

Article

An Algorithm for Nonparametric Estimation of A Multivariate Mixing Distribution with Applications to Population Pharmacokinetics

Walter M. Yamada ^{1,†} , Michael N. Neely ^{1,2,†}, Jay Bartroff ^{8,†}, David S. Bayard ^{1,3,†}, James V. Burke ^{4,†}, Mike van Guilder ^{1,†}, Roger W. Jelliffe ^{1,†}, Alona Kryshchenko ^{1,5,†}, Robert Leary ^{6,†}, Tatiana Tatarinova ^{7,†}, and Alan Schumitzky ^{1,8,*}

¹ Laboratory of Applied Pharmacokinetics and Bioinformatics, Children's Hospital of Los Angeles, Los Angeles, CA 90027, USA;

² Pediatric Infectious Diseases, Children's Hospital of Los Angeles, Keck School of Medicine, University of Southern California, Los Angeles, CA 90027, USA;

³ Jet Propulsion Laboratory, California Institute of Technology, Pasadena CA, 91109, USA;

⁴ Department of Mathematics, University of Washington, Seattle, WA 98195, USA;

⁵ Department of Mathematics, California State University Channel Islands, University Dr, Camarillo, CA 93012, USA;

⁶ Certara, Raleigh, NC 27606, USA;

⁷ Department of Biology, University of La Verne, La Verne, CA 91750, USA;

⁸ Department of Mathematics, University of Southern California, Los Angeles, CA 90089, USA;

* Correspondence: schum@usc.edu; Tel.: +1 818-249-9444

† These authors contributed equally to this work.

Abstract:

In this paper we describe a nonparametric maximum likelihood (NPML) method for estimating multivariate mixing distributions. Given N independent observations, convexity theory shows that the NPML estimator is discrete with at most N support points. The original infinite NPML problem then becomes the finite dimensional problem of finding the location and probability of the support points. The probability of the support points is found by a Primal-Dual Interior-Point method; the location of the support points is found by an Adaptive Grid method. Our method is able to handle high-dimensional and complex multivariate mixture models. An important application is discussed for the problem of population pharmacokinetics and a non-trivial example is treated. Our algorithm has been successfully applied in hundreds of published pharmacometric studies. In addition to population pharmacokinetics, this research also applies to empirical Bayes estimation and many other areas of applied mathematics.

Keywords: mixture distribution; mixture model; high dimensional statistics; nonparametric maximum likelihood; primal-dual interior-point method; adaptive grid

1. Introduction

Pharmacometric observations can be described statistically by a mixture model. In this case, the probability of random variable arguments (the PK population model) of the pharmacokinetic compartmental model are described by a mixing distribution. The problem of estimating the mixing distribution from a set of pharmacometric observations can be stated as follows. Let Y_1, \dots, Y_N be a sequence of independent but not necessarily identically distributed random vectors constructed from one or more observations from each of N subjects in the population. Let $\theta_1, \dots, \theta_N$ be a sequence of

independent and identically distributed random vectors belonging to a compact subset Θ of Euclidean space with common but *unknown* distribution F . The $\{\theta_i\}$ are not observed. It is assumed that the conditional densities $p(\mathbf{Y}_i|\theta_i)$ are known, for $i = 1, \dots, N$. The mixing distribution of \mathbf{Y}_i with respect to F is given by $p(\mathbf{Y}_i|F) = \int p(\mathbf{Y}_i|\theta_i)dF(\theta_i)$. Because of independence of the $\{\mathbf{Y}_i\}$, the mixing distribution of the $\{\mathbf{Y}_i\}$ with respect to F is given by

$$L(F) = p(\mathbf{Y}_1, \dots, \mathbf{Y}_N|F) = \prod_{i=1}^N \int p(\mathbf{Y}_i|\theta_i) dF(\theta_i) \quad (1)$$

16 The mixing distribution problem is to maximize the likelihood function $L(F)$ with respect to all probability
17 distributions F on Θ .

18 Remark. The distribution F^{ML} that maximizes $L(F)$ is a *consistent* estimator of the true mixing
19 distribution. This was proved originally by Kiefer and Wolfowitz in 1956 [1]. The consistency of F^{ML}
20 is especially important for our application to population pharmacokinetics where F^{ML} is used as a
21 prior distribution for Bayesian dosage regimen design.

22 The algorithm described in this paper differs from most other published methods in a number of
23 ways. Our algorithm allows for high dimensional Θ . Most published methods require the dimension of
24 Θ to be small and many require the dimension of Θ to be 1, see Section 1.1. We have treated examples
25 where the dimension of Θ is as high as 29, see Section 3.

26 Also most published algorithms require the $\{\mathbf{Y}_i\}$ to be identically distributed and assume that the
27 conditional densities $\{p(\mathbf{Y}_i|\theta_i)\}$ are rather simple, such as $p(\mathbf{Y}_i|\theta_i)$ is a multivariate normal density
28 with mean vector θ_i and covariance matrix Σ . Even if Σ is unknown and has to be estimated, the
29 structure of this model is straightforward. However, the estimation of Σ has to be done carefully to
30 avoid singularities, see Wang and Wang [2]. As will be described in Section 3, we allow $p(\mathbf{Y}_i|\theta_i)$ to be
31 calculated from a system of nonlinear ordinary differential-algebraic equations.

We now describe the details of our algorithm. It was proved by Lindsay [3] and Mallet [4],
under simple hypotheses on the conditional densities $\{p(\mathbf{Y}_i|\theta_i)\}$, that the global maximizer F^{ML} of
 $L(F)$ could be represented by a discrete distribution with at most N support points. This result leads
immediately to a finite dimensional optimization problem for F^{ML} , namely to maximize the likelihood
function

$$L(\lambda, \phi) = \prod_{i=1}^N \sum_{k=1}^K \lambda_k p(\mathbf{Y}_i|\phi_k) \quad (2)$$

32 with respect to the support points $\phi = (\phi_1, \dots, \phi_K)$ and weights $\lambda = (\lambda_1, \dots, \lambda_K)$ such that $\phi_k \in \Theta$, $\lambda_k \geq$
33 0 for $k = 1, \dots, K$, $K \leq N$ and $\sum_{k=1}^K \lambda_k = 1$.

In our algorithm $l(\lambda, \phi) = \log L(\lambda, \phi)$ is maximized, so that

$$l(\lambda, \phi) = \sum_{i=1}^N \log \sum_{k=1}^K \lambda_k p(\mathbf{Y}_i|\phi_k) \quad (3)$$

and the maximization problem becomes

$$\text{maximize } l(\lambda, \phi) \quad (4)$$

34 such that $\phi \in \Theta^K$, $\lambda = (\lambda_1, \dots, \lambda_K) \in \mathbb{R}_+^K$, $K \leq N$ and $\sum_{k=1}^K \lambda_k = 1$.

35 Although the maximization problem in Eq. (4) is finite dimensional, it is still high dimensional.
36 The dimension of the maximization problem in Eq. (4) is $N(\dim \Theta) + (N - 1)$.

37 The optimization problem in Eq. (4) is naturally divided into two problems:

38 Problem 1. Given a set of support points $\{\phi_k\}$, find the optimal weights $\{\lambda_k\}$.

39 Problem 2. Given the solution to Problem 1, find a better set of support points.

40 Problems 1 and 2 are solved cyclically until convergence, i.e. no significant improvement in
41 $l(\lambda, \phi)$.

Problem 1 is a convex programming problem. In our algorithm, we solve this problem by the Primal-Dual Interior-Point (PDIP) method. This type of method is standard in convex optimization theory, see Boyd and Vandenberghe [5]. However, the exact implementation for a specific problem varies from problem to problem. The exact details of our implementation is described in the Appendix. See also Bell [6], Baek [7] and Yamada *et al.* [8]. Our PDIP implementation is very fast and can easily handle thousands of variables.

Finding a better set of support points in Problem 2 is a more difficult problem. This location problem is a non-convex global optimization problem with many local extrema and whose dimension is potentially $N \times \dim \Theta$. The details of our algorithm, called the Adaptive Grid (AG) method, will be described in Section 2.3 and in Algorithm 1.

Roughly speaking, Problems 1 and 2 are solved as follows. An initial large grid of possible support points is defined in Θ . Problem 1 is solved on this large grid. After PDIP, most of the original grid points are removed due to near-zero weights leaving a smaller high-probability grid. Problem 1 is then solved on this smaller grid. Then the Adaptive Grid method for Problem 2 takes place. For each remaining grid point, up to $2 \times \dim \Theta$ new (daughter) support points are added. A daughter point outside the search space Θ or too close to a parent point is discarded. The new grid contains the current high-probability points plus the added daughter points. The algorithm is then ready for Problem 1, again. By construction, each iteration increases the value of $l(\lambda, \phi)$. This process continues until the function $l(\lambda, \phi)$ does not significantly change.

1.1. Other algorithms

Because of space limitations, in this section we only discuss NPML methods that optimize Eq. 4; methods that treat multivariate distributions; and methods which allow general conditional probabilities $\{P(Y_i, \theta_i)\}$. As explained in this paper, any such NPML algorithm has to address two problems: *locations* of support points and *weights* of support points. NPAG does *locations* by an Adaptive Grid method and *weights* by the Primal-Dual Interior-Point (PDIP) method.

The original methods of Lindsay [3] and Mallett [4] were based on algorithms of optimal design in the style of Fedorov [9]. In Schumitzky [10], an algorithm was proposed which did both *locations* and *weights* by the EM algorithm. It was very stable but also very slow.

In Lesperance and Kalbfleisch [11], a new method was introduced which did *weights* by the dual method described in Section 5 of Lindsay [3] and *locations* by what they called the Intra-Simplex Direction Method (ISDM). Even though, the Lesperance and Kalbfleisch paper was restricted to univariate distributions, the ISDM method has been generalized to the multivariate case. To briefly describe ISDM, let $D(\theta, F)$ be the directional derivative of $\log L(F)$ in the direction of the Dirac distribution δ_θ supported at $\theta \in \Theta$. (This function is defined in Section 4 below.) ISDM is an iterative algorithm. At stage k , let F^k be the current estimate F^{ML} . Then find all the local maxima of $D(\theta, F^k)$. These local maxima are added to the current set of support points and a new F^{k+1} is calculated. If there are no new local maxima, then the algorithm is done.

In Pilla, Bartolucci, and Lindsay [12], another new method was developed where the *locations* were found by an initial fine grid. But the *weights* were found by a dual version of the PDIP method.

In Savic, Kjellsson, and Karlsson [13], a nonparametric method was added to the popular NONMEM program. NONMEM-NP is a hybrid parametric-nonparametric approach. The *locations* of support points were found by a parametric maximum likelihood algorithm. Then the *weights* were found by maximizing Eq. (4) relative to the newly found support points. NONMEM-NP can handle high dimensional and complex multivariate distributions. An extension to NONMEM-NP was developed in Savic and Karlsson [14] where additional support points are added to the original set. A comparison between NONMEM-NP and NPAG is discussed in Leary [15].

In Wang and Wang [2], a new algorithm was developed for multivariate distributions. The *locations* were found by a combination of EM and a variant of ISDM. The *weights* were found by a

90 family of Quadratic Programs. In [2], examples are done for 8 and 13 dimensional multivariate mixing
91 distributions.

92 Note: The Quadratic Programming algorithm (QP) of Wang and Wang [2] has a very attractive
93 feature. For a prescribed set of support points, QP finds the zero probabilities exactly. Thus QP avoids
94 the Grid Condensation step where support points from PDIP with sufficiently low probabilities are
95 deleted. However, QP and PDIP are based on different numerical methods and a comparison of the
96 efficiency of both algorithms has not been determined.

97 We finally mention that the NPML problem is a special case of a finite mixture model problem
98 with unknown supports and weights. For a discussion of this approach see Tatarinova and Schumitzky
99 [16].

100 The algorithms which have shown by published examples to handle the highest dimensional
101 multivariate problems are NONMEM NP, Wang and Wang [2], and NPAG.

102 1.2. Benders Decomposition

103 For any set of grid points $\boldsymbol{\phi} = (\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_m)$ in Θ^m , let $\boldsymbol{\lambda} = \hat{\boldsymbol{\lambda}}(\boldsymbol{\phi})$ be the corresponding set of
104 optimal weights given by the PDIP method. Then the function $F(\boldsymbol{\phi}) = l(\hat{\boldsymbol{\lambda}}(\boldsymbol{\phi}), \boldsymbol{\phi})$ depends only
105 on $\boldsymbol{\phi}$ and can be maximized directly. For optimization methods, this technique is called Benders
106 Decomposition. The NPAG algorithm maximizes $F(\boldsymbol{\phi})$ by an adaptive search method. In a method
107 proposed by James Burke, $F(\boldsymbol{\phi})$ is maximized by a Newton type method. Since the function $F(\boldsymbol{\phi})$ is
108 not necessarily differentiable, a relaxed Newton method must be used similar to what is described in
109 the Appendix for the Primal-Dual Algorithm. For details of Benders Decomposition as applied to our
110 problem, see Bell [6], Baek [7] and Jordan-Squire [17].

111 2. Materials and Methods

112 2.1. Pmetrics

113 The simulations and NPAG optimizations in this paper can be duplicated in R, using programs in
114 the Pmetrics package [18]. R and Pmetrics are free software. R is available from many download sites.
115 Pmetrics is available from lapk.org. NPAG is run using the NPrun() command in Pmetrics. Sample
116 datasets and compartmental models are also available at lapk.org.

117 2.2. NPAG Subprograms

118 NPAG is a Fortran program consisting of a number of subroutines as described below. The
119 main program performs the Adaptive Grid (AG) method (consisting of expansion and compression
120 algorithms) and calls the Primal-Dual Interior-Point (PDIP) subprogram. The PDIP algorithm solves
121 the maximization problem of Eq. (4) for a fixed grid and is described precisely in the Appendix.

122 2.3. NPAG Implementation (NPAG - Algorithm 1)

123 For the purpose of this discussion, we can think of PDIP as a function $\hat{\boldsymbol{\lambda}}$ from Θ^m into the set $S^m =$
124 $\{\boldsymbol{\lambda} \in \mathbb{R}_+^m : \sum_{k=1}^m \lambda_k = 1\}$ defined as follows: If $\boldsymbol{\phi} = (\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_m)$ then $\hat{\boldsymbol{\lambda}}(\boldsymbol{\phi}) = (\hat{\lambda}_1, \dots, \hat{\lambda}_m)$ maximizes
125 Eq. (4) relative to the fixed set of grid points $(\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_m)$. In this case we write $G = (\boldsymbol{\phi}, \hat{\boldsymbol{\lambda}}(\boldsymbol{\phi}))$ and $l(G)$
126 $= l(\boldsymbol{\phi}, \hat{\boldsymbol{\lambda}}(\boldsymbol{\phi}))$.

127 In NPAG there are two types of grids: expanded and condensed. The expanded grids are the
128 initial grid and the grids after Grid Expansion (Algorithm 2). The condensed grids are generated by
129 Grid Condensation (Algorithm 3). Each cycle of NPAG begins with an expanded grid. The likelihood
130 calculation is done on the condensed grids.

131 Now for the Adaptive Grid method. Assume that Θ is a bounded Q -dimensional hyper-rectangle.
132 Initially we let $\boldsymbol{\phi}_{expanded}^0 = (\boldsymbol{\phi}_1^0, \dots, \boldsymbol{\phi}_M^0)$ be the set of M Faure grid points in Θ , see [19–21]. Alternatively,

133 we could initially let $\phi_{expanded}^0$ be generated by a uniform distribution on Θ or by a prior run of the
 134 program.

135 Remark. The Faure grid points for a hyper-rectangle Θ are a low-discrepancy set which in some
 136 sense optimally and uniformly covers Θ . In our implementation of NPAG, the Faure point sets come in
 137 discrete sizes which nest with each other. (Allowable number of points equals 2129, 5003, 10007, 20011,
 138 40009, 80021, and multiples of 80021.) This nesting property is useful for checking the optimality of
 139 F^{ML} , see Section 4. We have found that replacing the initial Faure set by a set generated by a uniform
 140 distribution on Θ increases the time to convergence but results in the same optimal distribution.

141 Now set $G_{expanded}^0 = (\phi^0, \hat{\lambda}(\phi^0))$. Our approach is to generate a sequence of solutions G^n to Eq.
 142 (4) of increasingly greater likelihood, where unless otherwise specified, G^n refers to the condensed
 143 grid at the n^{th} cycle of the algorithm. If G^n has log likelihood negligibly different than G^{n-1} , then G^n
 144 is considered the optimal solution to Eq. (4) and is relabeled F^{ML} . If not, then the process continues
 145 using the ϕ^n as the new seed. This loop is repeated until F^{ML} is found.

146 The stopping conditions for NPAG are defined precisely in Algorithm 1. If the stopping conditions
 147 are not met prior to a set maximum number of iterations, the program will exit after writing the last
 148 calculated G^n into a file.

149 2.4. Grid Expansion (EXPAND - Algorithm 2)

150 The crux of the Adaptive Grid method is how to go from G^0 to G^1 or more generally, from G^n to
 151 G^{n+1} . The details of doing this are now explained roughly below and precisely in Algorithm 1.

152 Let Q be the dimension of Θ . Suppose at stage n we have a grid of high-probability support
 153 points ϕ^n . We then add $2Q$ daughter points for each support point $\phi_k \in \phi^n$. The daughter points are
 154 the vertices of a small hyper-rectangle centered at each ϕ_k with size proportional to the original size of
 155 the hyper-rectangle defining Θ . The size of this small hyper rectangle decreases as the accuracy of the
 156 estimates increases. (See Algorithm 2.)

157 Let $\phi_{expanded}^{n+1} = \phi^n \cup \text{Daughter-Points}$. Then the PDIP subprogram is applied to $\phi_{expanded}^{n+1}$ resulting
 158 in the new solution set $G_{expanded}^{n+1} = (\phi_{expanded}^{n+1}, \hat{\lambda}(\phi_{expanded}^{n+1}))$; see Algorithm 1. The solution set $G_{expanded}^{n+1}$
 159 is now ready for grid condensation.

160 2.5. Grid Condensation (CONDENSE - Algorithm 3)

161 The above solution set $G_{expanded}^{n+1}$ may have many support points with very low probability. We
 162 remove all support points which have corresponding probability less than $(\max \lambda) \Delta_\lambda$, where λ is the
 163 vector of current probabilities and the default for Δ_λ is 10^{-3} . (Note that at this point the remaining
 164 probabilities are not normalized.) The probabilities of the remaining support points are normalized
 165 by a second call to the PDIP subprogram. This second call to PDIP is very fast. The likelihood
 166 associated with these remaining support points and normalized probabilities is then used to update the
 167 program control parameters and check for convergence (Algorithm 1 and Section 2.7). If convergence
 168 is attained, then the output of this second call to PDIP provides the support points and probabilities
 169 of the final solution. If convergence is not attained, then the remaining support points are sent to the
 170 Grid Expansion subprogram (Algorithm 2), initializing the next cycle.

171 At the end of the program, the output of this second call to PDIP provides the location and
 172 weights of the final solution.

173 2.6. PDIP Subprogram - See Appendix A

174 The PDIP subprogram finds the optimal solution to Eq. 4 with respect to λ for fixed ϕ . PDIP
 175 employs a primal-dual interior-point method that uses a relaxed Newton method to solve the
 176 corresponding Karush-Kuhn-Tucker equations. (See Eqs. 14-17 of Appendix A.)

177 For any $Y=(Y_1, \dots, Y_N)$ and any $\phi=(\phi_1, \dots, \phi_K) \in \Theta^K$, the input to the PDIP subprogram is the
 178 $N \times K$ matrix $\{p(Y_i|\phi_k)\}$. The output consists of the optimal weights $\hat{\lambda}(\phi)$ and the corresponding

179 log-likelihood $l(\hat{\lambda}(\boldsymbol{\phi}), \boldsymbol{\phi})$. An in-depth description of the PDIP algorithm and its implementation is
 180 presented in Appendix A. See also [6–8].

181 2.7. NPAG Stopping Conditions

182 As explained above, a *potential* solution to F^{ML} is not accepted as a global optimum until successive
 183 sequences of G^n produce final distributions evaluating to sufficiently close log likelihood. The various
 184 upper and lower bounds Δ for NPAG control and stopping conditions are defined below and are used
 185 in Algorithms 1, 2, and 3.

186 Δ_L Primary upper bound on the allowable difference between two successive estimated
 187 Log-Likelihoods; the default initialization is 10^{-4} .

188 Δ_F Secondary upper bound on the allowable difference between two successive estimated
 189 Log-Likelihoods of *potential* F^{ML} ; the default initialization is 10^{-2} .

190 Δ_e Sets an upper bound on the accuracy variable *eps* of Algorithm 1. The default initialization for
 191 Δ_e is 10^{-4} . The default initialization for *eps* is 0.2 and is stepped down until $eps \leq \Delta_e$

192 Δ_F and Δ_e define the two stopping conditions for Algorithm 1.

193 Δ_D Sets a lower bound on how close two support points can get; the default initialization is 10^{-4} .

194 Δ_λ Sets a lower bound factor on the probabilities of the weights λ ; the default initialization is 10^{-3} .

195 2.8. Calculation of $p(\mathbf{Y}_i|\boldsymbol{\phi}_k)$

196 Given observations \mathbf{Y}_i , $i = 1, \dots, N$ and grid points $\boldsymbol{\phi}_k$, $k = 1, \dots, K$, the PDIP subprogram only
 197 depends on the $N \times K$ matrix $\{p(\mathbf{Y}_i|\boldsymbol{\phi}_k)\}$. NPAG can be used for any problem once this matrix is
 198 defined. However, the default setting of NPAG is for the problem of population pharmacokinetics. For
 199 a good background of population pharmacokinetics see Davidian and Giltinan [22,23].

In population pharmacokinetics, generally $\mathbf{Y}_i = (\mathbf{y}_{i,1}, \dots, \mathbf{y}_{i,M})$ is a matrix of vector observations
 for the i -th subject. Since NPAG allows multiple outputs, each $\mathbf{y}_{i,m}$ is itself a q -dimensional vector $\mathbf{y}_{i,m}$
 $= (y_{i,m,1}, \dots, y_{i,m,q})$. The observations $y_{i,m,j}$, are then typically given by a regression equation of the
 form:

$$y_{i,m,j} = f_{i,m,j}(\boldsymbol{\theta}_i) + v_{i,m,j}, \quad j = 1, \dots, q \quad (5)$$

$$v_{i,m,j} \sim N(0, (\sigma_{i,m,j}(\boldsymbol{\theta}_i))^2)$$

$\boldsymbol{\theta}_i$ are unobserved parameters specific for \mathbf{Y}_i

200 In the above Eq. 5, $f_{i,m,j}$ is a known nonlinear function depending on the model structure, the dosage
 201 regimen, the sampling schedule, all covariates and of course the subject-specific parameter vector
 202 $\boldsymbol{\theta}_i$. Except for simple models, $f_{i,m,j}$ requires the solution of (possibly nonlinear) ordinary differential
 203 equations.

In the current implementation of NPAG, it is assumed that the $(\mathbf{y}_{i,1}, \dots, \mathbf{y}_{i,M})$ are independent.
 Then

$$p(\mathbf{Y}_i|\boldsymbol{\phi}_k) = \frac{\exp\left(-\frac{1}{2} \sum_{m=1}^M (\mathbf{y}_{i,m} - \mathbf{f}_{i,m}(\boldsymbol{\phi}_k)) \boldsymbol{\Sigma}_{i,m}^{-1}(\boldsymbol{\phi}_k) (\mathbf{y}_{i,m} - \mathbf{f}_{i,m}(\boldsymbol{\phi}_k))^T\right)}{\prod_{m=1}^M \sqrt{(2\pi)^q \det \boldsymbol{\Sigma}_{i,m}(\boldsymbol{\phi}_k)}} \quad (6)$$

204 where $\mathbf{f}_{i,m} = (f_{i,m,1}, \dots, f_{i,m,q})$ and $\boldsymbol{\Sigma}_{i,m} = \text{diag}(\sigma_{i,m,1}^2, \dots, \sigma_{i,m,q}^2)$. For the purposes of matrix multiplication
 205 in Eq. 6, we think of $\mathbf{y}_{i,m}$ and $\mathbf{f}_{i,m}$ as q -dimensional row vectors.

To complete the description of Eq. 6 we need to model the standard deviation terms $\sigma_{i,m,j}$ of the
 assay noise. In our implementation of NPAG, four different models are allowed. Let

$$\alpha_{i,m,j}(\boldsymbol{\phi}_k) = c_0 + c_1 f_{i,m,j}(\boldsymbol{\phi}_k) + c_2 f_{i,m,j}^2(\boldsymbol{\phi}_k) + c_3 f_{i,m,j}^3(\boldsymbol{\phi}_k) \quad (7)$$

and set

$$\sigma_{i,m,j} = \begin{cases} \alpha_{i,m,j} & \text{assay error polynomial only} \\ \gamma\alpha_{i,m,j} & \text{multiplicative error} \\ \sqrt{\alpha_{i,m,j}^2 + \gamma^2} & \text{additive error} \\ \gamma & \text{constant level of error} \end{cases} \quad (8)$$

206 The parameter γ in Eq. 8 is a variance factor. Artificially increasing the variance during the first
 207 several cycles of NPAG increases the likelihood for each ϕ , allowing the algorithm to use these cycles
 208 to find a better initial state from which to begin optimization. NPAG also has an option to “optimize”
 209 γ . This changes NPAG from a nonparametric method to a “semiparametric” method and will not be
 210 discussed here. The interested reader can consult [8].

Next if $c_0 = 0$ in Eq. 7, then $\alpha_{i,m,j}$ can become very small for certain values of ϕ that in early iterations can be far from optimal. This in turn causes numerical problems as the likelihood is infinite if $\sigma_{i,m,j} = 0$. One way to avoid this problem is to take $\sigma_{i,m,j} = \text{constant}$. Another way would be to assume that $\alpha_{i,m,j}$ is *known* and is given by

$$\alpha_{i,m,j} = c_0 + c_1 y_{i,m,j} + c_2 y_{i,m,j}^2 + c_3 y_{i,m,j}^3 \quad (9)$$

211 That is, to approximate σ by using a polynomial of the observed values rather than model predicted
 212 values. In our experience with NPAG, the approximation of Eq. 9 is useful for ensuring computational
 213 stability (especially during the early cycles of the algorithm). However, from a theoretical perspective,
 214 this change violates the conditions of maximum likelihood and will not be discussed here. Again the
 215 interested reader can consult [8].

216 2.9. Convergence

217 For a given initial grid ϕ^0 , the NPAG algorithm is only guaranteed to find a local maximum of
 218 $L(F)$. More precisely, if ϕ^* is the final grid of NPAG starting from ϕ^0 , then $\hat{\lambda}(\phi^*)$ is a global maximum
 219 on ϕ^* but the support points ϕ^* may be only a local maximum.

220 Global convergence of a nonparametric maximum likelihood method for estimation of a
 221 multivariate mixing distribution is very difficult. For one-dimensional distributions the problem
 222 is straightforward. The idea of proof goes back to at least Fedorov [9] in 1972, which involves the use
 223 of *Directional Derivatives*.

224 Let F be any distribution on Θ . Then the directional derivative of $\log L(F)$ in the direction of the
 225 Dirac distribution δ_θ supported at θ is defined by

226 $D(\theta, F) = [\sum_{i=1}^N P(Y_i|\theta) / P(Y_i|F)] - N$, $\theta \in \Theta$, where $p(Y_i|F) = \int p(Y_i|\theta) dF(\theta)$. Let F_k be the
 227 current NPML estimate at iteration k . The Fedorov method involves maximizing $D(\theta, F_k)$ for $\theta \in \Theta$,
 228 at every iteration. Then the point at which the maximum occurs is added in an optimal way to F_k to
 229 give F_{k+1} . Under the assumptions of regularity, Fedorov shows that $L(F_k)$ converges to $L(F^{ML})$, see
 230 Fedorov [9], (Theorem 2.5.3). Many improvements to this method have been made. In Lesperance and
 231 Kalbfleisch [11] and Wang and Wang [2], instead of just adding the point at which $D(\theta, F_k)$ occurs,
 232 all the points where local maxima occur are added in an optimal way. Again under the assumptions
 233 of regularity, convergence as above is proved. In one-dimension these methods are very efficient. In
 234 higher dimensions, these methods are not computationally practical.

We now suggest a method to check whether the final distribution of NPAG is globally optimal and if not optimal, how close it is to the optimal. It also involves the use of the directional derivative $D(\theta, F)$, but only at the last iteration of NPAG. Now define

$$D(F) = \max_{\theta \in \Theta} D(\theta, F)$$

235 Note that the *max* in the above expression is only over Θ and not over Θ^N . It is proved in Lindsay [3]
 236 that F^* is a global maximum of $L(F)$, i.e. $F^*=F^{ML}$, if and only if $D(F^*) = 0$. Even if $D(F^*) \neq 0$, it is
 237 useful to make this computation as it is also proved in Lindsay [3] that $L(F^{ML}) - L(F^*) \leq D(F^*)$, so
 238 this last expression gives an estimate of the accuracy of the final NPAG result.

239 Now even though we said above it is not practical to calculate $D(F)$ at every iteration of an
 240 algorithm, we are just suggesting to make this calculation at the end of the algorithm. This calculation
 241 can be performed by a deterministic or stochastic optimization algorithm.

242 3. Examples

243 First of all, the NPAG program has been used successfully in high-dimensional and very complex
 244 pharmacokinetic-pharmacodynamic models. In Ramos-Martin *et al.* [24], the NPAG program was
 245 used for a population model of the pharmacodynamics of vancomycin for CoNS infection in neonates.
 246 (Vancomycin is an antibiotic used to treat a number of serious bacterial infections. Coagulase-negative
 247 staphylococci (CoNS) are the most commonly isolated pathogens in the neonatal intensive care unit.)
 248 This model had 7 nonlinear differential equations and 11 random parameters. The population was
 249 a combination of 300 experimental and animal subjects. In Drusano *et al.* [25], the NPAG program
 250 was used for a population model of two drugs for the treatment of tuberculosis. This model had 5
 251 nonlinear differential equations, 3 nonlinear algebraic equations, 1671 observations from 6 outputs
 252 and 29 random parameters. In the algebraic equations, the state variables were only defined implicitly
 253 and had to be solved for by an iterative method.

The above two examples are too complex to use for simulation purposes. Consequently we
 present here a simpler model which has an analytic solution and which can be checked by other
 algorithms. Nevertheless, the estimation of parameters in this model is not trivial. We consider a
 three-compartment PK model with a continuous IV infusion into the central compartment and a bolus
 input into the absorption compartment. The individual subject model is described by the following
 differential equations:

$$\begin{aligned} \frac{dx_1}{dt} &= -K_a x_1, & x_1(t) &= \begin{cases} 0 & \text{for } 0 \leq t < 5 \\ b & \text{if } t = 5 \end{cases} \\ \frac{dx_2}{dt} &= K_a x_1 - (K_{el} + K_{cp}) x_2 + K_{pc} x_3 + r(t), & x_2(0) &= 0 \\ \frac{dx_3}{dt} &= K_{cp} x_2 - K_{pc} x_3, & x_3(0) &= 0 \end{aligned}$$

254 and output equation

$$255 \quad y_1(t) = x_2(t)/V_c + w(t), w(t) \sim N(0, \sigma^2), \sigma = 5.5$$

256 The inputs are a bolus $b = 2000$ at $t = 5$ and a continuous infusion $r(t) = 500$, for $t \geq 0$. This model
 257 has 5 random parameters ($V, K_a, K_{el}, K_{cp}, K_{pc}$). A diagram of this model is given in Figure 1. It is
 258 known that this model is structurally identifiable, see Godfrey [26]. However, we have found that
 259 for a continuous IV infusion, the parameters K_{cp} and K_{pc} are very difficult to estimate in a noisy
 260 environment.

261 The details of the simulation are as follows. There were 300 simulated subjects. The random
 262 variables (V, K_a, K_{cp}, K_{pc}) were independently simulated from normal distributions with means
 263 respectively equal to (1.2, 0.8, 0.2, 2.0) and standard deviations equal to 25% coefficient of variation.

264 The random variable K_{el} was independently simulated from a bimodal mixture of two normal
 265 distributions with means respectively equal to 0.5 and 1.5, with standard deviations equal to 10%
 266 coefficient of variation, and with weights equal to 0.2 and 0.8. This distribution would apply to an
 267 elimination rate constant with a bimodal distribution where 80% of the subjects have a mean of 1.5,
 268 and only 20% have a mean of 0.5. The power of the nonparametric method allows the detection of the
 269 20% group.

Table 1. Simulation versus optimization. Row 1: True simulated means for each parameter. Row 2: NPAG estimates of corresponding means. Row 3: True simulated variances for each parameter. Row 4: NPAG estimated variances for each parameter.

	Kel	Vc	Ka	Kcp	Kpc
μ_{SIM}	1.305	1.194	0.800	0.205	0.408
μ_{NPAG}	1.308	1.189	0.798	0.209	0.410
σ_{SIM}^2	0.170	0.093	0.042	0.002	0.010
σ_{NPAG}^2	0.173	0.086	0.040	0.003	0.011

270 Twelve observations were taken at times

271 $t = 1.1, 5.4, 6.1, 6.5, 6.7, 7.8, 8.4, 9.2, 13.5, 15.3, 15.5, 15.8.$

272 These sampling times were chosen in an ad hoc fashion and are not to be considered optimal. In Figure
273 2 we show the profiles of the 300 noisy model outputs y_1 . These profiles are plotted as piecewise linear
274 functions with nodes at the observation times.

275 The initial Faure set had 80,321 support points. After the first iteration of the NPAG algorithm,
276 the number of support points was down to 300, where it essentially stayed for the rest of the algorithm.
277 After 100 iterations NPAG was stopped based on the convergence criteria of Section 3.5.

278 The simulated and estimated marginal distributions are shown in Figures 3 and 4. It is seen that
279 the estimated marginal distributions were quite accurate. when compared to the simulated histograms.
280 In particular the bimodal shape of K_{el} was uncovered.

281 NPAG is designed to estimate the whole joint distribution of the parameters. As mentioned earlier,
282 the estimate F^{ML} is especially important for our application to population pharmacokinetics where
283 F^{ML} is used as a prior distribution for Bayesian dosage regimen design. However, F^{ML} is a consistent
284 estimator of the true mixing distribution and consequently, the moments of F^{ML} should be consistent
285 estimators of the true moments. Means and variances of parameter estimates for F^{ML} can be easily
286 obtained by integrating the corresponding marginal distributions. So as a check of this fact, in Table 1,
287 the comparisons of estimated versus simulated means and variances are shown. Again, results are
288 quite accurate, see Table 1.

289 Finally, in Figure 5 we include a graph of Predicted versus Observed values which shows the all
290 around good fit of the data. The predicted values are gotten as follows: For each subject, the Bayesian
291 mean estimate of the parameters are found using the final NPAG distribution as a prior and that
292 subject's observations. Then based on these parameter means, the subject's concentration profile is
293 calculated.

294 4. Final Remarks and Conclusions

295 4.1. Final Remarks

296 The NPAG program was developed at the USC Laboratory of Applied Pharmacokinetics. James
297 Burke (University of Washington) developed the Primal-Dual Interior-Point method discussed in the
298 Appendix. Robert Leary (Pharsight Corporation) developed the Adaptive Grid method and wrote the
299 original Fortran program for NPAG. Michael Neely, MD (USC Children's Hospital of Los Angeles)
300 developed the program package Pmetrics which contains NPAG as a subprogram. Pmetrics is an R
301 package for nonparametric and parametric population modeling and simulation and is available at
302 www.lapk.org, see Neely *et al.* [18].

303 4.2. Conclusions

304 We have described a nonparametric maximum likelihood method called NPAG for estimating
 305 multivariate mixing distributions. NPAG is based on an iterative algorithm employing the Primal-Dual
 306 Interior-Point method and an Adaptive Grid method. Our method is able to handle high-dimensional
 307 and complex mixture models. Other methods are discussed. A detailed description of NPAG is given.
 308 The important application to population pharmacokinetics is described and a non-trivial example is
 309 given.

310 In addition to population pharmacokinetics, this research also applies to empirical Bayes
 311 estimation, see Koenker and Mizera [27] and to many other areas of applied mathematics, see Banks
 312 *et al.* [28].

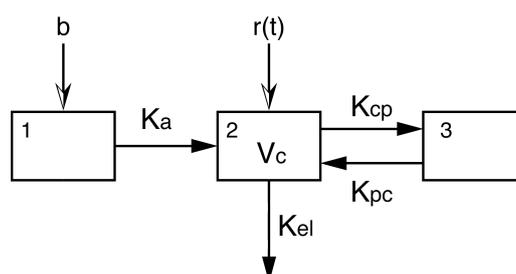


Figure 1. Model.

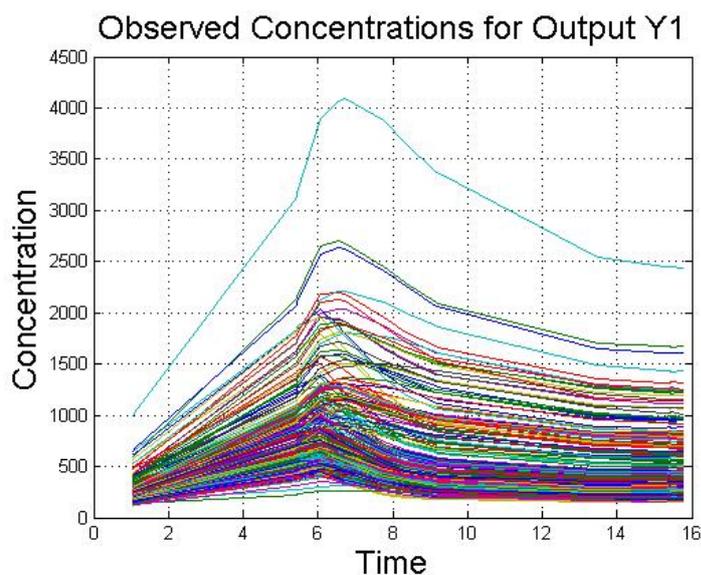


Figure 2. True simulated model profiles.

313 **Author Contributions:** conceptualization, J.Burke, R.L and A.S.; data curation, M.v.G., M.N. and W.Y.; formal
 314 analysis, J.Burke and A.S.; funding acquisition, J.Burke, R.J., M.N. investigation, M.N. and R.J.; methodology, R.L.,
 315 T.T. and A.S., ; project administration, A.S.; resources, M.N. and R.J; software, M.v.G., A.K. and W.Y.; supervision,
 316 A.S.; validation, J.Bartroff, D.B., J.Burke, A.K., R.L., T.T., A.S. and W.Y.; visualization, W.Y.; writing–original draft
 317 preparation, W.Y.; writing–review and editing, J.Bartroff, D.B., J.Burke, A.K., R.L., T.T., A.S. and W.Y.

318 **Funding:**

319 This work was supported in part by grants from NIH: RR11526, GM65619, GM068968, EB005803, EB001978,
 320 HD070886. JB was supported in part by NSF/DMS-0505712.

321 **Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the
 322 study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to
 323 publish the results’.

324 **Abbreviations**

Algorithm 1 NPAG Algorithm. Input: $(Y, \phi^0, a, b, \Delta_D, \Delta_L, \Delta_F, \Delta_e, \Delta_\lambda)$, a and b are the lists of lower and upper bounds, respectively, of Θ ; Δ_D is the minimum distance allowable between points in the estimated F^{ML} . Δ_x see §2.7. Output: $(\phi, \lambda, l(\lambda, \phi))$.

```

1: procedure NPAG( $Y, \phi^0, a, b, \Delta_D$ ) ▷ Estimate  $F^{ML}$  given  $Y$ 
2:   Initialization:  $\phi = \phi^0, \text{LogLike} = -10^{30}, F_0 = 10^{30}, F_1 = 2 * F_0, \text{eps} = 0.2, \Delta_e = 10^{-4},$ 
    $\Delta_F = 10^{-2}, \Delta_L = 10^{-4}, \Delta_\lambda = 10^{-3}, n = 0$ 
3:   while  $\text{eps} \geq \Delta_e$  or  $|F_1 - F_0| \geq \Delta_F$  do
4:     Calculate  $\Psi(\phi)$  ▷  $N \times K$  matrix  $\{p(Y_i|\phi_k)\}$ 
5:      $[\hat{\lambda}(\phi), l(\hat{\lambda}(\phi), \phi)] \leftarrow \text{PDIP}(\Psi(\phi))$  ▷ Appendix A
6:     if (MAXCYCLES == 0) then
7:        $F_{est}^{ML} \leftarrow l(\hat{\lambda}(\phi), \phi)$ 
8:        $\lambda \leftarrow \hat{\lambda}(\phi)$ 
9:       return  $[\phi, \lambda, F_{est}^{ML}]$ 
10:    end if
11:     $n \leftarrow n + 1$ 
12:     $\phi^c \leftarrow \text{CONDENSE}(\phi, \hat{\lambda}(\phi), \Delta_\lambda)$  ▷ Alg. 3
13:     $[\hat{\lambda}(\phi^c), l(\hat{\lambda}(\phi^c), \phi^c)] \leftarrow \text{PDIP}(\Psi(\phi^c))$  ▷ PDIP returns  $G^n$ 
14:     $\text{NewLogLike} = l(\hat{\lambda}(\phi^c), \phi^c)$ 
15:    if ( $n > \text{MAXCYCLES}$ ) then
16:       $F_{est}^{ML} \leftarrow l(\hat{\lambda}(\phi^c), \phi^c)$ 
17:       $\lambda \leftarrow \hat{\lambda}(\phi^c)$ 
18:      return  $[\phi, \lambda, F_{est}^{ML}]$ 
19:    end if
20:    if  $|\text{NewLogLike} - \text{LogLike}| \leq \Delta_L$  and  $\text{eps} > \Delta_e$  then
21:       $\text{eps} = \text{eps}/2$  ▷ Adjust precision
22:    end if
23:    if  $\text{eps} \leq \Delta_e$  then ▷ check EXIT conditions
24:       $F_1 = \text{NewLogLike}$ 
25:      if  $|F_1 - F_0| \leq \Delta_F$  then
26:         $F_{est}^{ML} \leftarrow F_1$ 
27:         $\phi \leftarrow \phi^c; \lambda \leftarrow \hat{\lambda}(\phi^c)$ 
28:        return  $[\phi, \lambda, F_{est}^{ML}]$ 
29:      else
30:         $F_0 = F_1; \text{eps} = 0.2$  ▷ Reset Algorithm
31:      end if
32:    end if
33:     $\phi \leftarrow \phi^c \leftarrow \text{EXPAND}(\phi^c, \text{eps}, a, b, \Delta_D)$  ▷ Alg. 2
34:     $\text{LogLike} \leftarrow \text{NewLogLike}$ 
35:  end while
36: end procedure

```

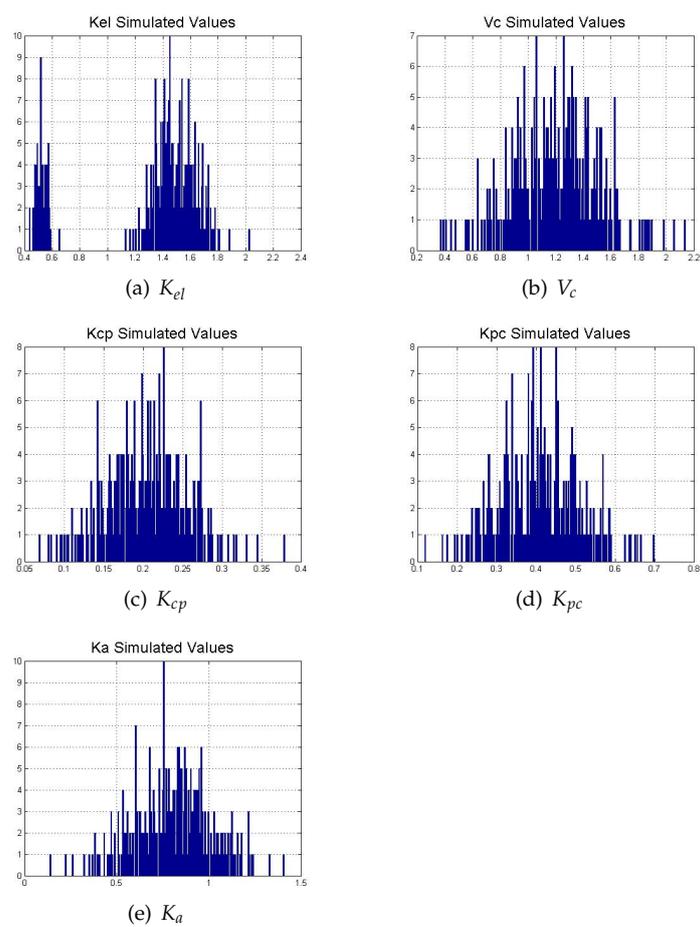


Figure 3. Histogram of simulated PK parameters.

Algorithm 2 EXPAND. Input: $\phi = (\phi_1, \dots, \phi_K)$, Δ_G , $\Theta = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_Q, b_Q]$, $\mathbf{a} = [a_1, \dots, a_Q]$, $\mathbf{b} = [b_1, \dots, b_Q]$, Δ_D . Output: $\phi' = (\phi'_1, \dots, \phi'_M)$, where $M \leq K(1 + 2Q)$. Note: In this algorithm, $\phi = (\phi_1, \dots, \phi_K)$ is a $Q \times K$ matrix, with $Q = \dim \Theta$.

```

function EXPAND( $\phi$ ,  $\Delta_G$ ,  $\mathbf{a}$ ,  $\mathbf{b}$ ,  $\Delta_D$ )
2:   Initialize:  $[Q, K] = \text{size}(\phi)$ ,  $\mathbf{I} = \mathbf{Q} \times \mathbf{Q}$  Identity matrix,  $\text{new}\phi \leftarrow \phi$ 
      for  $k = 1, \dots, K$  do                                      $\triangleright K = \text{number of input support points}$ 
4:     for  $d = 1, \dots, Q$  do                                      $\triangleright Q = \dim \Theta$ 
           $T(d) = \Delta_G(b(d) - a(d))$ 
6:     if  $\phi(d, k) + T(d) \leq b(d)$  then                        $\triangleright$  Check upper boundary
           $\phi^+ = \phi(:, k) + T(d)\mathbf{I}(:, d)$ 
8:      $dist = 10^{30}$ 
          end if
10:    for  $k_{in} = 1 : \text{length}(\text{new}\phi)$  do
           $\text{newdist} = \sum \text{abs}(\phi^+ - \text{new}\phi(:, k_{in})) ./ (\mathbf{b} - \mathbf{a})$     $\triangleright x ./ y$  done component-wise
12:     $dist = \min(dist, \text{newdist})$ 
          end for
14:    if  $dist \geq \Delta_D$  then                                  $\triangleright$  Check distance to new support point
           $\text{new}\phi \leftarrow [\text{new}\phi, \phi^+]$ 
16:    end if
          if  $\phi(d, k) - T(d) \geq a(d)$  then                        $\triangleright$  Check lower boundary
18:     $\phi^- = \phi(:, k) - T(d)\mathbf{I}(:, d)$ 
           $dist = 10^{30}$ 
20:    end if
          for  $k_{in} = 1 : \text{length}(\text{new}\phi(1, :))$  do
22:     $\text{newdist} = \sum (\text{abs}(\phi^- - \text{new}\phi(:, k_{in})) ./ (\mathbf{b} - \mathbf{a}))$     $\triangleright x ./ y$  done component-wise
           $dist = \min(dist, \text{newdist})$ 
24:    end for
          if  $dist \geq \Delta_D$  then                                  $\triangleright$  Check distance to new support point
26:     $\text{new}\phi \leftarrow [\text{new}\phi, \phi^-]$ 
          end if
28:    end for
      end for
30:    $\phi \leftarrow \text{new}\phi$ 
end function

```

Algorithm 3 Condense Algorithm. Input: $(\phi, \lambda, \Delta_\lambda)$, Output: ϕ^c Note: ϕ^c is considered a subset of ϕ

```

function CONDENSE( $\phi$ ,  $\lambda$ ,  $\Delta_\lambda$ )
    $\text{ind} = \text{find}(\lambda > (\max \lambda)\Delta_\lambda)$     $\triangleright$  Inequality and max are performed component-wise
    $\phi^c = \phi(:, \text{ind})$ 
   return  $\phi^c$ 
end function

```

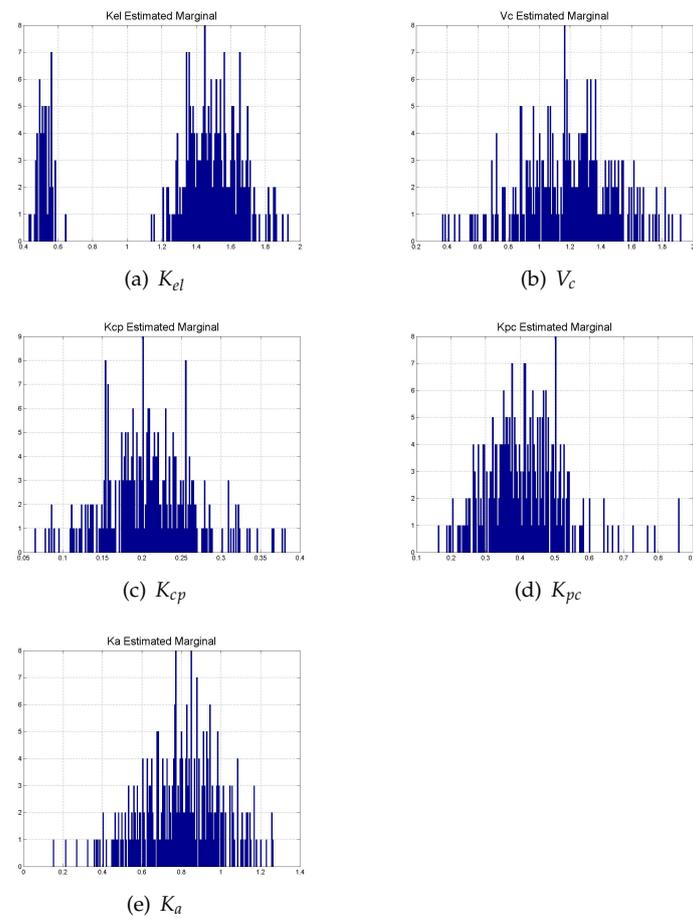


Figure 4. Estimated Marginals of PK parameters.

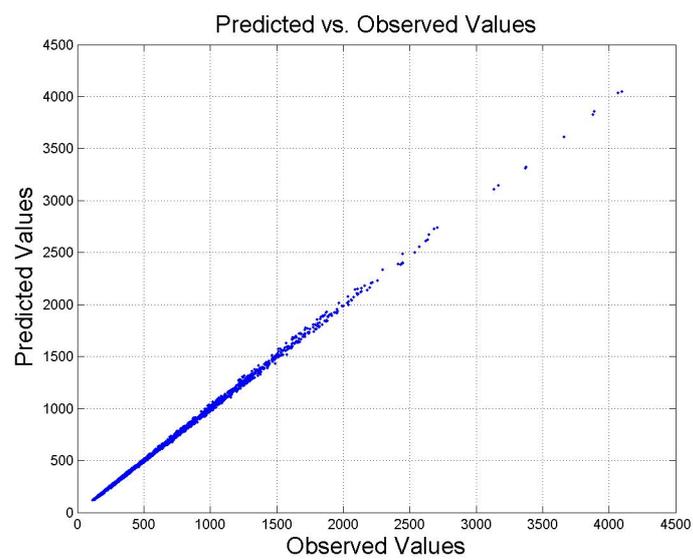


Figure 5. Predicted vs. Observed.

³²⁵ The following abbreviations are used in this manuscript:

AG Adaptive Grid

ISDM Intrasimplex direction method

NPAG Nonparametric adaptive grid algorithm

³²⁶ NPML Nonparametric maximum likelihood

PDIP Primal-dual interior point method

QP Quadratic programming

327 Appendix A. A Primal-Dual Interior-Point Algorithm (PDIP)

328 To make this paper self-contained, we outline here the PDIP algorithm which was written by
 329 James Burke. This algorithm is a FORTRAN subroutine of NPAG. The description below is based on
 330 the Matlab and C++ codes found in Bradley Bell's website, see [6]. Definition of general terms and
 331 theorems can be found in Boyd and Vandenberghe [5].

332 Appendix A.1. Duality Theory and the Basic Problem

Given a set of support points $\{\phi_k\}$, the problem of finding the optimal weights $\{\lambda_k\}$ in Eq. 4 can be posed as the following optimization problem

$$\mathcal{P} \quad \min \Phi(\Psi\lambda) \quad \text{s.t. } 0 \leq \lambda, \quad e^\top \lambda = 1,$$

where $\Psi \in \mathbb{R}^{n \times m}$ is the matrix whose (i, j) entry is $p(y_i | \phi_j)$ and where in general, the function $\Phi : \mathbb{R}^k \mapsto \mathbb{R} \cup \{+\infty\}$ is given by

$$\Phi(z) = \begin{cases} -\sum_{i=1}^k \log z_i & , 0 < z, \text{ and} \\ +\infty & , \text{otherwise.} \end{cases} \quad (\text{A1})$$

333 The symbol e is always to be interpreted as the vector of all ones of the appropriate dimension.

The problem \mathcal{P} is a convex programming problem since the objective function Φ is convex and the constraining region is a convex set. The Fenchel-Rockafellar dual of the convex program \mathcal{P} is the problem

$$\mathcal{D} \quad \min \Phi(\omega) \quad \text{s.t. } L^\top \omega \leq me.$$

334 From Boyd we obtain the following Karush-Kuhn-Tucker (KKT) equations relating the solutions to the
 335 problem \mathcal{P} and \mathcal{D} .

$$me = \Psi^\top w + y \quad (\text{A2})$$

$$e = W\Psi\lambda \quad (\text{A3})$$

$$0 = \Lambda Y e \quad (\text{A4})$$

336 where for any vector x , we define X to be the diagonal matrix having x along the diagonal.

337 Appendix A.2. An Interior-Point Path-Following Algorithm

The relaxed KKT is given by

$$me = \Psi^\top w + y \quad (\text{A5})$$

$$e = W\Psi\lambda \quad (\text{A6})$$

$$\mu e = \Lambda Y e \quad (\text{A7})$$

$$0 \leq \lambda, \quad 0 \leq w, \quad 0 \leq y, \quad (\text{A8})$$

338 for $\mu > 0$. (μ is the relaxation parameter.) A damped Newton's method is used to solve the above
 339 system.

Consider the function $F : \mathbb{R}^{2m+n} \mapsto \mathbb{R}^{2m+n}$ given by

$$F(\lambda, w, y) = \begin{bmatrix} \Psi^\top w + y \\ W\Psi\lambda \\ \Lambda Y e \end{bmatrix}.$$

A triple (λ, w, y) solves Eqs. A.A5 to A.A8 if and only if

$$F(\lambda, w, y) = \begin{pmatrix} me \\ e \\ \mu e \end{pmatrix} \quad (\text{A9})$$

and $0 \leq \lambda, 0 \leq w$, and $0 \leq y$. Path-following algorithms attempt to solve A9 by applying Newton's method for progressively smaller values of the relaxation parameter μ . We first need the derivative of F . It follows

$$F'(\lambda, w, y) = \begin{bmatrix} 0 & \Psi^T & I \\ W\Psi & Z & 0 \\ Y & 0 & \Lambda \end{bmatrix}$$

340 where $z = \Psi\lambda$.

At the k th iteration of the algorithm, the Newton step is given by the solution to the nonsingular linear system

$$F(\lambda^k, w^k, y^k) + F'(\lambda^k, w^k, y^k) * [\Lambda^k, W^k, Y^k]^T = [e_m, e_n, \mu^k e_m]^T \quad (\text{A10})$$

341 where y is constrained to satisfy the first KKT condition $y^k = e_m - \Psi^T w^k$.

The above set of equations can be reduced by standard techniques. It follows:

$$\Delta w = H^{-1} r_2 \quad (\text{A11})$$

$$\Delta y = -\Psi \Delta w \quad (\text{A12})$$

$$\Delta \lambda = r_1 - \lambda - D_1 \Delta y \quad (\text{A13})$$

342 where $H = D_2 - \Psi D_1 \Psi^T$, $D_2 = ZW^{-1}$, $D_1 = \Lambda Y^{-1}$, $r_1 = \mu Y^{-1} e$, $r_2 = W^{-1} e - \Psi r_1$ where the
343 superscript k is suppressed for simplicity.

344 Appendix A.3. The Algorithm

345 To describe the algorithm we need to define the variables:

$$346 \quad q = \frac{1}{m} \sum_{i=1}^m \lambda_i y_i$$

$$347 \quad \rho = \|e - WZe\|_\infty$$

348 and the scaled duality gap

$$349 \quad \gamma = \frac{|\Phi(w) + \Phi(\Psi\lambda)|}{1 + |\Phi(\Psi\lambda)|}.$$

350 (Initialization)

351 Initially choose $\lambda^0 = e_m/m$, $w^0 = e_n/\Psi\lambda^0$, and $y^0 = e_m - \Psi^T w^0$. (Division of two vectors is
352 performed component-wise.) Set $\varepsilon = 10^{-8}$.

353 (Iteration)

354 At iteration $k+1$, set

$$355 \quad \mu^{k+1} = \sigma^k q^k$$

where the reduction factor σ is defined by

$$\sigma = \begin{cases} 1 & , \text{if } \mu \leq \varepsilon \text{ and } \rho > \varepsilon, \\ \min(0.3, (1 - \delta_1)^2), (1 - \delta_2)^2, \frac{|\rho - \mu|}{\rho + 100\varepsilon} & , \text{otherwise.} \end{cases}$$

The next iterates are given by $\lambda^{k+1} = \lambda^k + \delta_1[\Delta\lambda^k]$, $\omega^{k+1} = \omega^k + \delta_2[\Delta\omega^k]$ and $y^{k+1} = y^k + \delta_2[\Delta y^k]$, where the “damping” factors δ_1 and δ_2 are defined by

$$\delta_{1,0} = - \left[\min(\min(\Lambda^{-1}\Delta\lambda), -\frac{1}{2}) \right]^{-1}$$

$$\delta_{2,0} = - \left[\min(\min(Y^{-1}\Delta y), \min(W^{-1}\Delta\omega), -\frac{1}{2}) \right]^{-1}$$

$$\delta_1 = \min(1, 0.99995\delta_{1,0})$$

$$\delta_2 = \min(1, 0.99995\delta_{2,0})$$

356 (Exit Conditions)

357 Iterate Eqs. A.11-A-13 until

358 $\mu \leq \varepsilon$ and $\rho \leq \varepsilon$ and $\gamma \leq \varepsilon$.

359 If these conditions are not satisfied after a set number of iterations, then write “PDIP did not converge
360 in the given number of iterations.”

361 References

- 362 1. Kiefer, J.; Wofowitz, J. Consistency of the Maximum Likelihood Estimator in the Presence of Infinitely
363 many Incidental Parameters. *Ann. Math. Statist.* **1956**, *27*, 887–906.
- 364 2. Wang, X.; Wang, Y. Nonparametric multivariate density estimation using mixtures. *Stat Comput* **2015**,
365 *25*, 33–43.
- 366 3. Lindsay, B.G. The Geometry of Mixture Likelihoods: A general theory. *Ann. Statist.* **1983**, *11*, 86–94.
- 367 4. Mallet, A. A Maximum Likelihood Estimation Method for Random Coefficient Regression Models.
368 *Biometrika* **1986**, *73*, 645–656.
- 369 5. Boyd, S.; Vandenberghe, L. *Convex Optimization*; Cambridge University Press, 2004.
- 370 6. Bell, B. Non-Parametric Population Analysis.
371 http://moby.ihme.washington.edu/bradb主ell/non_par/non_par.xml, 2012.
- 372 7. Baek, Y. An Interior Point Approach to Constrained Nonparametric Mixture Models. PhD dissertation,
373 University of Washington, Department of Mathematics, 2006.
- 374 8. Yamada, W.; Bartroff, J.; Bayard, D.; Burke, J.; van Guilder, M.; Jelliffe, R.; Leary, R.; Neely, M.; Kryshchenko,
375 A.; Schumitzky, A. The Nonparametric Adaptive Grid Algorithm for Population Pharmacokinetic
376 Modeling. Technical Report TR-2014-1, Children’s Hospital Los Angeles, Los Angeles, CA, 2014.
- 377 9. Fedorov, V.V. *Theory of Optimal Experiments*; Academic Press, 1972. edited and translated by W.J. Studden
378 and E.M. Klimko.
- 379 10. Schumitzky, A. Nonparametric EM Algorithms For Estimating Prior Distributions. *Applied Mathematics*
380 *and Computation* **1991**, *45*, 143–157.
- 381 11. Lesperance, M.L.; Kalbfleisch, J.D. An algorithm for computing the nonparametric MLE of a mixing
382 distribution. *J. Am. Stat. Assoc.* **1992**, *87*, 120–126.
- 383 12. Pilla, R.S.; Bartolucci, F.; Lindsay, B.G. Model building for semiparametric mixtures. *arXiv* **2006**.
- 384 13. Savic, R.M.; Kjellsson, M.C.; Karlsson, M.O. Evaluation of the nonparametric estimation method in
385 NONMEM VI. *European Journal of Pharmaceutical Sciences* **2009**, *37*, 27–35.
- 386 14. Savic, R.M.; Karlsson, M.O. Evaluation of an extended grid method for estimation using nonparametric
387 distributions. *AAPS J* **2009**, *11*, 615–627.
- 388 15. Leary, R. An overview of nonparametric estimation methods used in population analysis. Abstracts of the
389 Annual Meeting of the Population Approach Group in Europe; PAGE: Population Analysis Group Europe,
390 , 2017; Number Abstract 7383, p. 26.
- 391 16. Tatarinova, T.; Schumitzky, A. *Nonlinear Mixture Models: A Bayesian Approach*; Imperial College Press,
392 2015.
- 393 17. Jordan-Squire, C. Convex Optimization over Probability Measures. PhD dissertation, University of
394 Washington, Department of Mathematics, 2015.

- 395 18. Neely, M.; van Guilder, M.; Yamada, W.; Schumitzky, A.; Jelliffe, R. Accurate Detection of Outliers
396 and Subpopulations with Pmetrics: a non-parametric and parametric pharmacometric package for R.
397 *Therapeutic Drug Monitoring* **2012**, *34*, 467–476.
- 398 19. Faure, H. Discr pance de suites associ es   un syst me de num ration (en dimension s). *Acta Arithmetica*
399 **1982**, *41*, 337–351.
- 400 20. Bratley, P.; Fox, B.L. Algorithm 659: Implementing Sobol’s Quasirandom Sequence Generator. *ACM*
401 *Transactions on Mathematical Software* **1988**, *14*, 88–100. <http://www.netlib.org/toms/659>.
- 402 21. Fox, B.L. Algorithm 647: Implementation and Relative Efficiency of Quasirandom Sequence Generators.
403 *ACM Transactions on Mathematical Software* **1986**, *12*, 362–376. <http://www.netlib.org/toms/647>.
- 404 22. Davidian, M.; Giltinan, D.M. *Nonlinear Models for Repeated Measurement Data*; Chapman and Hall/CRC
405 Press, 1995.
- 406 23. Davidian, M.; Giltinan, D.M. Nonlinear Models for Repeated Measurement Data: An overview and update.
407 *Journal of Agricultural, Biological, and Environmental Statistics* **2003**, *8*, 387–419.
- 408 24. Ramos-Martin, V.; Johnson, A.; Livermore, J.; McEntee, L.; Goodwin, J.; Whalley, F.; Docobo-Perez, F.;
409 Felton, T.W.; Zhao, W.; Jacqz-Aigrain, E.; Sharland, M.; Turner, M.; Hope, W.W. Pharmacodynamics
410 of vancomycin for CoNS infection: experimental basis for optimal use of vancomycin in neonates. *J*
411 *Antimicrob Chemother* **2016**, *71*, 992–1002.
- 412 25. Drusano, G.; Neely, M.; van Guilder, M.; Schumitzky, A.; Brown, D.; Fikes, S.; Peloquin, C.; Louie, A.
413 Analysis of combination drug therapy to develop
414 regimens with shortened duration treatment for
415 tuberculosis. *PLoS ONE* **2014**, *9*, e101311.
- 416 26. Godfrey, K.R. The identifiability of parametric models used in biomedicine. *Math Model* **1986**, *7*, 1195–1214.
- 417 27. Koenker, R.; Mizera, I. Convex optimization, shape constraints, compound decisions, and empirical Bayes
418 rules. *J. Am. Stat. Assoc.* **2014**, *109*, 674–85. [http://www.econ.uiuc.edu/~sim\\$roger/research/ebayes/
419 brown.pdf](http://www.econ.uiuc.edu/~sim$roger/research/ebayes/brown.pdf).
- 420 28. Banks, H.T.; Kenz, Z.R.; Thompson, W.C. A Review of Selected Techniques in Inverse Problem
421 Nonparametric Probability Distribution Estimation. *Journal of Inverse and Ill-posed Problems* **2012**,
422 *20*, 429–460.