

1 **BayICE: A hierarchical Bayesian deconvolution** 2 **model with stochastic search variable selection**

3 An-Shun Tai¹, George C. Tseng² and Wen-Ping Hsieh^{1*}

- 4
5 1. Institute of Statistics, National Tsing Hua University, Hsinchu, Taiwan, R.O.C.
6 2. Department of Biostatistics, University of Pittsburgh, Pennsylvania, USA

7
8 * Corresponding author

9 Wen-Ping Hsieh

10 Institute of Statistics, National Tsing Hua University, Hsin-Chu, Taiwan

11 101 sec. 2, Kwang-Fu Rd.,

12 Hsinchu City, Taiwan, R.O.C. 30013,

13 TEL: 886-35715131-33188; FAX:886-35728318;

14 E-mail: wphsieh@stat.nthu.edu.tw

15 16 **Abstract**

17 Gene expression deconvolution is a powerful tool for exploring the microenvironment of
18 complex tissues comprised of multiple cell groups using transcriptomic data. Characterizing
19 cell activities for a particular condition has been regarded as a primary mission against diseases.
20 For example, cancer immunology aims to clarify the role of the immune system in the
21 progression and development of cancer through analyzing the immune cell components of
22 tumors. To that end, many deconvolution methods have been proposed for inferring cell
23 subpopulations within tissues. Nevertheless, two problems limit the practicality of current
24 approaches. First, all approaches use external purified data to preselect cell type-specific genes
25 that contribute to deconvolution. However, some types of cells cannot be found in purified
26 profiles and the genes specifically over- or under-expressed in them cannot be identified. This
27 is particularly a problem in cancer studies. Hence, a preselection strategy that is independent
28 from deconvolution is inappropriate. The second problem is that existing approaches do not
29 recover the expression profiles of unknown cells present in bulk tissues, which results in biased
30 estimation of unknown cell proportions. Furthermore, it causes the shift-invariant property of
31 deconvolution to fail, which then affects the estimation performance. To address these two
32 problems, we propose a novel deconvolution approach, BayICE, which employs hierarchical
33 Bayesian modeling with stochastic search variable selection. We develop a comprehensive
34 Markov chain Monte Carlo procedure through Gibbs sampling to estimate cell proportions,
35 gene expression profiles, and signature genes. Simulation and validation studies illustrate that
36 BayICE outperforms existing deconvolution approaches in estimating cell proportions.
37 Subsequently, we demonstrate an application of BayICE in the RNA sequencing of patients
38 with non-small cell lung cancer. The model is implemented in the R package “BayICE” and
39 the algorithm is available for download.

41 1 Introduction

42 Exploring the cellular components of heterogeneous tissues from their gene expression profiles
43 is an essential work for revealing molecular mechanisms across different cell types. For
44 instance, increasing evidence suggests that levels of tumor-infiltrating immune cells are
45 associated with tumor progression, response to therapy, and patient survival (Dieu-Nosjean, et
46 al., 2014; Fridman, et al., 2012; Fridman, et al., 2017). Thus, powerful technologies for single-
47 cell isolation, such as laser microdissection and flow cytometry, have been employed to
48 quantify the numbers of malignant and normal cells in tissue (Hu, et al., 2016). However, these
49 physical approaches to isolating cells of interest at the gene expression level are costly and
50 time-consuming, resulting in drastically reduced biological-content yields. In contrast to
51 single-cell technologies, RNA-seq and microarrays yield bulk gene expression from hundreds
52 of thousands of cells. In heterogeneous tissues, where more than one cell type is present, the
53 expression profile from bulk RNA-seq or microarrays is from cell mixtures; thus, to correctly
54 interpret these data, gene expression deconvolution approaches are required to recover cell
55 type-specific expression and the distinct cellular proportions within complex tissues.

56 In the study of gene expression deconvolution, numerous computational and statistical
57 approaches have been proposed to characterize cell subpopulations within tissues (Anghel, et
58 al., 2015; Becht, et al., 2016; Gong, et al., 2011; Li, et al., 2016; Newman, et al., 2015;
59 Ogundijo and Wang, 2017; Racle, et al., 2017; Xie, et al., 2018; Zhong, et al., 2013).
60 Expression data from a heterogeneous tissue can be modeled as a linear combination of the
61 distinct expression profiles of the cells present in that tissue, weighted by the corresponding
62 cell fractions. These approaches can be grouped into one of three categories depending on
63 whether they use a prior database of cell type-specific expression profiles in the deconvolution
64 procedure: reference-free, reference-based, and semi-reference-based methods. Reference-free
65 approaches aim to directly perform expression deconvolution without cell type-specific
66 references, and their most significant feature is that they estimate the relative cellular
67 proportions and simultaneously disentangle their expression profiles. For instance, many
68 studies have been leveraged on non-negative matrix factorization to decompose mixed gene
69 expression matrices into cell fractions and their corresponding expression profiles (Gaujoux
70 and Seoighe, 2012; Prassas, et al., 2012). Although reference-free models are valuable in the
71 exploration of an uncharacterized cell population, such as tumor subclones (Xie, et al., 2018),
72 relating the cellular components they identify to specific cell types of interest is difficult. Hence,
73 the results of reference-free approaches are unable to clarify the association between a
74 particular cell type and disease progression.

75 By contrast, reference-based methods incorporate external expression profiles of pure
76 cell samples for deconvolution. For example, the analytical tool CIBERSORT successfully
77 borrows cell type-specific information to predict the immune cell components in blood tissues
78 and tumors through v -support vector regression (v – SVR) (Charoentong, et al., 2017; Newman,
79 et al., 2015). A fundamental assumption about reference-based models is that all types of cells
80 present in the target tissues are included in the reference set, and the cellular proportions should
81 sum up to one. Unfortunately, the pure expression profile of malignant cells, a key component
82 in tumors, is a great challenge because of the high genetic heterogeneity of tumors. Hence,
83 reference-based models can only derive the relative cell proportions concerning the reference
84 set rather than the exact proportions concerning the microenvironment. Therefore, the relative
85 cell proportions are not comparable across samples. To overcome this problem, TIMER adopts
86 a series of deconvolution procedures to adjust the relative cell proportions with tumor purity,
87 which is the content of malignant cells in a tumor (Li, et al., 2016).

88 The abovementioned limitation forms the main incentive for developing semi-
89 reference-based deconvolution approaches. In 2017, Racle et al. proposed a framework for
90 estimating the proportion of immune and cancer cells (EPIC) for RNA-seq data (Racle, et al.,
91 2017). EPIC applies least-squares regression with a non-negativity constraint to the
92 deconvolution problem, and requires that the sum of all cell proportions in each tissue must be
93 less than or equal to one. When the sum is not equal to one, one minus the sum of the estimated
94 cell proportions represents the fraction of uncharacterized cells in a tissue that is not accounted
95 for by the reference set; this number is interpreted as the malignant cell proportion in a bulk
96 tumor.

97 Although semi-reference-based models demonstrate the advantages of incorporating
98 cell-specific information and simultaneously extracting the uncharacterized cell types present
99 in tissues, two problems behind these models should be addressed to complete the framework
100 of gene expression deconvolution. First, signature gene selection is critical to the performance
101 of gene expression deconvolution. In some studies, the incorporation of a preselected signature
102 gene set has successfully improved the accuracy of immune cell deconvolution (Chen, et al.,
103 2017; Newman, et al., 2015; Racle, et al., 2017; Wu, et al., 2018). However, the gene activities
104 of a particular cell type usually vary across different tissue microenvironments. Hence, the use
105 of a preselected gene set might cause the loss of data-dependent information and lead to less
106 deconvolution power. The second problem concerns the natural characteristics of
107 deconvolution. In a reasonable strategy for deconvolution, shifts in the mean level of reference
108 samples and tumor samples should not change the estimation of cell proportions. We refer to
109 this as the shift-invariant property for deconvolution. However, the constrained model
110 implemented for EPIC does not maintain the shift-invariant property in deconvolution, and the
111 estimation of cellular components is unstable because the sequencing depth changes across
112 experiments. More specifically, the uncharacterized cell fractions estimated using the
113 constrained least-squares approach are biased toward zero, which will be demonstrated in our
114 results.

115 Therefore, to address the aforementioned problems, we propose a new model based on
116 a hierarchical Bayesian framework for intracellular component exploration. It is called BayICE
117 and is a semi-reference-based approach. Under the Gaussian assumption, we first considered
118 stochastic search variable selection (SSVS) for novel signature gene selection (George and
119 McCulloch, 1993). The SSVS approach has been widely used in transcriptome analyses to
120 select significant genes. For example, Ishwaran and Rao introduced a rescaled Bayesian model
121 for selecting differentially expressed genes through multi-group microarray data (Ishwaran and
122 Rao, 2005). To the best of our knowledge, BayICE is the first attempt to incorporate the
123 mechanism of feature selection for inferring the cellular components of bulk tissues. Moreover,
124 we claim that the BayICE model possesses the shift-invariant property of deconvolution, which
125 yields unbiased estimates of cellular proportions. The model with the shift-invariant property
126 further guarantees that it can recover the expression profiles of uncharacterized cells using
127 posterior mean inference. For the purpose of inference, we applied Gibbs sampling and the
128 Metropolis–Hastings as the sampling procedure in the estimation. In brief, BayICE performs
129 cellular component estimation, uncharacterized cell expression profile estimation, and a novel
130 strategy for signature gene selection.

131 The remainder of this paper is organized as follows. Section 2 introduces the
132 deconvolution in gene expression and states the shift-invariant property. Section 3 introduces
133 the statistical modeling of BayICE for gene expression deconvolution and proposes a Markov
134 chain Monte Carlo (MCMC) algorithm for simulating the posterior distributions of parameters.
135 To assess the model's performance, Section 4 presents simulation studies that investigate gene
136 expression deconvolution, gene selection, and model robustness compared to two existing

137 methods. Section 5 presents applications to two real datasets where underlying true cell
138 proportions are known and performance can be benchmarked. Section 6 describes application
139 of BayICE to 199 non-small cell lung cancer RNA-seq samples and exploration of the cell
140 components present in the microenvironments of lung tumors. Finally, Section 7 provides final
141 discussion and conclusions.

142 2 Deconvolution

143 The deconvolution in gene expression can be formalized as an optimization problem in which
144 the parameter of interest is the cellular proportion $\mathbf{W} = (w_1, \dots, w_K)'$, and the estimates of \mathbf{W}
145 are obtained by

$$(2.1) \quad \hat{\mathbf{W}} = \operatorname{argmin}_{\mathbf{W}} D \left(\begin{pmatrix} y_1 \\ \vdots \\ y_G \end{pmatrix}, \begin{pmatrix} B_{11} & \cdots & B_{1K} \\ \vdots & \ddots & \vdots \\ B_{G1} & \cdots & B_{GK} \end{pmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_K \end{pmatrix} \right),$$

146 where y_g is the gene expression of gene g , B_{gk} is the expression of the k -th cell type in gene g ,
147 and D is the distance metric. In general, D is the Euclidean distance. As mentioned above, the
148 reference-free deconvolution approaches assume the cell type-specific expression (B_{gk}) is
149 unobserved, and the reference-based methods, by contrast, require B_{gk} as input for the
150 optimization problem. Additionally, the semi-reference-based approaches allow that one of the
151 cell types is absent in the reference. We now define the shift-invariant property to characterize
152 different deconvolution models.

153

154 Definition 1 (Shift-invariant property)

155 Let \mathbf{Y} be the mixed expression from bulk tissue, and \mathbf{B} be the cell type-specific expression
156 matrix. If the cellular proportion estimate of a deconvolution method M is invariant when the
157 expression distribution shifts with a constant (location parameter), then M has the **shift-**
158 **invariant property**. That is, M satisfies

$$159 \quad \operatorname{argmin}_{\mathbf{W}} D(\mathbf{Y}, \mathbf{B}\mathbf{W}) = \operatorname{argmin}_{\mathbf{W}} D(\mathbf{Y} + \mathbf{c}, (\mathbf{B} + \mathbf{c})\mathbf{W})$$

160 for any constant \mathbf{c} , where D is the distance metric used by M .

161

162 The shift-invariant property in Definition 1 is essential for evaluating the deconvolution
163 methods, especially in the genomic study. Since the protocol of a gene expression experiment
164 is typically designed for each study, the mean read depth varies across experiments. More
165 specifically, the different experimental protocols cause the location parameters of expression
166 distributions to change. If the location parameter affects the estimation of the same composition,
167 then it is not reasonable to compare results across studies. However, the shift-invariant property
168 guarantees that a deconvolution method with this property can estimate the proportions
169 precisely when the location parameter is changed, and thus, the between-study comparison is
170 valid. In the supplementary file, we have shown that the reference-free and reference-based
171 deconvolution approaches follow the shift-invariant property. By contrast, EPIC adopted the
172 inequality-constrained optimization method to derive cell proportions lacks the shift-invariant
173 property, and hence the following simulation result reveals that the estimates of EPIC are
174 biased.

175 To address the issue of shift-invariant property for the semi-reference-based
176 deconvolution approaches, we proposed a Bayesian deconvolution model which is more robust
177 to the change in the location parameter. The Bayesian hierarchy architecture facilitates the

178 construction of equality-constrained objective function for the semi-reference-based
 179 deconvolution problem via likelihood approach. The proof details in the supplementary
 180 material, and the model construction will be detailed in the next section.

181 3 BayICE Deconvolution Model

182 In this section, we present the proposed hierarchical Bayesian deconvolution model for
 183 intracellular component exploration with novel signature gene selection. Figure 1 provides a
 184 graphical representation of the BayICE hierarchical model. We first describe the input data and
 185 establish the statistical modeling for reference samples and tumor samples. Subsequently, a
 186 stochastic search method, the Bayesian false discovery rate, and an inflation factor are
 187 introduced for signature gene selection. Finally, we adopt the Gibbs sampling approach and
 188 the Metropolis–Hastings approach to develop a comprehensive sampling procedure.

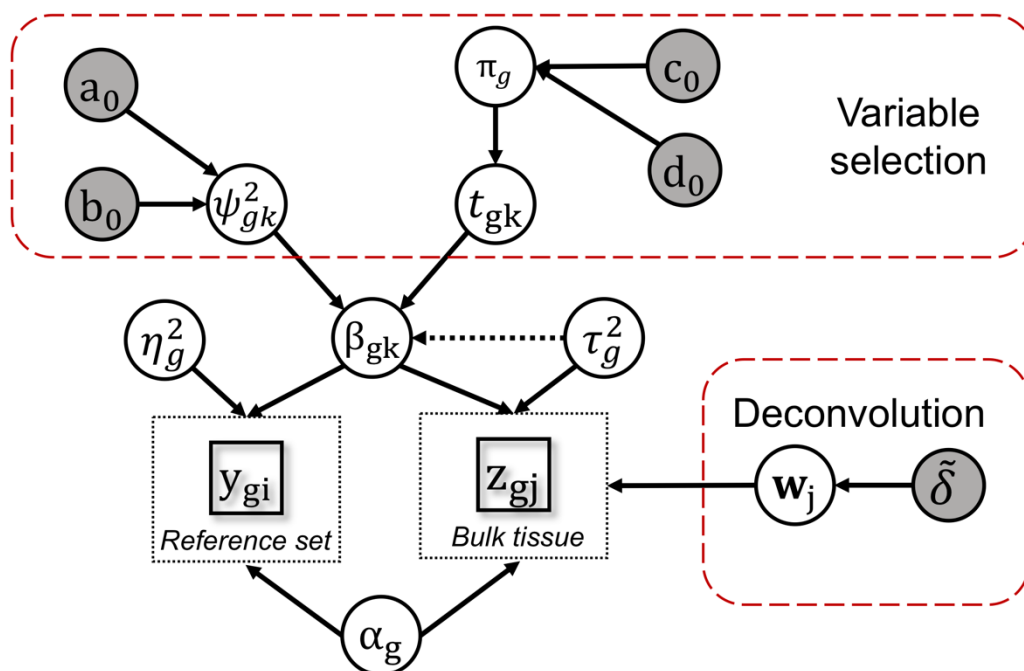


Figure 1. Bayesian hierarchical model of BayICE

The square symbols represent observed data from the reference set and bulk tissues. The white circles indicate priors and the grey circles are hyperparameters.

189 3.1 Input data and normalization

190 BayICE is a statistical framework designed for gene expression deconvolution. We first assume
 191 that the gene expression data are available from two sets of samples, namely heterogeneous
 192 tissues with a given clinical condition and a reference set consisting of several groups of
 193 samples with pure cell types. Although BayICE is a quantitative model-based approach,
 194 microarray and RNA-seq expression data are both legal inputs for BayICE. For the microarray
 195 data, we implement a locally weighted scatterplot smoothing algorithm (Yang, et al., 2002) for
 196 normalization. For the read count RNA-seq data, we recommend two different strategies for
 197 normalization. The first is to use new gene expression units, called transcripts per million
 198 (TPMs), which were proposed by Wagner et al. in 2012. The main feature of TPM
 199 normalization is to make the sum of all TPMs equal across samples, which facilitates fair
 200 comparisons between samples. For public data, TPM normalization approach is not always
 201 available because the information of sequencing depth or gene length required for TPM could
 202 be missing. Thus, we adopt subsampling normalization as the second strategy. Subsampling
 203 normalization applies a binomial sampler to resample the read count, and it aims to maintain

204 internal associations between genes and can simultaneously adjust external variance between
205 samples. The details of how we implement the two normalization approaches are provided in
206 the supplementary materials. After normalization by the first strategy or the second strategy,
207 we consider a log transformation by $\log(\text{count} + 1)$. The log transformation of count data has
208 been widely applied in RNA-seq studies.

209 **3.2 Statistical modeling**

210 The problem of gene expression deconvolution can be formulated as a system of linear
211 equations that describes the expression of a given gene in a bulk tissue as the weighted sum of
212 the expression values from multiple cell types present in the tissue. To maximize the
213 deconvolution power, BayICE incorporates a reference set comprising cell-specific expression
214 profiles into the inference of cellular components in bulk tissues. In this study, the reference
215 set contains various types of immune cells, such as T cells, B cells, natural killer cells,
216 monocytes, neutrophils, and normal tissue cells. In addition to these nonmalignant cells,
217 malignant cells are a major cellular component of tissues in cancer deconvolution studies.
218 Unfortunately, the high genetic heterogeneity of malignant cells hinders the possibility of
219 constructing a predefined cancerous cell expression profile that can be applied to every bulk
220 sample. To address this problem, BayICE takes the advantages of the flexibility of hierarchical
221 Bayesian inference to extract malignant cell profiles directly from bulk tissues. Moreover,
222 BayICE infers cellular proportions with respect to the reference cell types and unknown cell
223 types of each tissue by integrating with a Bayesian signature gene selection approach.

224 We first introduce the statistical model for the reference set consisting of N purified
225 samples with K cell types. We denote an observation in this set as y_{gi} , which is the normalized
226 value for gene- g ($g = 1, \dots, G$) in the i -th purified sample ($i = 1, \dots, N$). The N purified samples
227 belong to K nonmalignant cell types that exist in the bulk tissues. We introduce a binary vector
228 variable $x_i = (x_{i1}, \dots, x_{iK})$ to represent the cell type for the i -th purified sample. More
229 specifically, if the i -th purified sample belongs to cell type k , then $x_{ik} = 1$ and $x_{it} = 0$ for all
230 $t \neq k$. BayICE assumes that y_{gi} follows the Gaussian distribution with a mean level μ_{gi} and a
231 gene-specific variance denoted by $\sigma_{y,g}^2$. For the purpose of signature gene identification, we
232 consider the mean structure μ_{gi} to comprise a gene-specific baseline (α_g) and cell type-specific
233 effects (β_{gk}). Thus, the modeling of reference data can be written as

$$(3.1) \quad (y_{gi} | \alpha_g, \{\beta_{gk}\}, \{x_{ik}\}, \sigma_{y,g}^2) \sim N(\mu_{gi}, \sigma_{y,g}^2),$$

where $\mu_{gi} = \alpha_g + \sum_{k=1}^K x_{ik} \beta_{gk}$.

234 To search for signature genes that exhibit a differential effect among cell types, we define a
235 baseline cell type against which changes in expression levels are measured. We set β_{g1} to be 0
236 for all gene- g .

237 The cell type-specific effects ($\beta_{g1}, \dots, \beta_{gK}$) in (3.1) are shared in constructing the mean
238 structure of gene expression in bulk tissues. It has been observed in the literature that gene
239 expression of a certain cell type changed when it went through cell sorting (Richardson, et al.,
240 2015; van den Brink, et al., 2017). To accommodate the effect of the changes induced by cell
241 separation, we adopt the joint modeling of $\beta_{g1}, \dots, \beta_{gK}$ for pure cell samples and bulk samples.

242 For the model of the bulk tissue, it is a linear combination of K nonmalignant cell-
243 specific expression profiles and malignant cell gene expression. The expression of gene- g from
244 the j -th bulk tissue ($j = 1, \dots, M$) is denoted by z_{gj} . The mean of z_{gj} , called v_{gj} , can be modeled
245 as the weighted sum of nonmalignant cell-specific effects $\{\beta_{gk}\}$ and one malignant cell-
246 specific effect β_{g0} through a linear regression, given by

$$(3.2) \quad (z_{gj} | \alpha_g, \{\beta_{gk}\}, \{w_{jk}\}, \sigma_{z,g}^2) \sim N(v_{gj}, \sigma_{z,g}^2),$$

$$v_{gj} = \alpha_g + \sum_{k=1}^K w_{jk} \beta_{gk} + w_{j0} \beta_{g0}.$$

247 The weights, w_{jk} , are the proportions of expression attributable to normal cell type k in the j -
 248 th tumor, and the weight, w_{j0} , represents tumor purity, which is the percentage of malignant
 249 cells in a tumor tissue. Notably, BayICE can be used to explore not only tumors but also other
 250 noncancerous tissues. For noncancerous tissues, w_{j0} can refer to the proportion of one
 251 unknown cell type that is uncharacterized by the reference set. In our model, a natural constraint
 252 for these cell proportions is that the sum of weights across cell types should be one (i.e.,
 253 $\sum_{k=1}^K w_{jk} + w_{j0} = 1$ for each j). Furthermore, to characterize the gene expression pattern in
 254 malignant cells, we introduce a tumor-specific parameter β_{g0} to represent the effect size of
 255 gene- g in cancer. In a real application, multiple unknown cell types or cancerous cell types
 256 might be present in a bulk tissue. In this case, BayICE treats these uncharacterized cell types
 257 as a whole, and w_{j0} represents the proportion of unknown class in a bulk tissue.

258 3.3 Novel gene selection using the SSVS approach

259 Identifying signature genes that are expressed in a particular cell type is essential to the success
 260 of expression deconvolution. Although a preselected signature gene set could be easily applied
 261 to data analysis, the application of external signature genes could lose data-dependent
 262 information for deconvolution. Thus, we incorporate the stochastic search variable selection
 263 (SSVS) approach into our Bayesian deconvolution model for integrating expression
 264 deconvolution with novel signature gene selection. The SSVS approach, introduced by George
 265 and McCulloch (1993), specifies a spike-and-slab mixture prior, which uses data to extract the
 266 potential features of the true model by inferring posterior probability. The spike component,
 267 which concentrates its mass at values close to zero, shrinks small effects to zero, whereas the
 268 slab component spreads its mass over a wide range of possible values for the effect size.

269 The proposed prior structure of BayICE on effect size exhibits a bimodal distribution
 270 on the variance of the coefficients that result in a spike-and-slab type prior on the effects
 271 themselves (Ishwaran, et al., 2010; Ishwaran and Rao, 2005). For each effect size β_{gk} , the prior
 272 structure is given by

$$(3.3) \quad (\beta_{gk} | \psi_{gk}^2, \pi_k) \sim N(0, \psi_{gk}^2 t_{gk}),$$

$$(t_{gk} | \pi_k) \sim \pi_k I_1(t_{gk}) + (1 - \pi_k) I_c(t_{gk})$$

$$(\psi_{gk}^2 | a_0, b_0) \sim Inv\Gamma(a_0, b_0),$$

$$(\pi_k | c_0, d_0) \sim Beta(c_0, d_0).$$

273 where $I_c(t_{gk})$ is a mass function that is 1 at $t_{gk} = c$ and 0 everywhere else. We set the value c
 274 as a small positive number in this study, such as $c = 10^{-5}$, and thus the random variable t_{gk} is
 275 1 with probability π_k and close to zero with probability $1 - \pi_k$. When the g -th gene is
 276 differentially expressed between the k -th cell type and the other types, β_{gk} is more likely
 277 generated from the slab component and t_{gk} equals one. By contrast, $t_{gk} = c$ indicates that the
 278 g -th gene is irrelevant to cell types and its effect size is from the spike component. The
 279 hypervariance ψ_{gk}^2 is sampled from an inverse gamma with two given hyperparameters, a_0 and
 280 b_0 . Following Ishwaran and Rao (2005), a_0 and b_0 are set as 5 and 50, respectively. The
 281 proportion of genes differentially expressed in cell type k is controlled by π_k , and we assume
 282 that π_k follows a beta distribution with $c_0 = 0.1$ and $d_0 = 0.1$.

283 To borrow variance information across samples, we modify the variance structure of
 284 effect size β_{gk} in (3.3) by considering the gene-specific variance $\sigma_{z,g}^2$ of mixture samples in
 285 (3.2) and the modified prior structure for β_{gk} is given by

$$(3.4) \quad \begin{aligned} (\beta_{gk} | \psi_{gk}^2, t_{gk}, \sigma_{z,g}^2) &\sim N(0, \psi_{gk}^2 t_{gk} \sigma_{z,g}^2), \\ (t_{gk} | \pi_k) &\sim \pi_k I_1(t_{gk}) + (1 - \pi_k) I_c(t_{gk}) \\ (\psi_{gk}^2 | a_0, b_0) &\sim \text{Inv}\Gamma(a_0, b_0), \\ (\pi_k | c_0, d_0) &\sim \text{Beta}(c_0, d_0). \end{aligned}$$

286 The role of the gene-specific variance $\sigma_{z,g}^2$ that appears in (3.4) can be intuitively interpreted
 287 as the baseline in the feature selection procedure. The modified prior structure considers the
 288 trade-off between the value of effect size and gene-specific variance to facilitate the
 289 establishment of feature selection.

290 3.4 Bayesian false discovery rate

291 In frequentist approaches to the test multiplicity problem, controlling the false discovery rate
 292 (FDR) has been widely applied to more adequately control genome-wide false positives.
 293 Whittemore in 2007 introduced a Bayesian FDR associated analogously with the frequentist
 294 FDR (Whittemore, 2007) as follows:

$$(3.5) \quad \phi_k(r) = \frac{\sum_{g=1}^G P_{gk}(H_0|Y, Z) D_{gk}(r)}{\sum_{g=1}^G D_{gk}(r)}$$

295 where $P_{gk}(H_0|Y, Z)$ is the posterior probability that gene- g is not associated with cell type k
 296 (H_0) given observation (Y, Z) , and $D_{gk}(r)$ is the rejection rule defined by $I(P_{gk}(H_0|Y, Z) < r)$.
 297 The tuning parameter r can be adjusted to control the Bayesian FDR at a certain α level. In the
 298 following simulations and applications, the Bayesian FDR is used to address the multiplicity
 299 problem.

300 3.5 Inflation factor

301 In a Bayesian framework, the influence of priors on posterior always vanishes as the
 302 sample size increases. This phenomenon limits Bayesian variable selection because the
 303 mechanism of such selection requires an effective prior setting. To overcome this disadvantage,
 304 Ishwaran and Rao (2005) proposed a rescaling approach to enable estimation invariant to the
 305 sample size by setting the prior as a function of the sample size. They applied data rescaling to
 306 the gene selection framework for multi-group microarray data. Furthermore, to achieve
 307 invariance to sample size, Ishwaran and Rao performed a sample size-related transformation
 308 of gene expression through multiplication by the global inflation factor

$$309 \quad \sqrt{(\text{total sample size}) / (\text{estimate of total variance})}.$$

310 Multiplication by this global inflation factor has been shown to ensure that the prior has a
 311 nonvanishing effect. Hence, following the concept of data transformation, we rescale the
 312 observations in our reference set and bulk tissue set as follows:

$$(3.6) \quad y_{gi}^* = \sqrt{\frac{N+M}{\hat{S}_g^2}} y_{gi} \quad \text{and} \quad z_{gj}^* = \sqrt{\frac{N+M}{\hat{S}_g^2}} z_{gj},$$

313 where $\hat{S}_g^2 = \frac{1}{N-K} \sum_i \sum_k (y_{gi} x_{ik} - \bar{y}_{gk})^2$ is an unbiased estimator of σ_g^2 calculated from the
 314 reference set. Although we assume that the variance σ_g^2 of reference data $\{y_{gi}\}$ is shared with
 315 tissue data $\{z_{gj}\}$, the calculation of an unbiased estimator using both $\{y_{gi}\}$ and $\{z_{gj}\}$ data is a

316 difficult task because of the convolution structure in $\{z_{gj}\}$. Note that the multiplier in (3.6) is
 317 a gene-specific inflation factor rather than the abovementioned global inflation factor since the
 318 inflation factor is composed of the total sample size and a gene-related variance. Therefore, the
 319 use of gene-specific factors can simultaneously achieve sample size invariance and gene-scale
 320 consistency.

321 After rescaling, the corresponding distributions for the transformed data $\{y_{gi}^*\}$ and
 322 $\{z_{gj}^*\}$ are modified as follows:

$$(3.7) \quad \begin{aligned} (y_{gi}^* | \alpha_g, \{\beta_{gk}\}, \{x_{ik}\}, \eta_g^2) &\sim N(\mu_{gi}, (N + M)\eta_g^2), \\ (z_{gj}^* | \alpha_g, \{\beta_{gk}\}, \{w_{jk}\}, \tau_g^2) &\sim N(v_{gj}, (N + M)\tau_g^2), \end{aligned}$$

323 The new variances of the transformed data are adjusted as sample size-related parameters, and
 324 this adjustment can be interpreted as a penalization shrinkage effect of the posterior mean.
 325 After the adjustment with inflation factors, the variances of y_{gi}^* and z_{gj}^* are asymptotically
 326 equal to $N+M$. For the purpose of variable selection, we introduce two further parameters η_g^2
 327 and τ_g^2 for variances in (3.7) to keep the flexibility of modeling.

328 3.6 MCMC sampling procedure

329 Next, based on the transformed data $\{y_{gi}^*\}$ and $\{z_{gj}^*\}$, we complete the structure of the
 330 hierarchical Bayesian model in BayICE and then establish a sampling procedure to achieve
 331 signature gene selection and cell component inference. Following the abovementioned
 332 specification, the BayICE model is given by

$$(3.8) \quad \begin{aligned} (y_{gi}^* | \alpha_g, \{\beta_{gk}\}, \{x_{ik}\}, \eta_g^2) &\sim N(\mu_{gi}, (N + M)\eta_g^2), \quad i = 1, \dots, N, g = 1, \dots, G, \\ \mu_{gi} &= \alpha_g + \sum_{k=1}^K x_{ik}\beta_{gk}, \\ (z_{gj}^* | \alpha_g, \{\beta_{gk}\}, \{w_{jk}\}, \tau_g^2) &\sim N(v_{gj}, (N + M)\tau_g^2), \quad j = 1, \dots, M, g = 1, \dots, G, \\ v_{gj} &= \alpha_g + \sum_{k=1}^K w_{jk}\beta_{gk} + w_{j0}\beta_{g0}, \\ \mathbf{w}_j &= (w_{j1}, \dots, w_{jK}, w_{j0}) \sim \text{Dirichlet}(\delta, \dots, \delta), \quad j = 1, \dots, M, \\ (\alpha_g, \eta_g^2, \tau_g^2) &\sim f(\eta_g^2, \tau_g^2) \propto \frac{1}{\eta_g^2} \times \frac{1}{\tau_g^2} \\ (\beta_{gk} | \psi_{gk}^2, t_{gk}, \tau_g^2) &\sim N(0, \psi_{gk}^2 t_{gk} \tau_g^2), \quad g = 1, \dots, G, k = 1, \dots, K, \\ (t_{gk} | \pi_k) &\sim \pi_k I_1(t_{gk}) + (1 - \pi_k) I_c(t_{gk}) \\ (\psi_{gk}^2 | a_0, b_0) &\sim \text{Inv}\Gamma(a_0, b_0), \\ (\pi_k | c_0, d_0) &\sim \text{Beta}(c_0, d_0), \end{aligned}$$

333 where hyperparameters are specified as $a_0 = 5, b_0 = 50, c_0 = 1, d_0 = 1$, and $\delta = 1$ in this
 334 study.

335 Subsequently, based on the hierarchical prior setting, we apply Gibbs sampling and the
 336 Metropolis–Hastings approach to simulate the posterior value from

$$337 \quad (\boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{T}, \boldsymbol{\pi}, \mathbf{W}, \boldsymbol{\eta}, \boldsymbol{\tau}^2, \boldsymbol{\psi} | \mathbf{Y}^*, \mathbf{Z}^*, \mathbf{X})$$

338 where $\boldsymbol{\alpha} = \{\alpha_g\}$, $\boldsymbol{\beta} = \{\beta_{gk}\}$, $\mathbf{T} = \{t_{gk}\}$, $\boldsymbol{\pi} = \{\pi_k\}$, $\mathbf{W} = \{\mathbf{w}_j\}$, $\boldsymbol{\eta} = \{\eta_g^2\}$, $\boldsymbol{\tau}^2 = \{\tau_g^2\}$, $\boldsymbol{\psi} =$
 339 $\{\psi_{gk}^2\}$, $\mathbf{Y}^* = \{y_{gi}^*\}$, $\mathbf{Z}^* = \{z_{gj}^*\}$, and $\mathbf{X} = \{x_{ik}\}$. The Gibbs sampler in BayICE works are shown
 340 in the supplementary file. After a number of iterations, we specify the burn-in period and the
 341 thinning interval to obtain the posterior distribution of the parameters and then perform an

342 analysis using the posterior mean or posterior median. In BayICE, the default length of burn-
343 in is 0.6 times the number of iterations, and the thinning interval has a length of 3 iterations.
344 The convergence property is discussed in the supplementary materials.

345 **4 Simulation Study**

346 This section presents simulation results based on synthetic datasets to benchmark the
347 performance of BayICE. We consider two types of data: array-based data, which is generated
348 by a normal simulator, and sequencing-based data, which can be produced by a multinomial
349 simulator or negative binomial simulator. We mainly use the multinomial simulator to
350 demonstrate the estimation of cellular proportions, detection of signature genes, and the ability
351 to recover the expression profile of unknown or malignant cells. Furthermore, we apply three
352 different simulators to demonstrate the robustness of BayICE, and the results of the robustness
353 study are shown in the supplementary file.

354 **4.1 Multinomial simulator settings**

355 We first perform a sequencing-based simulation using a multinomial simulator for expression
356 deconvolution. Thus, we consider a scenario in which five distinct cell subpopulations are
357 present in a tissue, and one cell type is absent from our reference set. This simulation includes
358 5000 genes, with 300 genes designed as cell type-related genes. These 300 genes are divided
359 randomly into five disjoint groups, and the genes assigned to a particular group are associated
360 with one cell type. For each cell type in a reference set, we have 20 replicates; hence, the
361 reference set includes 80 samples. In the reference set, an 80×4 matrix of binary variables
362 $X = \{x_{ik}\}$ is used to record the cell type of samples. Furthermore, in the bulk dataset, we
363 simulate the expression data of 90 mixed samples with different cell proportions. The
364 simulation procedure is described as follows.

365 To account for the fact that the expression level varies across genes, we simulated data
366 according to a set of real data from purified samples. We collected 19 RNA-seq samples of
367 normal lung tissues from the Gene Expression Omnibus database with accession number
368 GSE81089 (Mezheyeuski, et al., 2018), and then took the average gene expression across 19
369 samples to obtain the baseline a_g^{sim} of the g -th gene. Among 17,775 genes, we randomly
370 picked 5000 genes for our simulation study. Based on each gene-specific expression level, we
371 define the cell-specific effect as follows:

$$372 \quad b_g^{sim} = (\text{NES})^{c_g^{sim}} (-\text{NES})^{1-c_g^{sim}} a_g^{sim},$$

373 where c_g^{sim} is a binary variable, which is 1 for upregulated status and 0 for downregulated
374 status, and NES is the normalized effect size set as one of the numbers $\{0.1, 0.2, \dots, 0.6\}$. The
375 number $\{c_g^{sim}\}$ is randomly generated with a probability of 0.5 for 0 and 1. We further sample
376 a series of values from Uniform(0.9,1.1), called $\{e_i^{sim}\}$, as sample-specific effects because of
377 the sample heterogeneity. Finally, we generate cellular proportions $\{w_{j1}^{sim}, \dots, w_{jK}^{sim}, w_{j0}^{sim}\}$. To
378 explore the effect of unknown cell content on the model performance, we assign a fixed number
379 to w_{j0}^{sim} , and the remaining components $\{w_{j1}^{sim}, \dots, w_{jK}^{sim}\}$ are generated from the Dirichlet
380 distribution with parameters $(1, \dots, 1)$. Because of the sum-to-one constraint on cellular
381 proportions, $\{w_{j1}^{sim}, \dots, w_{jK}^{sim}, w_{j0}^{sim}\}$ should be normalized by dividing their sum.

382 The mean expression structure is determined for both purified and mixed samples
383 according to the abovementioned parameter settings. Notably, the use of a multinomial model
384 in our simulation is to simulate the sequence alignment procedure that maps the reads against
385 the reference sequence. The number of trials in the multinomial model is related to the total

386 reads, and we set it as 5×10^6 ; in other words, the average read depth is designed as 1000. The
387 probability of the multinomial model controls the expression levels across genes, and therefore,
388 we use the relative mean expression as the probability value. The complete sampling procedure
389 details in the supplementary file.

390 **4.2 Assessing the inference of deconvolution**

391 To assess the performance of BayICE deconvolution, we include two semi-reference-based
392 approaches for comparison: EPIC and non-negative least-squares (NNLS). NNLS is a general
393 approach for solving the constrained least-squares problem where the coefficients are not
394 allowed to become negative. We modify the NNLS approach by restricting the sum of
395 coefficients to less than one for incomplete reference data deconvolution. We first examine the
396 cellular proportion estimations obtained using EPIC, NNLS, and BayICE. Notably, EPIC and
397 NNLS both require an external step to identify signature genes before deconvolution. In this
398 case, we applied the marker genes identified by BayICE into EPIC and NNLS for a fair
399 comparison.

400 According to the simulation setting, we generate 90 bulk samples per simulation in
401 which the unknown cell proportions vary from 0.1 to 0.9. A large unknown cell proportion
402 value indicates that the corresponding tissue is highly heterogeneous, such as in tumors; by
403 contrast, a small proportion simulates the microenvironment of normal tissues in which most
404 of the cell types can be purified. Additionally, we evaluate the performance of each model
405 under different normalized effect sizes of marker genes. Moreover, BayICE can recover the
406 underlying expression profile of the unknown cell type using the posterior mean. Because EPIC
407 and NNLS cannot infer unknown expression profiles, we directly compare the estimation of
408 uncharacterized cell profiles with the true mean expression. The results are shown in Figures
409 2 and 3.

410 We first evaluate the estimation of gene expression profiles for the unknown cell type
411 and gene identification from BayICE in Figure 2. Figure 2(A) is a scatter plot between the true
412 expression of the unknown cell type and the estimated expression. The results reveal the
413 correlation between the real and estimated values to be greater than 0.98, which implies that
414 BayICE can recover uncharacterized cell expression when one cell type is absent from the
415 reference set. In Figure 2(B), we adopt the receiver operating characteristic (ROC) curve to
416 quantify the results of gene identification under different normalized effect sizes which
417 represent the strength of cell type-specific activity in marker genes. The performance in terms
418 of area under the curve (AUC) is significantly improved with the increase of the effect size,
419 and more specifically, the AUC exceeds 0.94 when the effect is larger than 0.2. Subsequently,
420 we apply root-mean-square error (RMSE) to quantify the accuracy of cell proportion estimation
421 from BayICE, EPIC, and NNLS. Figure 3(A) illustrates the changes in RMSE of all estimates
422 with the increase of normalized effect size under the settings of unknown cell proportion = 0.1,
423 0.5, and 0.9. In Figure 3(B), we fix the normalized effect size at 0.1, 0.3, and 0.6 and then
424 evaluate those approaches across different proportions of unknown cells. As a result, it is clear
425 that the performance of all approaches decreases when the unknown cell content increases or
426 the effect size decreases. The abovementioned phenomenon reflects that these gene expression
427 deconvolution approaches are less stable in their inference for tissues with highly
428 uncharacterized content or weak cell type-specific signal. However, the simulation shows that
429 the effects of high unknown cell content and low effect size on BayICE estimation are less
430 severe, and overall, BayICE outperforms the other methods.

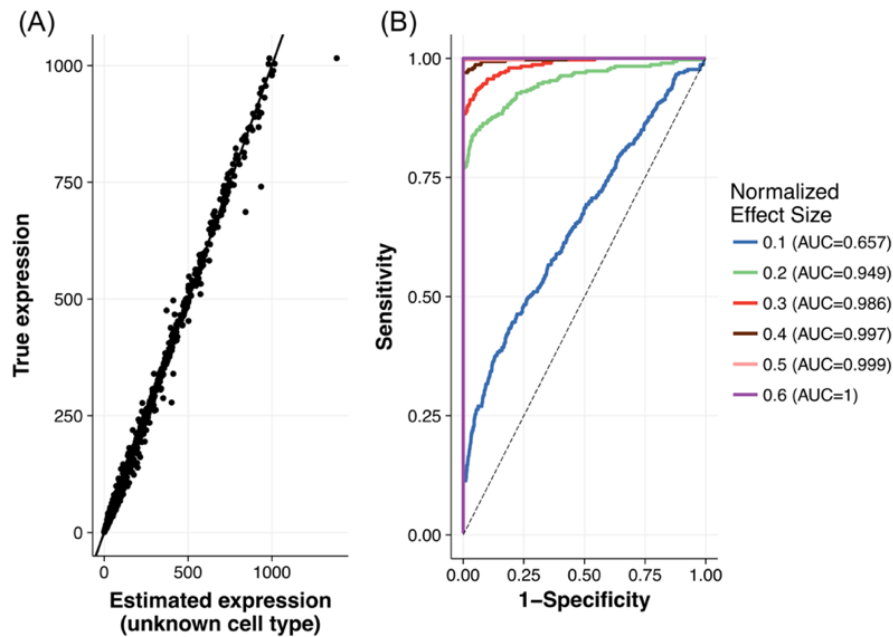


Figure 2. Results of unknown profile estimation and marker gene selection.

(A) Scatter plot of gene expression of the unknown cell type between the truth and estimation. (B) ROC curves to evaluate the gene selection of BayICE under different normalized effect sizes.

431

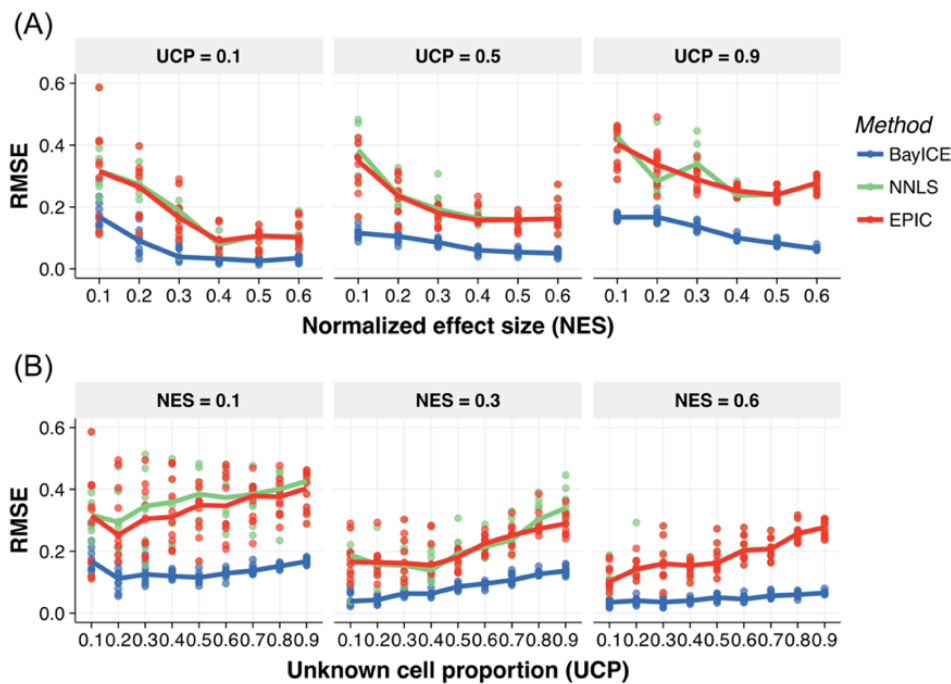


Figure 3. Deconvolution results of cell proportion estimation.

Root-mean-square error (RMSE) between the true and estimated cellular proportions under (A) different effect sizes or (B) different proportions of unknown cells. Ten random sets were generated for each condition, and the RMSE in a set was calculated with five pairs of numbers. The medians of the 10 random sets were connected as the line in the figure.

432 4.3 Evaluating gene detection accuracy

433 The spike-and-slab prior in our model provides a natural consequence of gene selection. In
 434 contrast, the existing models require an external tool to select differentially expressed genes
 435 before deconvolution, and DESeq and edgeR are two popular gene selection tools adopted by
 436 the deconvolution models. To evaluate the accuracy of gene detection, we compare with
 437 DESeq and edgeR. It is worth noted that the genes specifically expressed in the unknown cell
 438 type cannot be identified using DESeq and edgeR, and hence the unknown cell type-related

439 genes are excluded in the evaluation. According to the simulation setting, each gene set is
440 differentially expressed in one cell type. For simplicity, the effect size of marker genes is fixed
441 at 0.2, and the cell proportions are randomly decided across samples. In this simulation study,
442 we also assume some genes express inconsistently between purified cells and the same cell
443 type in bulk samples. For each cell type, we randomly select 100 genes from the gene pool to
444 be inconsistent genes, and disturbed the expression of these inconsistent genes by multiplying
445 an inconsistency level. The values of inconsistency level are 0.1, 0.2, 0.3, 0.4, and 0.5, and the
446 mean of the disturbed gene expression is the original mean in Section 4.1 multiplied by 0.9,
447 0.8, 0.7, 0.6, or 0.5.

448 In this partial comparison, we evaluate the accuracy of identifying genes related to the
449 four cell types present in the reference. AUC is used to assess the gene detection under different
450 levels of inconsistency, and the results are shown in Figure 4. For the low inconsistency level,
451 the AUCs among three methods are close, and it implies that BayICE is comparable to the
452 other approaches designed explicitly for gene selection. In the case of the considerable
453 inconsistency in gene expression, the AUC of BayICE outperforms DESeq2 and edgeR, and
454 it shows that BayICE succeeds in borrowing information from mixed bulk samples for gene
455 detection. Consequently, this partial comparison reveals that BayICE can efficiently recover
456 the findings of the two-step approaches when cell activity causes the difference in gene
457 expression between pure cells and tumors.

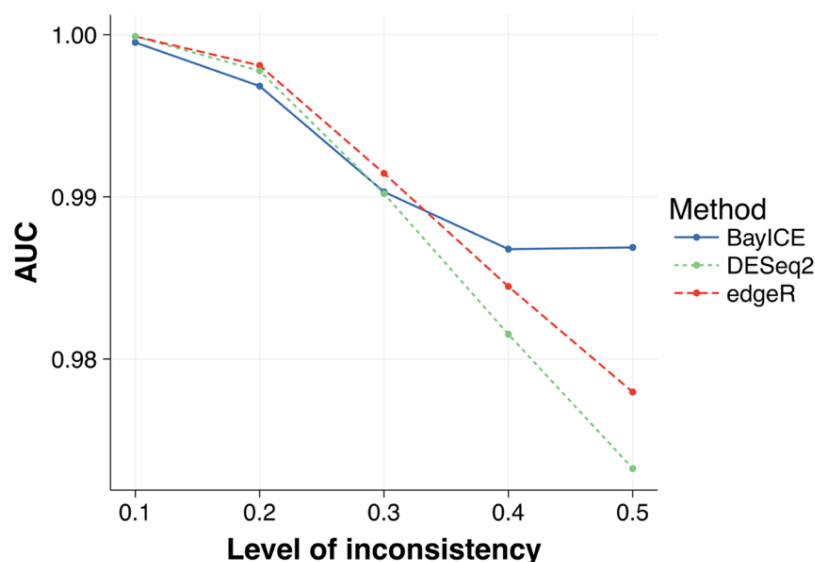


Figure 4. AUC for gene identification.

Comparison of the AUC performance on gene identification. Measuring the performances of BayICE, DESeq2, and edgeR at identifying marker genes as measured by the Area Under the ROC Curve (AUC).

458 **5 Validation in real data with true proportions**

459 In this section, we consider two microarray mixture experiments for validation. They can be
460 downloaded from the Gene Expression Omnibus database using the accession numbers
461 GSE19830 and GSE11058 (Abbas, et al., 2009; Shen-Orr, et al., 2010). The samples from
462 GSE19830 were mixed with rat brain, liver, and lung tissue derived from the same animal in
463 different proportions. The samples in GSE11058 were mixed with four immune cell lines,
464 Jurkat, IM-9, Raji, and THP-1 at various proportions. To validate the performance in the case
465 of the incomplete reference set, we masked one cell type from the reference set and treated the
466 excluded type as the unknown cell component in tissues. Although EPIC is designed for RNA-
467 seq data, the core concept of EPIC modeling is constrained least-squares optimization, and it
468 can be widely applied to various types of data. Hence, we also performed comparisons with
469 EPIC and NNLS on the abovementioned two real datasets.

470 Several microarray studies have confirmed that deconvolution on raw-scale data rather
 471 than on log-scale data can accurately reflect the underlying cellular components. The raw scale
 472 of the data was adopted in this validation study. The signature gene selection for the methods
 473 we compared in the microarray followed Hunt et al. (Hunt, et al., 2018), using a *t*-test between
 474 one cell type and all other cell types for each gene, and we selected the top 200 differentially
 475 expressed genes associated with each cell type.

476 Figures 5 and 6 compare the true and estimated cell proportions using different methods
 477 under the two evaluated datasets. For example, the scatter plot on the first row of Figure 5
 478 represents the deconvolution results when the expression profile of Jurkat cells was unknown
 479 to all the methods. According to the characteristics of the semi-reference-based approach, the
 480 proportions of Jurkat cells in mixed samples can be recovered through estimating unknown
 481 cell proportions. Notably, the constrained models, NNLS and EPIC, tend to assign a very small
 482 proportion to the unknown cell component. This phenomenon was observed in the
 483 supplementary material of robustness and can be attributed to the loss of the shift-invariant
 484 property, which causes the shrinkage of unknown cell proportion estimates.

485 Furthermore, we used the RMSE and Pearson correlation to assess the performance.
 486 RMSE is a measure of accuracy, and the average RMSE of BayICE is 0.093, significantly
 487 lower than those of NNLS (0.170) and EPIC (0.176). Similarly, the measurement of correlation
 488 value is used to monitor the relative order of cellular proportion estimates, and BayICE also
 489 outperformed NNLS and EPIC in terms of the relative size (BayICE = 0.82, NNLS = 0.59, and
 490 EPIC = 0.58). Apparently, BayICE possesses the advantage of incomplete data deconvolution.
 491 Furthermore, the difficulty of deconvolution increases when the abundance of the unknown
 492 cell increases. For instance, in Figure 6, the content of liver cells exceeds 50% of the mixed
 493 samples on average, and the deconvolution approaches perform relatively worse when liver
 494 cells are excluded from the reference set. The aforementioned observation is consistent with
 495 our simulation results.

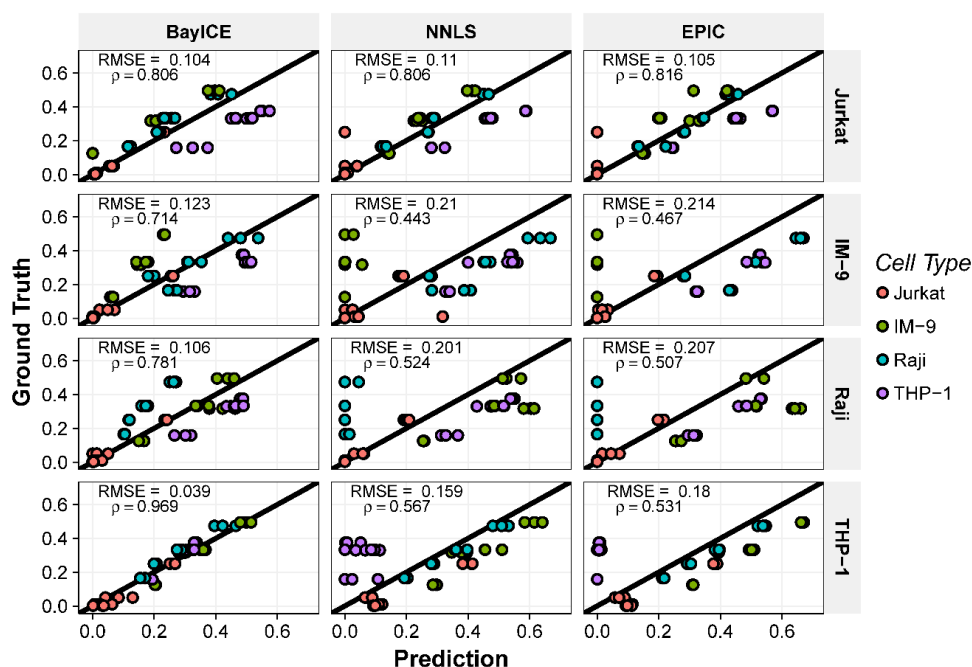


Figure 5. Validation by GSE11058.

Scatter plots comparing estimated cell proportions between the true and estimated proportions from the results of GSE11058. Each column represents a particular method of deconvolution. The row name indicates the cell type that was masked from the reference set and referred to as the unknown cell type. Each of the 12 mixed sample results in four estimates of weights and thus 4 points in each plot. The root-mean-square error and correlation between the ground truth and estimation are also provided in the upper-left corner of each plot.

496

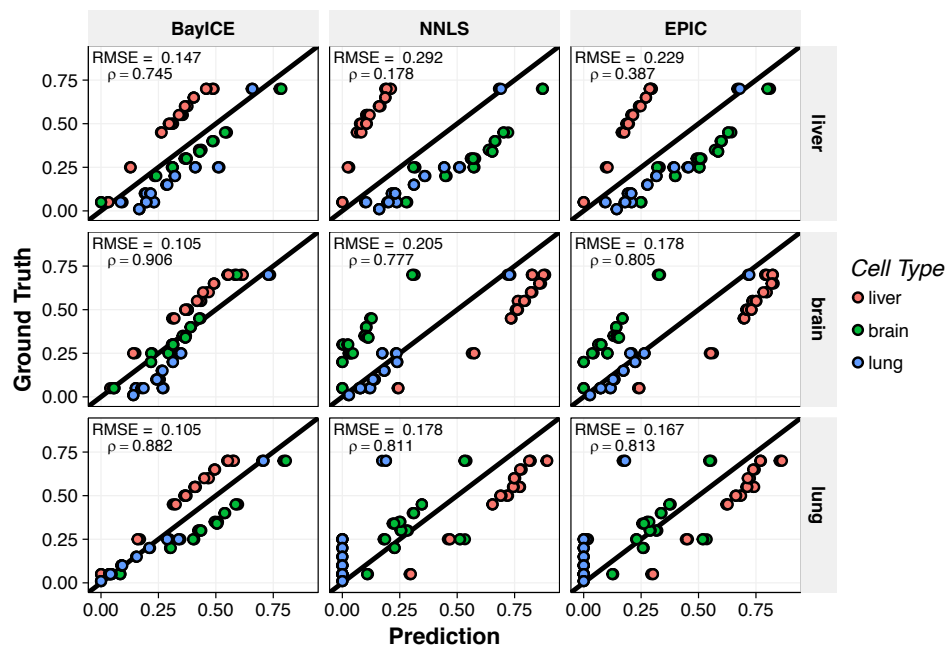


Figure 6. Validation by GSE19830.

Scatter plots comparing the true and estimated proportions from the results of GSE19830. Each column represents a particular method of deconvolution. The row name indicates the cell type that was masked from the reference set and referred to as the unknown cell type. Each of the 9 mixed sample results in three estimates of weights and thus three points in each plot. The root-mean-square error and correlation between the ground truth and estimation are also provided in the upper-left corner of each plot.

497 **6 Application to Non-Small Cell Lung Cancer**

498 To demonstrate an application of BayICE to a biological problem, we consider RNA-seq data
 499 of lung tissues from patients with non-small cell lung cancer (NSCLC). Because tumor-
 500 infiltrating lymphocytes play a critical role in cancer treatment, exploring the changes in
 501 components of immune cells across tumors is of interest. Thus, we consider 199 patients with
 502 NSCLC obtained from GSE81089, and apply BayICE to estimate the cellular components of
 503 the tumor-infiltrating lymphocytes in each tumor sample (Mezheyeuski, et al., 2018). To
 504 construct the reference set, we collect RNA-seq samples of B cells, T cells, granulocytes, and
 505 monocytes of blood tissues from GSE51984 (Pabst, et al., 2016). In addition to immune cell
 506 types, we include normal lung tissues in the reference set, and the expression profiles of 19
 507 normal lung tissues from GSE81089 are used to infer the contents of normal lung cells in
 508 tumors. Because the immune cells are also present in normal lung tissues, we first apply
 509 BayICE to normal lung samples to extract the purified gene expression of lung cells and
 510 estimate the immune cell components in normal samples. Additionally, we randomly divided
 511 19 normal samples into two sets: ten samples are used for extracting purified expression of
 512 lung cells in Section 6.1, and nine samples are analyzed along with tumors for validation in
 513 Section 6.2. As a result, we used the complete reference set consisting of immune cell profiles
 514 and purified lung cell profiles to recover the cellular proportions of each NSCLC sample.

515 **6.1 Deconvolution of normal lung tissues**

516 Figure 7 illustrates the estimated cellular proportions of ten normal lung samples based on the
 517 B cells, T cells, granulocytes, and monocytes. The results of normal lung tissue deconvolution
 518 reveal that monocytes are more prevalent than the other immune cells in lung tissues.
 519 Monocytes typically circulate through the blood for 1–3 days before migrating into tissues,
 520 where they become macrophages or dendritic cells. In lungs, monocytes migrate from the

521 bloodstream into the pulmonary alveolus and are specifically called alveolar macrophages,
522 which play a critical role in homeostasis, host defense, response to foreign substances, and
523 tissue remodeling (Kopf, et al., 2015).

524 In addition to the immune cell types collected in the reference set, BayICE could extract
525 the component of one uncharacterized cell type present in tissues. In this case, the unknown
526 cell type was presumably dominated by normal lung cells, and we calculated the mean
527 expression profile of the unknown cell type using the posterior mean. The next step was to
528 integrate the immune cell profiles with the estimated profile of normal lung cells to construct
529 a new reference set for tumor sample deconvolution.

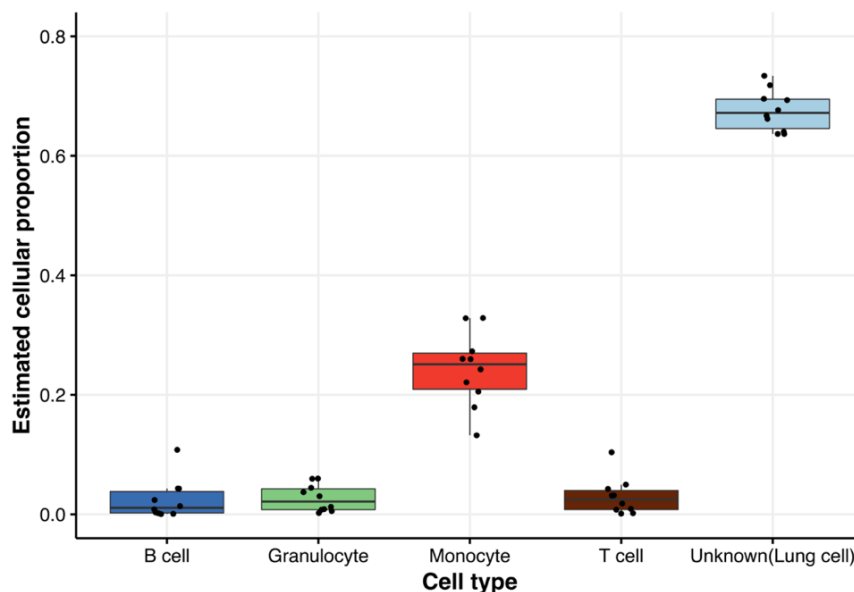


Figure 7. Cell proportion estimation of normal lung tissues.

Plot showing estimated proportions of different cell types from 10 normal lung tissues. The rightmost box represents the estimated unknown component from BayICE, and we refer to this component as normal lung cells.

530 6.2 Deconvolution of tumors

531 In this step, we investigated the 199 NSCLC samples and nine healthy samples according to
532 the new reference set. We applied BayICE to these 208 samples and estimated the fractions of
533 normal lung cells, B cells, T cells, granulocytes, and monocytes, as well as malignant cells.
534 Following the mechanism of normal tissue deconvolution, the malignant cell fractions were
535 defined as the unknown cell proportions in tumor deconvolution. To more effectively
536 understand the change in cellular components during tumor progression, we considered the
537 classification of a malignant tumor (TNM) staging system, which is a standard for classifying
538 the extent of cancer spread. Figure 8 is a boxplot of estimated cell proportions under different
539 TNM stages.

540 Two crucial observations were made from the deconvolution results. First, the
541 estimated cell proportions were more consistent between healthy samples than between
542 patients with NSCLC. For example, the interquartile ranges of estimated lung cell proportions
543 in normal tissues and tumors were 0.047 and 0.137, respectively. A three-fold difference of
544 dispersion between healthy samples and patients revealed that the homeostatic balance between
545 cells in the lung is disturbed when tissues are cancerous. The high fluctuation of cellular
546 components between patients with NSCLC can directly explain the inter-tumor heterogeneity,
547 which leads to the different attributes of different tumors despite the same diagnoses. Second,
548 our study showed that the number of immune cells in the tissue microenvironment increases
549 from normal stage to cancerous stage, which is in agreement with past studies (Banat, et al.,

550 2015; Li, et al., 2016; Ruffini, et al., 2009; Seo, et al., 2018). In particular, (Banat, et al., 2015)
551 comprehensively assessed the number of immune cells in lung cancer by directly counting cells
552 with cell-specific biomarkers and observed an increased number of immune cells in lung cancer
553 tissues compared with healthy donor lungs. In 2018, Seo et al. (2018) applied two approaches,
554 ESTIMATE and TIMER, to infer the cellular components in NSCLC (Li, et al., 2016; Seo, et
555 al., 2018; Yoshihara, et al., 2013). Their results showed a high abundance of dendritic cells
556 (derived from monocytes), which coincides with our observation of estimated monocyte
557 proportions.

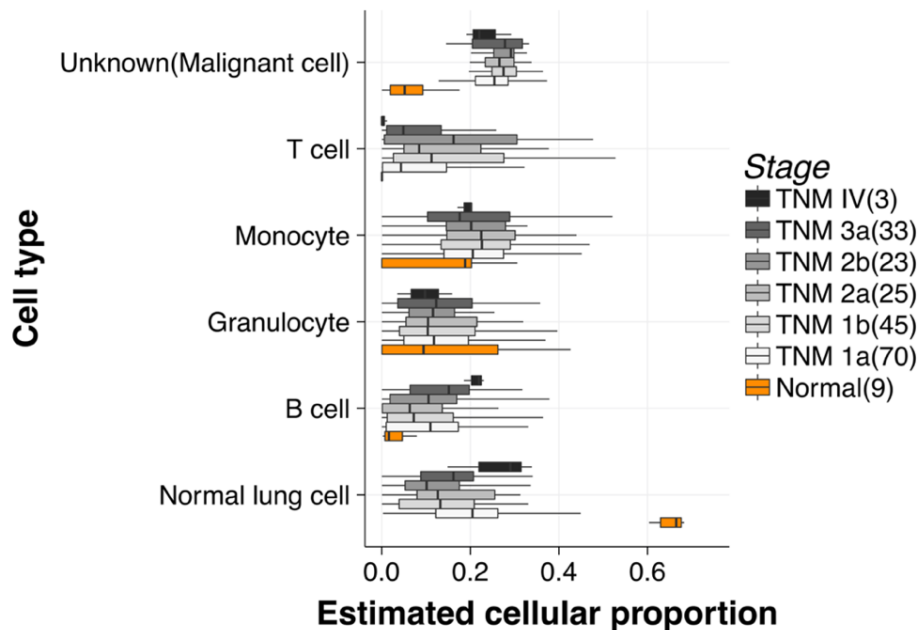


Figure 8. Box plot of cell proportion estimation.

Orange boxes refer to the results of the nine normal lung samples, and other boxes are results from 199 patients with non-small cell lung cancer. The tumor samples can be grouped into six stages and the sample number included for each stage is indicated in the parentheses. For samples of the same tumor stage, the estimated proportions of any specific cell type are summarized as one boxplot and the boxplots for the same cell type are plotted side-by-side.

558 7 Discussion and Conclusion

559 In this study, we developed a novel deconvolution model, BayICE, to infer the cellular
560 components of bulk tissues. BayICE is a semi-reference-based approach that aims to explore
561 cell populations characterized by an external reference set of purified samples and
562 simultaneously investigate the content of uncharacterized cells present in bulk tissues.
563 However, in contrast to constrained models, BayICE takes advantage of a hierarchical
564 Bayesian framework to not only estimate the unknown cell proportion but also recover its gene
565 expression profile. Furthermore, BayICE maintains the shift-invariant property, which
566 guarantees the accuracy of cell proportion estimation.

567 Other than the above-mentioned merits, there are two major contributions of this study.
568 First, BayICE integrates gene expression deconvolution and gene selection in the same model.
569 Most of the current deconvolution approaches require pre-analysis of an additional dataset to
570 identify signature genes for deconvolution, and the external gene selection might not be
571 consistent to the target dataset. Thus, BayICE incorporates SSVS, a Bayesian variable selection
572 approach, to implement internal gene selection. Second, BayICE adopts shared parameters
573 between the pure cells and tumor samples for cell-specific effect. It has been studied that the
574 cell-to-environment interaction causes some of the genes to be expressed inconsistently after
575 cell sorting. We have shown that the joint modeling of both pure cells and tumors in BayICE
576 is more resistant to the problem of inconsistent genes.

577 To evaluate the model's performance, we first conducted an analysis under several
578 simulation scenarios to investigate cell proportion estimation, gene identification, and the
579 robustness of the model. For proportion estimation, we compared two other semi-reference-
580 based approaches, EPIC and NNLS, with BayICE under different unknown cell contents. The
581 results revealed that BayICE significantly outperformed the other methods. We further
582 provided a partial comparison of gene identification with the well-known gene selection
583 approaches, DESeq and edgeR. We found that BayICE can decompose bulk data extremely
584 well and identify cell type-related genes. Moreover, a simulation study using three different
585 simulators revealed that BayICE is more robust with respect to data types. In addition to
586 simulation data, we further applied two real datasets with underlying true cell proportions for
587 validation. The validation of incomplete data deconvolution was consistent with our simulation
588 results, in which BayICE exhibited high accuracy in cell proportion estimation.

589 Our real data application presented an example using the data of 199 patients with
590 NSCLC. We first applied BayICE to healthy lung samples to investigate the cellular
591 components under a normal condition and extract the relatively purified expression profile of
592 lung cells. The deconvolution of normal tissues succeeded in capturing the primary component
593 of immune cells. Next, we formed a new reference set consisting of the immune cell profiles
594 and estimated normal lung cell profiles to analyze the patients with NSCLC. From the analysis,
595 we observed inter-tumor heterogeneity in NSCLC according to the high variation in cell
596 proportions across tumors. In addition, when comparing immune cell proportions between
597 normal and NSCLC samples, the increased immune cell content in tumors revealed that the
598 immune system was highly activated in the cancerous microenvironment. The inference of
599 NSCLC deconvolution coincides with not only the analytic deconvolution results from other
600 studies but also the observations from an immunohistochemical experiment.

601 BayICE has thoroughly addressed technical problems of deconvolution to investigate
602 cellular components, but one issue relating to cell activity remains. In real applications, the
603 activities of cell-to-cell communication and cell-to-environment interaction cause some of the
604 genes to be expressed inconsistently between reference and bulk tissue samples. This
605 phenomenon increases the difficulty of selecting marker genes, and an inappropriate gene set
606 limits the ability to explore tissue environments. Although the joint modeling technique of
607 BayICE can adjust the biased gene expression induced by cell sorting, integrating the biological
608 information of cell activity with deconvolution is believed to be more efficient in estimating
609 cell proportions of bulk tissues. Single-cell RNA-seq has emerged as a powerful new set of
610 technologies for characterizing cell interaction, and the primary goal of our future studies is
611 incorporating single-cell RNA data to modify the prior structure with cell interaction. This
612 extension will generate new insights into the deconvolution framework.

613 Moreover, we will also investigate the unknown cell proportions estimated by BayICE.
614 The proportion of the unknown cell class can be further decomposed if the unknown class is
615 comprised of multiple cancerous cell types. Instead of estimating immune cell proportions, the
616 study of tumor clonal evolution focuses on exploring the contents of different tumor cell types
617 to understand tumor progression. Hence, after BayICE dissects the part belonged to tumor, we
618 can further decompose tumor cells to evaluate the size of tumor clones.

619 All of the work in this study was implemented on R, and the R package, BayICE, is
620 publicly available for deconvolution analysis (<https://github.com/AshTai/BayICE>).

621 **Acknowledgments**

622 This work was supported by the Ministry of Science and Technology [MOST xxxx].

623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674

References

- ABBAS, A.R., WOLSLEGEL, K., SESHASAYEE, D., MODRUSAN, Z. and CLARK, H.F. (2009). Deconvolution of blood microarray data identifies cellular activation patterns in systemic lupus erythematosus. *PloS one*; 4(7):e6098.
- ANGHEL, C.V., QUON, G., HAIDER, S., NGUYEN, F., DESHWAR, A.G., MORRIS, Q.D. and BOUTROS, P.C. (2015). ISOpureR: an R implementation of a computational purification algorithm of mixed tumour profiles. *BMC bioinformatics*; 16(1):156.
- BANAT, G.-A., TRETYN, A., PULLAMSETTI, S.S., WILHELM, J., WEIGERT, A., OLESCH, C., EBEL, K., STIEWE, T., GRIMMINGER, F. and SEEGER, W. (2015). Immune and inflammatory cell composition of human lung cancer stroma. *PLoS One*; 10(9):e0139073.
- BECHT, E., GIRALDO, N.A., LACROIX, L., BUTTARD, B., ELAROUCI, N., PETITPREZ, F., SELVES, J., LAURENT-PUIG, P., SAUTÈS-FRIDMAN, C. and FRIDMAN, W.H. (2016). Estimating the population abundance of tissue-infiltrating immune and stromal cell populations using gene expression. *Genome biology*; 17(1):218.
- CHAROENTONG, P., FINOTELLO, F., ANGELOVA, M., MAYER, C., EFREMOVA, M., RIEDER, D., HACKL, H. and TRAJANOSKI, Z. (2017). Pan-cancer immunogenomic analyses reveal genotype-immunophenotype relationships and predictors of response to checkpoint blockade. *Cell reports*; 18(1):248-262.
- CHEN, Z., HUANG, A., SUN, J., JIANG, T., QIN, F.X.-F. and WU, A. (2017). Inference of immune cell composition on the expression profiles of mouse tissue. *Scientific reports*; 7:40508.
- DIEU-NOSJEAN, M.-C., GOC, J., GIRALDO, N.A., SAUTÈS-FRIDMAN, C. and FRIDMAN, W.H. (2014). Tertiary lymphoid structures in cancer and beyond. *Trends in immunology*; 35(11):571-580.
- FRIDMAN, W.H., PAGÈS, F., SAUTÈS-FRIDMAN, C. and GALON, J. (2012). The immune contexture in human tumours: impact on clinical outcome. *Nature Reviews Cancer*; 12(4):nrc3245.
- FRIDMAN, W.H., ZITVOGEL, L., SAUTÈS-FRIDMAN, C. and KROEMER, G. (2017). The immune contexture in cancer prognosis and treatment. *Nature reviews Clinical oncology*; 14(12):717.
- GAUJOUX, R. and SEOIGHE, C. (2012). Semi-supervised Nonnegative Matrix Factorization for gene expression deconvolution: a case study. *Infection, Genetics and Evolution*; 12(5):913-921.
- GEORGE, E.I. and MCCULLOCH, R.E. (1993). Variable selection via Gibbs sampling. *Journal of the American Statistical Association*; 88(423):881-889.
- GONG, T., HARTMANN, N., KOHANE, I.S., BRINKMANN, V., STAEDTLER, F., LETZKUS, M., BONGIOVANNI, S. and SZUSTAKOWSKI, J.D. (2011). Optimal deconvolution of transcriptional profiling data using quadratic programming with application to complex clinical blood samples. *PloS one*; 6(11):e27156.
- HU, P., ZHANG, W., XIN, H. and DENG, G. (2016). Single cell isolation and analysis. *Frontiers in cell and developmental biology*; 4:116.
- HUNT, G.J., FREYTAG, S., BAHLO, M. and GAGNON-BARTSCH, J.A. (2018). dtangle: accurate and robust cell type deconvolution. *Bioinformatics*; 35(12):2093-2099.
- ISHWARAN, H., KOGALUR, U.B. and RAO, J.S. (2010). spikeslab: Prediction and Variable Selection Using Spike and Slab Regression. *R Journal*; 2(2).
- ISHWARAN, H. and RAO, J.S. (2005). Spike and slab gene selection for multigroup microarray data. *Journal of the American Statistical Association*; 100(471):764-780.
- ISHWARAN, H. and RAO, J.S. (2005). Spike and slab variable selection: frequentist and Bayesian strategies. *The Annals of Statistics*; 33(2):730-773.
- KOPF, M., SCHNEIDER, C. and NOBS, S.P. (2015). The development and function of lung-resident macrophages and dendritic cells. *Nature immunology*; 16(1):36.
- LI, B., SEVERSON, E., PIGNON, J.-C., ZHAO, H., LI, T., NOVAK, J., JIANG, P., SHEN, H., ASTER, J.C. and RODIG, S. (2016). Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome biology*; 17(1):174.

- 675 MEZHEYEUSKI, A., BERGSLAND, C.H., BACKMAN, M., DJUREINOVIC, D., SJÖBLOM,
676 T., BRUUN, J. and MICKÉ, P. (2018). Multispectral imaging for quantitative and compartment-
677 specific immune infiltrates reveals distinct immune profiles that classify lung cancer patients.
678 *The Journal of pathology*; 244(4):421-431.
- 679 NEWMAN, A.M., LIU, C.L., GREEN, M.R., GENTLES, A.J., FENG, W., XU, Y., HOANG, C.D.,
680 DIEHN, M. and ALIZADEH, A.A. (2015). Robust enumeration of cell subsets from tissue
681 expression profiles. *Nature methods*; 12(5):453.
- 682 OGUNDIJO, O.E. and WANG, X. (2017). A sequential Monte Carlo approach to gene expression
683 deconvolution. *PLoS one*; 12(10):e0186167.
- 684 PABST, C., BERGERON, A., LAVALLÉE, V.-P., YEH, J., GENDRON, P., NORDDAHL, G.L.,
685 KROSL, J., BOIVIN, I., DENEULT, E. and SIMARD, J. (2016). GPR56 identifies primary
686 human acute myeloid leukemia cells with high repopulating potential in vivo. *Blood*:blood-
687 2015-2011-683649.
- 688 PRASSAS, I., CHRYSOJA, C.C., MAKAWITA, S. and DIAMANDIS, E.P. (2012).
689 Bioinformatic identification of proteins with tissue-specific expression for biomarker discovery.
690 *BMC medicine*; 10(1):39.
- 691 RACLE, J., DE JONGE, K., BAUMGAERTNER, P., SPEISER, D.E. and GFELLER, D. (2017).
692 Simultaneous enumeration of cancer and immune cell types from bulk tumor gene expression
693 data. *eLife*; 6.
- 694 RICHARDSON, G.M., LANNIGAN, J. and MACARA, I.G. (2015). Does FACS perturb gene
695 expression? *Cytometry Part A*; 87(2):166-175.
- 696 RUFFINI, E., ASIOLI, S., FILOSSO, P.L., LYBERIS, P., BRUNA, M.C., MACRÌ, L., DANIELE,
697 L. and OLIARO, A. (2009). Clinical significance of tumor-infiltrating lymphocytes in lung
698 neoplasms. *The Annals of thoracic surgery*; 87(2):365-372.
- 699 SEO, J.-S., KIM, A., SHIN, J.-Y. and KIM, Y.T. (2018). Comprehensive analysis of the tumor
700 immune micro-environment in non-small cell lung cancer for efficacy of checkpoint inhibitor.
701 *Scientific reports*; 8(1):14576.
- 702 SHEN-ORR, S.S., TIBSHIRANI, R., KHATRI, P., BODIAN, D.L., STAEDTLER, F., PERRY,
703 N.M., HASTIE, T., SARWAL, M.M., DAVIS, M.M. and BUTTE, A.J. (2010). Cell type-
704 specific gene expression differences in complex tissues. *Nature methods*; 7(4):287.
- 705 VAN DEN BRINK, S.C., SAGE, F., VÉRTESY, Á., SPANJAARD, B., PETERSON-MADURO,
706 J., BARON, C.S., ROBIN, C. and VAN OUDENAARDEN, A. (2017). Single-cell sequencing
707 reveals dissociation-induced gene expression in tissue subpopulations. *Nature methods*;
708 14(10):935.
- 709 WHITTEMORE, A.S. (2007). A Bayesian false discovery rate for multiple testing. *Journal of*
710 *Applied Statistics*; 34(1):1-9.
- 711 WU, A., QIN, F., CHEN, Z., QUAN, L., HUANG, A., YUAN, Y., YUAN, X., SHEN, Q., SHANG,
712 J. and BEN, Y. (2018). seq-ImmuCC: Cell-centric view of tissue transcriptome measuring
713 cellular compositions of immune microenvironment from mouse RNA-Seq data. *Frontiers in*
714 *Immunology*; 9:1286.
- 715 XIE, F., ZHOU, M. and XU, Y. (2018). BayCount: A Bayesian decomposition method for inferring
716 tumor heterogeneity using RNA-Seq counts. *The Annals of Applied Statistics*; 12(3):1605-1627.
- 717 YANG, Y.H., DUDOIT, S., LUU, P., LIN, D.M., PENG, V., NGAI, J. and SPEED, T.P. (2002).
718 Normalization for cDNA microarray data: a robust composite method addressing single and
719 multiple slide systematic variation. *Nucleic acids research*; 30(4):e15-e15.
- 720 YOSHIHARA, K., SHAHMORADGOLI, M., MARTÍNEZ, E., VEGESNA, R., KIM, H.,
721 TORRES-GARCIA, W., TREVIÑO, V., SHEN, H., LAIRD, P.W. and LEVINE, D.A. (2013).
722 Inferring tumour purity and stromal and immune cell admixture from expression data. *Nature*
723 *communications*; 4:2612.
- 724 ZHONG, Y., WAN, Y.-W., PANG, K., CHOW, L.M. and LIU, Z. (2013). Digital sorting of
725 complex tissues for cell type-specific gene expression profiles. *BMC bioinformatics*; 14(1):89.
726