method with the reduced computational complexity of SGD is available for further study (Sohl-Dickstein, Poole, & Ganguli, 2014; Zhou, Wei, Zhang, & Zheng, 2021).

Expectation-Maximization Algorithm

There is no approach more ubiquitous across psychometrics than the expectation-maximization (EM) algorithm (Dempster, Laird, & Rubin, 1977). The EM algorithm is based on the idea that there exists some complete data, $\mathbf{z} = \{\mathbf{x}, \boldsymbol{\eta}\}$, and a complete data likelihood, $\mathcal{L}(\boldsymbol{\theta}; \mathbf{z}) = \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \boldsymbol{\eta})$, that could be maximized with respect to $\boldsymbol{\theta}$ if the complete data were observed. The problem is that only \mathbf{x} is observed and $\boldsymbol{\eta}$ are latent; hence the values of $\boldsymbol{\eta}$ represent missing data and so maximizing $\mathcal{L}(\boldsymbol{\theta}; \mathbf{z})$ cannot easily occur.

The EM algorithm ameliorates the missing data problem by using provisional, or conditional expected values, of $\boldsymbol{\eta}$ and then maximizing the complete data likelihood. Algorithm (3) outlines the alternating steps used to implement EM.

Algorithm 3 Expectation Maximization Algorithm

Input: Create starting values for $\boldsymbol{\theta}$.

- 1: E-step is to compute expected values given the observed data and the provisional parameter estimate, $Q(\boldsymbol{\theta}; \boldsymbol{\theta}_t) = \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \boldsymbol{\eta}) | \mathbf{x}, \boldsymbol{\theta}_t]$.
 - 2: M-step is to find $\arg\max_{\boldsymbol{\theta} \in \Omega} Q(\boldsymbol{\theta}; \boldsymbol{\theta}_t)$ where $\boldsymbol{\theta}$ is in the parameter space Ω .
 - 3: Iterate between steps (1) and (2) until criterion for a stopping rule has been satisfied.
-

There are many examples of the EM algorithm in the psychometric literature for interested readers to further explore (Hsu, Ackerman, & Fan, 1999). One very simple example that can be quickly illustrated is to reconsider estimating the linear mixed model of the form $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$. Here, the complete data are $\{\mathbf{y}, \mathbf{u}\}$ and if fully observed, the parameters of the mixed model could be found using the MLE variants in Equation (16) for the variance components and Equation (8) for estimates of the fixed effects³. The random effects are in reality missing, but the joint distribution of the complete data, $\{\mathbf{y}, \mathbf{u}\}$, is assumed multivariate normal and standard distribution theory for the multivariate normal provides that $\mathbb{E}(\mathbf{u} | \mathbf{y})$ can be easily determined (McCulloch & Searle, 2001) and used to update a new value for the M step.

There are some occasions where the M-step may not have a closed form or data might be so large the maximization step could be daunting even when working values for the missing data are easily obtained. In fact, it is somewhat challenging to use EM at scale without some additional considerations. One very scalable option is to separate the computations when possible. For instance, there are occasions in psychometrics where the expected values used in the E-step can be computed independently, representing what is known as an embarrassingly parallel problem and the E-step (and possibly also the M-step) can be executed in parallel (Lee, Leemaqz, & McLachlan, 2016).

An interesting variant of the EM algorithm would be to use SGD in the M-step to manage a large data problem if the E-step was easily available. For example, in the case of the mixed model above, assume $\mathbb{E}(\mathbf{u} | \mathbf{y})$ is available through an efficient E-step, but maximization via (8) is hard

³If the random effects were observed, the least squares solution is justified as $\mathbf{y} - \mathbf{Z}\mathbf{u} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$

because the dimensions of the model matrix \mathbf{X} are extremely large. Instead of computing the M-step by (8), the SGD could be used in the M-step as a way to manage the process allowing computation to be highly scalable. An illustration for how this could be implemented for mixed models is provided in the appendix. Other occasions may arise, for instance, in generalized linear mixed models that involve an integral for the M-step and one of the tools discussed above for numerically approximating that integral might be used. In other cases, maximization may require an iterative solution, thus pulling in the Newton-Raphson procedure as a way to solve the M-step.

Applications to Psychometric Problems

The previously discussed methods are computational templates that can be used to build and scale psychometric applications. Three examples are provided here with supporting R code in a Github repository available upon request. The intent is not to center attention on these specific examples; rather, these examples are used because they illustrate how the preceding templates can be used to compute something complex that is not routinely found in standard software. For this reason, they are attractive for exploring computationally efficient approaches.

Example 1: Estimating the Error in Variables Linear Model

Estimates of examinee performance from a psychometric instrument are noisy reflections of a true, latent trait. The observed measures are subsequently used in various ways and often used as inputs in regression models of some form. When using these observed estimates, the assumptions of the Gauss-Markov theorem do not hold and the parameter estimates are then inconsistent unless corrections are made (Doran, 2014; Lockwood & McCaffrey, 2020; Nab, van Smeden, Keogh, & Groenwold, 2021).

These models are chosen as an example because with different assumptions about the nature of the measurement error (e.g., homoscedastic versus heteroscedastic), the models may take different forms and custom software is often written to obtain parameter estimates. One such example is the `mecor` (Nab et al., 2021) package in R and a code review of its implementation shows that the algebraic representation and computational representation are one and the same.

One form of the linear EiV regression is commonly shown to be (StataCorp, n.d.)

$$(\mathbf{X}'\mathbf{X} - \mathbf{S})\boldsymbol{\beta} = \mathbf{X}\mathbf{y} \quad (25)$$

where, in the case of known estimates for the reliability, the matrix $\mathbf{S} = \text{diag}(N(1-r_1)\sigma_1^2, \dots, N(1-r_p)\sigma_p^2)$ where r_p is the reliability of the p th variable, N is the number of individuals, and σ_p^2 is the sample variance of the p th variable in the model matrix \mathbf{X} . Now, find $\mathbf{X}'\mathbf{X} - \mathbf{S} = \mathbf{L}\mathbf{L}'$ and obtain estimates using (8). It's useful to note that (25) can be extended for the linear mixed model and the computational approach previously shown can be applied.

Example 2: Estimating a High-Dimensional EAP

Suppose we are interested in estimating the ability of an examinee using a P -dimensional EAP estimator. This is chosen as an example because it is a desirable estimator that continues to be chal-

lenging for psychometricians to implement (Chalmers, 2012; Ferrando & Lorenzo-Seva, 2016). Let the expected value of the posterior be

$$\hat{\theta} = \frac{\int \cdots \int \theta \mathcal{L}(\theta; \mathbf{z}) \psi(\theta; \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\theta}{\int \cdots \int \mathcal{L}(\theta; \mathbf{z}) \psi(\theta; \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\theta} \quad (26)$$

where $\psi(\theta; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the multivariate normal density function and the likelihood function (see appendix) is a product of domain-specific likelihoods as $\mathcal{L}(\theta; \mathbf{z}) = \prod_{p=1}^P \mathcal{L}(\theta_p; \mathbf{z}_p)$. Individuals are assumed to have some idiosyncratic ability contributing to their response on each dimension and that some test items are uniquely associated with dimension 1 (θ_1), others with dimension 2 (θ_2), and so on for the P th dimension, $\theta = (\theta_1, \theta_2, \dots, \theta_P)'$.

The likelihood function and its corresponding population distribution in (26) are not in the same parametric family, as such the integral cannot be evaluated analytically. Here three options are explored that make for possible implementation that vary in the number of times the function $\mathcal{L}(\theta; \mathbf{z})$ is computed. Once the likelihood is evaluated, then step 2 is effectively direct implementation of the sum over all rows of \mathbf{G} . Sparse grids or “pruning” some values are options for reducing the number of points over which the function is evaluated, but those are not explored here.

The first canonical option is to use (23) directly. This involves using the Cholesky factor to adjust the nodes to absorb the population covariance structure and then, possibly shift the nodes over a new centered point. This works quite well; however, has one downside worth noting that makes this computationally somewhat less attractive.

Initially (before rotating with the Cholesky factor), the nodes in \mathbf{G} are patterned such that marginally (i.e., each column in \mathbf{G}) the Q unique nodes are repeated over and over. When the rows of this matrix are rotated via $\mathbf{L}(\theta_{q1}, \dots, \theta_{qp})'$, then the uniqueness is lost exponentially over each column. That is, column 1 has Q unique nodes repeated, column 2 has only Q^2 unique nodes, column 3 would have Q^3 unique nodes and so on up to the P th column.

This reduced pattern of unique nodes creates a larger burden in that we must compute the values of the likelihood at each node and that IRT likelihood is somewhat expensive as it involves computing exponents. In column 1, it is feasible to compute the value of the likelihood at the Q unique points and then simply copy those values into the other positions in column 1 where the same value of the node appears. That is, we would compute the likelihood only Q times and then recycle those values by copying them into other positions where the same node appears in the grid. In column 2, we must compute the likelihood at all Q^2 unique nodes and then the values could then be copied into the positions of this column where the same node appears again. The upside here is that the method yields a good approximation of the integral. The downside is the computational expense is very large. In this case, we would evaluate the likelihood $Q + Q^2 + Q^3, \dots, Q^P$ times.

A second option for computing this integral could be to use a stochastic quadrature routine. In this process the nodes are random draws from $\mathcal{N}_p(\mathbf{M}, \mathbf{I}_P)$, where \mathbf{I}_P is a P -dimensional identity matrix again creating the grid \mathbf{G} , but in this instance it would have dimensions $R \times P$ where R is the number of random variates chosen and it is not an exponential function of R . In this manner, the variates in \mathbf{G} would reflect the population covariance after being premultiplied by the Cholesky factor (see introductory section) when the summation in (23) is used to evaluate the integral. This stochastic approach makes the computation reasonable and its precision is highly dependent on the law of large numbers such that $\hat{\theta}_p \xrightarrow{a.s.} \theta_p$ as $R \rightarrow \infty$ ⁴. In this case, we would evaluate the

⁴The asymptotics here also depend on the number of test items being very large and not only the number of variates

likelihood $R \times P$ times.

The GHR tends to behave well when the function being evaluated approximates a low order polynomial (Liu & Pierce, 1994). However, in situations where an examinee responds to many test items, the function could spike in some areas causing for GHR to perform badly as few quadrature points could live within a region of interest. Perhaps the integral within (26) can be rewritten so that it behaves more like a low order polynomial (Rabe-Hesketh, Skrondal, & Pickles, 2002) and a third approach can be explored as

$$\int \cdots \int \mathcal{L}(\boldsymbol{\theta}; \mathbf{z}) \psi(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) d\boldsymbol{\theta} = \int \cdots \int \frac{\mathcal{L}(\boldsymbol{\theta}; \mathbf{z}) \psi(\boldsymbol{\theta}; \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\tilde{\psi}(\boldsymbol{\theta}; \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})} \tilde{\psi}(\boldsymbol{\theta}; \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) d\boldsymbol{\theta} \quad (27)$$

and then inserting back into (26) gives the approximation

$$\hat{\theta}_p \simeq \frac{\sum_{q=1,1}^Q \cdots \sum_{q=p,1}^Q \theta_{p,q} \mathcal{L}(\theta_{1,q}; \mathbf{z}_1) \mathcal{L}(\theta_{2,q}; \mathbf{z}_2) \cdots \mathcal{L}(\theta_{P,q}; \mathbf{z}_p) \phi(\boldsymbol{\theta}_q; \boldsymbol{\mu}', \boldsymbol{\Sigma}') w_q^*}{\sum_{q=1,1}^Q \cdots \sum_{q=p,1}^Q \mathcal{L}(\theta_{1,q}; \mathbf{z}_1) \mathcal{L}(\theta_{2,q}; \mathbf{z}_2) \cdots \mathcal{L}(\theta_{P,q}; \mathbf{z}_p) \phi(\boldsymbol{\theta}_q; \boldsymbol{\mu}', \boldsymbol{\Sigma}') w_q^*} \quad (28)$$

where

$$\phi(\boldsymbol{\theta}_q; \boldsymbol{\mu}', \boldsymbol{\Sigma}') = \frac{\psi(\boldsymbol{\theta}_q; \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\tilde{\psi}(\boldsymbol{\theta}_q; \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})} \quad (29)$$

and $\tilde{\psi}(\cdot)$ is a proposal distribution. This expression resembles an importance sampling approach (Gelman, Carlin, Stern, & Rubin, 2004) and has been termed importance Gauss Hermite (IGH) method (Elvira, Martino, & Closas, 2021) and other applications have used similar ideas (Antal & Oranje, 2007; Tuerlinckx, Rijmen, Verbeke, & De Boeck, 2006; Pinheiro & Bates, 1994). The value of the IGH approach is that the weights are modified to have greater importance for certain values of the integrand, in which case fewer quadrature points may be needed to evaluate the integral. The trade-off is that some calculations involving the densities in $\phi(\cdot)$ is required, though that additional cost is seemingly smaller than if more quadrature points were used and importance-based methods may be valuable for reducing computational burden (Ackerberg, 2009).

The additional benefit is that the original (unrotated) Gauss Hermite nodes are used in (28). We can exploit the uniqueness in the repeating values of the nodes in each column and compute the value of the likelihood only Q times in each dimension and then copy those values into other positions of the column where they repeat. As a result, the likelihood is evaluated only $Q \times P$ times and not $Q^P \times P$ times! Assuming $Q < R$, this approach evaluates the likelihood the smallest number of times of all three cases.

Naturally, the question of what values are used in $\tilde{\psi}(\cdot)$ needs clarification. This choice is what has the largest impact relative to the computational burden, yet unfortunately, some experimentation is needed for choosing these values. The R code implementation available for this example uses $\tilde{\psi}(\boldsymbol{\theta}; \mathbf{0}, \mathbf{I})$ as a template and some tuning of that distribution can be explored.

One final consideration for reducing some computational burden in problems of this nature might be to use an iterative process as suggested by Naylor and Smith (1982). If values for the parameters of interest were known, or at least approximately known, the nodes could be recentered via the transformation $\sqrt{2}\theta_{p,q}\sigma_p + \hat{\theta}_p$ and take another pass over the integral with the transformed nodes and weights at an increased number of quadrature points for precision could be used.

Some additional computational benefits could easily be realized for this problem. For instance, computing the likelihood values over each column in the node array is independent and consequently is embarrassingly parallel, and the results could then be joined to take summations. In fact, test scoring applications in psychometrics are an embarrassingly parallel problem as the test score for examinee k does not depend on the scores for others, it depends only on the item parameters and their item responses, both of which are treated as fixed.

Example 3: A Simple ML Classifier with SGD

Imagine a group of test item developers would benefit from knowing whether the new items being developed are likely to have examinee response times that are fast or slow before those items have any field test and timing statistics. Such forecasting systems could guide item development plans and item writing workshops and be helpful in an array of test construction activities. Machine learning applications are premised on the idea that computers can be trained to learn from existing data and subsequently apply the results to form future classifications, or predictions, and such an application is suitable here.

For purposes of the example, assume response time data from a testing program exist for each test item and the data could be concatenated over years, over grades, over test forms, and over all examinees yielding an extremely large data file forming N total observations. In addition, assume all items in the set have a collection of observable features that are predictive of response time, such as readability indices, stem word length, content domain, cognitive complexity, item type and so on. Further, assume a binary outcome is associated with item response time such that $y_i = 1$ if the response to the i th item is fast and $y_i = 0$ otherwise and that the many item features can be used to form the model matrix \mathbf{X} to serve as covariates in the regression. Response time is of course a real positive number, but here we bin into two categories to build an illustration that differs from the prior linear regression model using SGD.

The example described here yields an extremely large data file from which we intend to learn patterns and then make future predictions on other data not included in the training sample. Because the available data to be used for training is so large, logistic regression using something like iteratively reweighted least squares is not feasible. In this case, an alternative classifier can be constructed using SGD for training allowing for computation to occur in an efficient way.

To show a contrast, in batch gradient descent with N observations, L parameters, and using K epochs, the gradient would be evaluated $N \times K \times L$ times. The number for K is generally arbitrary, but is often selected as a very large value. Instead, SGD can reduce this computational burden and evaluate only $K \times L$ gradient terms to achieve comparable results. In fact, because SGD randomly samples from the full set of observations, it may be helpful to only randomly stream in portions of the full available data and sample from within those portions repeatedly. This could alleviate some overhead space needed to be reserved for computing. If mini-batch gradient descent were used, then $N_i \times K \times L$ gradient terms would be evaluated, where N_i represents a subset of N . The reduction in the number of gradient evaluations is impressive and makes very large problems feasible. Simply to be conceptual, imagine $N = 1,000,000$, $K = 1000$, and $L = 3$; then, batch gradient descent would evaluate 3 billion gradient terms whereas SGD would only evaluate 3,000 gradient terms!

Here, let the objective function for the logistic model be $\mathcal{J}_i(\boldsymbol{\theta}) = -y_i \log(z_i) + (1 - y_i)(1 - z_i)$

and then

$$\nabla \mathcal{J}_i(\boldsymbol{\theta}) = (z_i - y_i)' \mathbf{X}_i \quad (30)$$

where the following familiar logistic (sigmoid) function is used

$$\Pr(y_i = 1 | \mathbf{X}_i, \boldsymbol{\beta}) = z_i = \frac{1}{1 + \exp(-(\mathbf{X}_i \boldsymbol{\beta}))}. \quad (31)$$

The gradient in (30) can be used in Equation (12) and the steps in Algorithm (1) would be implemented for optimization. Generalizations of this concept could be further developed to explore IRT forecasting problems where a sense of the future item parameters could be obtained from features associated with items or in training a learning management system to associate targeted instructional content to an examinee based on their test performance.

Discussion

The ideas here are but a subset of useful numerical methods. However, the subset represents a core set of methods that can be combined or generalized and broadly applied in psychometric application building. Hopefully, these recipes encourage even further study and conversations in methods that are necessary ingredients to support the ideas of computational psychometrics. While some specific recipes for computing are presented here, it's important not to lose sight of the primary motivation and aim of this work. Specifically, scalable psychometric solutions will be more viable with better advancement towards efficient computational recipes and no longer conflating algebraic representations to also be the computational solutions. The aim of any computational psychometrician should be to exploit any special structures and find reduceable solutions to large problems.

There are a few possible ways in which the work described here might be useful. First, many psychometricians find jobs in industry where they are charged with creating psychometric applications to support their organizational infrastructure and client requirements. The applications used by assessment companies deploy tests to many hundreds of thousands (even millions) of examinees and the operational systems used for psychometric work must be fast, accurate, and highly scalable. Second, perhaps in academia, the ideas here can supplement other course materials borrowed from computational statistics that are more mathematical in nature. This work is intended to create connections to computational psychometrics and may be a helpful way to bridge the knowledge base between general numerical methods and psychometric theory. Third, many commercial psychometric applications need some retrofitting to better support larger computational demands and the ideas here can be potentially used to support those upgrades. Last, psychometricians setting out on their own course of study to improve their computational acumen might benefit from implementing the ideas proposed here as test cases to explore and replicate.

This work is also intended to fit into a larger conversation of the future of psychometrics. The skill set for building scalable psychometric applications means innovators need foundational numerical approaches, but also skills in software development to prototype, finalize, and deploy at scale applications to support those ideas. Commercial cloud-based infrastructures are more readily available and provide platforms to take advantage of distributed computing. Nonetheless, the first step should be to build an efficient method and then secondly scale that method with computational advantages offered.

Appendix A: Computational Formulae

Proof that $(Z'ZD + I)^{-1} = Q_1\lambda^{*-1}Q_1'$

Let $Z'Z = Q_1\lambda_1Q_1'$ and $D = Q_2\lambda_2Q_2' = \lambda_2 I$ given that $D = dI$, then

$$Z'ZD + I = (Q_1\lambda_1Q_1')(Q_2\lambda_2Q_2') + I \quad (32)$$

$$= (Q_1\lambda_1Q_1')\lambda_2 + I \quad (33)$$

The commutative property for matrices holds if and only if dI , hence $Q_1'\lambda_2 = \lambda_2Q_1'$. Then, $(Q_1\lambda_1Q_1')\lambda_2 + I = (Q_1\lambda_1\lambda_2Q_1') + I$. Finally, we can take advantage of the identity matrix and co-diagonalize its elements alongside the eigenvalues as

$$Q_1\lambda_1\lambda_2Q_1' + I = Q_1(\lambda_1\lambda_2 + I)Q_1' = Q_1\lambda^*Q_1' \quad (34)$$

where $\lambda^* = \text{diag}\{\lambda_{11}d + 1, \lambda_{12}d + 1, \dots, \lambda_{1n}d + 1\}$ and λ_{1n} are the elements of λ_1 . Finally, the inverse of the sums is simplified to

$$(Z'ZD + I)^{-1} = Q_1\lambda^{*-1}Q_1'. \quad (35)$$

Equation (35) is useful in iterative algorithms as it precomputes an expensive component, Q_1 , initially and then it is repeatedly reused in subsequent iterations. Because updating λ^* is simple and it is diagonal, it permits a trivial inverse. Collectively, the special structures in this instance are exploited to reduce a larger computational burden into one that is much less expensive within iterative algorithms.

General Remark on the Proof

It is important to note that the preceding proof applies under conditions when $D = dI$, in which case the commutative property holds between D and some other (conformable) matrix. However, when that condition is not true, we arrive at a slightly different problem. Suppose D is diagonal, but the elements along the diagonal are not constant. Then, we may write $(Z'ZD + I)^{-1} = (Q\lambda Q' + I)^{-1} = Q(\lambda + I)^{-1}Q'$ where now $Z'ZD = Q\lambda Q'$.

This is interesting, but not entirely helpful within an iterative algorithm as D cannot be factored out of the decomposition. Hence, it requires computing the decomposition $Z'ZD$ at each iteration even though only D is changing. At the current time, there is no accepted method for the inverse of matrix sums under this condition. There are some concepts that may later be considered using eigenvalue perturbations. That is, if we can first compute the decomposition $Z'Z$, then update this with the changing values of D to yield an approximation of $Q\lambda Q'$, then we might be able to efficiently find the sum of the inverse when D is diagonal, but not with constant elements. This remains a numerical challenge to explore.

Likelihood Function Details

Each individual domain likelihood, $\mathcal{L}(\theta_P; \mathbf{z}_p)$, is composed of a mixture of binary and polytomous test items with known item parameters. Let z_{ijp} denote the observed response of the i th examinee to the j th item in the p th dimension, then

$$\mathcal{L}(\theta_P; \mathbf{z}_p) = \mathcal{L}_1(\theta_P; \mathbf{z}_p) \mathcal{L}_2(\theta_P; \mathbf{z}_p)$$

$$\mathcal{L}_1(\theta_P; \mathbf{z}_p) = \prod_{j \in p} \left[c_j + \frac{1 - c_j}{1 + \exp[-Da_j(\theta_P - b_j)]} \right]^{z_{ijp}} \left[1 - \left(c_j + \frac{1 - c_j}{1 + \exp[-Da_j(\theta_P - b_j)]} \right) \right]^{1 - z_{ijp}}$$

where $j \in p$ is used to mean there is a collection of j items uniquely associated with dimension p , c_j is the lower asymptote of the item response curve (i.e., the guessing parameter), a_j is the slope of the item response curve (i.e., the discrimination parameter), b_j is the location parameter, and D is a constant, by default fixed at 1.7. Then the items scored in multiple categories takes the form of the graded response model as

$$\mathcal{L}_2(\theta_P; \mathbf{z}_p) = \prod_{j \in p} \Pr(z_{ijp} | \theta_P)$$

$$\Pr(z_{ijp} | \theta) = \begin{cases} \frac{1}{1 + e^{Da_j(\theta_P - b_{j1})}}, & z_{ijp} = 0 \\ \frac{\frac{1}{1 + e^{Da_j(\theta_P - b_{j,z+1})}} - \frac{1}{1 + e^{Da_j(\theta_P - b_{jz})}}}{\frac{1}{1 + e^{-Da_j(\theta_P - b_{jK})}}}, & 0 < z_{ijp} < K \\ \frac{1}{1 + e^{-Da_j(\theta_P - b_{jK})}}, & z_{ijp} = K \end{cases} \quad (36)$$

where b_K denotes the k th step and the other definitions used for the binary model apply here.

Efficient Parallel E-Step Example for Linear Mixed Model

Let the joint distribution of the fixed and random effects be written as a multivariate normal

$$\begin{bmatrix} \mathbf{u} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left[\begin{bmatrix} \mathbf{0} \\ \mathbf{X}\beta \end{bmatrix}, \begin{bmatrix} \mathbf{G} & \mathbf{GZ}' \\ \mathbf{ZG} & \mathbf{V} \end{bmatrix} \right] \quad (37)$$

where $\mathbf{V} = \mathbf{ZGZ}' + \mathbf{\Omega}$. Standard distribution theory provides that the conditional values of the random effects from this joint multivariate normal are $\mathbb{E}[\mathbf{u} | \mathbf{y}] = \mathbf{GZ}'\mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\beta)$. For two-level nested models with random intercepts only, this can be written as follows for the j th group with N_j denoting the number of units in the group

$$\mathbb{E}[\mathbf{u}_j | \mathbf{y}_j] = \tilde{\mathbf{u}}_j = \left[\frac{\sigma_q^2}{\sigma_e^2 + N_j \sigma_q^2} \right] \sum_{i \in j} (y_i - \mathbf{X}_i \beta) \quad (38)$$

where \mathbf{X}_i is the i th row in the model matrix \mathbf{X} and σ_e^2 and σ_q^2 are the residual variance and the marginal variance of the random effects at level q , respectively. Equation (38) can be computed as an embarrassingly parallel problem for each of the j groups. Assembling all $\tilde{\mathbf{u}}_j$ to use as the provisional working values for the missing data, we can now proceed with the M-step using stochastic gradient descent for least squares as shown in the section on linear models.

References

- Ackerberg, D. (2009). A new use of importance sampling to reduce computational burden in simulation estimation. *Quant Mark Econ*, 7, 343-376.
- ACTNext. (2019). *Computational psychometrics: A field guide* (Tech. Rep.). <https://actnext.org/research-and-projects/computational-psychometrics-field-guide/>: ACT.
- Allen, N., Donoghue, J., & Schoeps, T. (2001). *The NAEP 1998 technical report* (Tech. Rep.). <https://nces.ed.gov/nationsreportcard/pubs/main1998/2001509.asp>: U.S. Department of Education. Office of Educational Research and Improvement. National Center for Education Statistics.
- Antal, T., & Oranje, A. (2007). *Adaptive numerical integration for item response theory* (Tech. Rep.). <https://files.eric.ed.gov/fulltext/EJ1111562.pdf>: ETS.
- Bates, D. (2004a, June). Least squares calculations in R. *R News*, 4(1), 17–20.
- Bates, D. (2004b). *Sparse matrix representations of linear mixed models* (Tech. Rep.). <http://www.stat.wisc.edu/bates/reports/MixedEffects.pdf>: University of Wisconsin-Madison.
- Bock, R. D., & Mislevy, R. J. (1982). Adaptive EAP estimation of ability in a microcomputer environment. *Applied Psychological Measurement*, 6(4), 431-444.
- Cai, L., Yang, J., & Hansen, M. (2011). Generalized full-information item bifactor analysis. *Psychological methods*, 16(3), 221-248.
- Chalmers, R. (2012, 04). Mirt: A multidimensional item response theory package for the R environment. *JSS Journal of Statistical Software*, 48.
- Chowdhary, K., Salloum, M., Debusschere, B., & Larson, V. E. (2015). Quadrature methods for the calculation of subgrid microphysics moments. *Monthly Weather Review*, 143(7), 2955-2972.
- Cizek, P., & Cizkova, L. (2004). Iterative methods for solving linear systems. In J. Gentle, W. Härdle, & Y. Mori (Eds.), *Handbook of computational statistics* (p. 120-126). New York: Springer.
- Davis, T., & Hager, W. (2005, 01). Row modifications of a sparse cholesky factorization. *SIAM J. Matrix Analysis Applications*, 26, 621-639.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B*, 39, 1-38.

- Doran, H. (2014). Methods for incorporating measurement error in value-added models and teacher classifications. *Statistics and Public Policy*, 1(1), 114-119. Retrieved from <https://doi.org/10.1080/2330443X.2014.955228>
- Doran, H., Bailey, P., Buehler, E., & joo Lee, S. (2021). Dire: Linear regressions with a latent outcome variable [Computer software manual]. Retrieved from <https://CRAN.R-project.org/package=Dire> (R package version 1.0.3)
- Elvira, V., Martino, L., & Closas, P. (2021). Importance gaussian quadrature. *IEEE Transactions on Signal Processing*, 69, 474-488. Retrieved from <http://dx.doi.org/10.1109/TSP.2020.3045526>
- Ferrando, P., & Lorenzo-Seva, U. (2016). A note on improving EAP trait estimation in oblique factor-analytic and item response theory models. *Psicologica*, 37(2), 235-247.
- Fletcher, R. (1987). *Practical methods of optimization* (Second ed.). New York, NY, USA: John Wiley & Sons.
- Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (2004). *Bayesian data analysis* (2nd ed. ed.). Chapman and Hall/CRC.
- Genz, A., & Kass, R. E. (1997). Subregion-adaptive integration of functions having a dominant peak. *Journal of Computational and Graphical Statistics*, 6(1), 92-111.
- Hsu, Y., Ackerman, T. A., & Fan, M. (1999). The relationship between the Bock-Aitkin procedure and the EM algorithm for IRT model estimation..
- Jäckel, P. (2005). *A note on multivariate gauss-hermite quadrature*. (Tech. Rep.). <http://www.jaeckel.org/ANoteOnMultivariateGaussHermiteQuadrature.pdf>.
- Judd, K. L., Maliar, L., & Maliar, S. (2011). Numerically stable and accurate stochastic simulation approaches for solving dynamic economic models. *Quantitative Economics*, 2(2), 173-210. Retrieved from <https://onlinelibrary.wiley.com/doi/abs/10.3982/QE14>
- Kessy, A., Lewin, A., & Strimmer, K. (2018). Optimal whitening and decorrelation. *The American Statistician*, 72(4), 309-314. Retrieved from <https://doi.org/10.1080/00031305.2016.1277159>
- King, G. (1998). *Unifying political methodology: The likelihood theory of statistical inference*. Ann Arbor: University of Michigan Press.
- Laird, N. M., & Ware, J. H. (1982). Random-effects models for longitudinal data. *Biometrics*, 38(4), 963-974. Retrieved from <http://www.jstor.org/stable/2529876>

- Lee, S. X., Leemaqz, K. L., & McLachlan, G. J. (2016). A simple parallel EM algorithm for statistical learning via mixture models. In *2016 international conference on digital image computing: Techniques and applications (dicta)* (p. 1-8).
- Lesaffre, E., & Spiessens, B. (2001). On the effect of the number of quadrature points in a logistic random-effects model: An example. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 50(3), 325–335. Retrieved from <http://www.jstor.org/stable/2680873>
- Lira, M., Iyer, R., Trindade, A., & Howle, V. (2016). QR versus cholesky: A probabilistic analysis. *International Journal of Numerical Analysis and Modeling*, 13(1), 114–121. (Publisher Copyright: © 2016 Institute for Scientific Computing and Information.)
- Liu, Q., & Pierce, D. A. (1994). A note on gauss-hermite quadrature. *Biometrika*, 81(3), 624–629. Retrieved from <http://www.jstor.org/stable/2337136>
- Lockwood, J. R., & McCaffrey, D. F. (2020). Recommendations about estimating errors-in-variables regression in stata. *The Stata Journal*, 20(1), 116-130. Retrieved from <https://doi.org/10.1177/1536867X20909692>
- McCulloch, C. E., & Searle, S. R. (2001). *Generalized, linear, and mixed models*. New York: John Wiley and Sons.
- McLean, R. A., Sanders, W. L., & Stroup, W. W. (1991). A unified approach to mixed linear models. *The American Statistician*, 45(1), pp. 54-64. Retrieved from <http://www.jstor.org/stable/2685241>
- Mislevy, R. J. (1984). Estimating latent distributions. *Psychometrika*, 49(3), 359-381.
- Nab, L., van Smeden, M., Keogh, R. H., & Groenwold, R. H. (2021). Mecor: An R package for measurement error correction in linear regression models with a continuous outcome. *Computer Methods and Programs in Biomedicine*, 208, 106238.
- National Research Council. (2013). *Frontiers in massive data analysis*. Washington, DC: The National Academies Press. ("<https://www.nap.edu/catalog/18374/frontiers-in-massive-data-analysis>") doi: 10.17226/18374
- Pinheiro, J., & Bates, D. (1994, 09). Approximations to the log-likelihood function in the nonlinear mixed-effects model. *Journal of Computational and Graphical Statistics*, 4. doi: 10.1080/10618600.1995.10474663
- Quarteroni, A., Saleri, F., & Gervasio, P. (2010). *Scientific computing with matlab and octave. 3rd ed.* Springer.
- Rabe-Hesketh, S., Skrondal, A., & Pickles, A. (2002). Reliable estimation of generalized linear

- mixed models using adaptive quadrature. *The Stata Journal*, 2(1), 1-21.
- Rosseel, Y. (2021). Evaluating the observed log-likelihood function in two-level structural equation modeling with missing data: From formulas to R code. *Psych*, 3(2), 197–232. Retrieved from <https://www.mdpi.com/2624-8611/3/2/17>
- SAS Documentation 14.2. (n.d.). Estimating fixed and random effects in the mixed model [Computer software manual]. (Online Help Manual)
- Searle, S. (1982). *Matrix algebra useful for statistics*. New York: John Wiley and Sons.
- Shamir, O., & Zhang, T. (2013, 17–19 Jun). Stochastic gradient descent for non-smooth optimization: Convergence results and optimal averaging schemes. In S. Dasgupta & D. McAllester (Eds.), *Proceedings of the 30th international conference on machine learning* (Vol. 28, pp. 71–79). Atlanta, Georgia, USA: PMLR. Retrieved from <https://proceedings.mlr.press/v28/shamir13.html>
- Sohl-Dickstein, J., Poole, B., & Ganguli, S. (2014). Fast large-scale optimization by unifying stochastic gradient and quasi-newton methods. In *Proceedings of the 31th international conference on machine learning* (Vol. 32). Beijing, China. Retrieved from <http://proceedings.mlr.press/v32/sohl-dicksteinb14.pdf>
- StataCorp. (n.d.). Errors-in-variables regression [Computer software manual]. Retrieved from <https://www.stata.com/manuals/reivreg.pdf> (Online Help Manual)
- Stringer, A. (2021). *Implementing approximate bayesian inference using adaptive quadrature: the aghq package*.
- Teng, S.-H. (2016). Scalable algorithms for data and network analysis. *Foundations and Trends in Theoretical Computer Science*, 12(1-2), 1-274. Retrieved from <http://dx.doi.org/10.1561/04000000051> doi: 10.1561/04000000051
- Tran, D., Toulis, P., & Airolidi, E. (2015, 09). Stochastic gradient descent methods for estimation with large data sets. *Journal of Statistical Software*.
- Tuerlinckx, F., Rijmen, F., Verbeke, G., & De Boeck, P. (2006, 12). Statistical inference in generalized linear mixed models: A review. *The British journal of mathematical and statistical psychology*, 59, 225-55.
- Woodbury, M. A. (1950). Inverting Modified Matrices. In J. Kuntzmann (Ed.), . Princeton, NJ: Princeton University.
- Zhang, F. (1999). *Matrix theory: Basic results and techniques*. New York: Springer.
- Zhou, J., Wei, W., Zhang, R., & Zheng, Z. (2021). Damped newton stochastic gradient descent

method for neural networks training. *Mathematics*, 9(13). Retrieved from <https://www.mdpi.com/2227-7390/9/13/1533>