

Inference in generalized bilinear models

Jeff Miller

Joint work with Scott L. Carter

Harvard T.H. Chan School of Public Health
Department of Biostatistics

IBEST-IMCI Seminar
Univ of Idaho, Apr 29, 2021

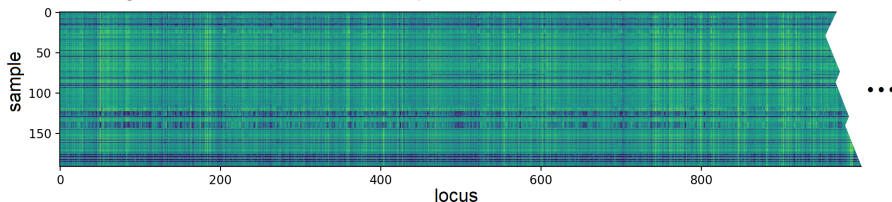
Preprint: <https://arxiv.org/abs/2010.04896>

Background

Modern high-throughput sequencing yields large matrices of counts.

- Copy ratio estimation in cancer genomics
 - ▶ whole-exome or whole-genome sequencing data
- Copy number variation in genetics
 - ▶ whole-exome or whole-genome sequencing data
- Gene expression analysis in biology/medicine
 - ▶ RNA-seq data for transcript abundance

log counts for a whole-exome seq data set of 191 samples \times 171523 loci



Background

- Latent factor models are widely used to discover and adjust for hidden variation in these applications and many others.
- Estimation and inference in latent factor models is challenging.
- Consequently, most methods do not fully account for uncertainty in the latent factors, which can lead to miscalibrated inferences such as overconfident p-values.

This talk

- Generalized bilinear models (GBMs) are a flexible extension of generalized linear models (GLMs) to include latent factors as well as row covariates, column covariates, and interactions.
- We propose fast and accurate methods for GBM estimation and inference (i.e., uncertainty quantification).
- We introduce *delta propagation*, a novel technique for propagating uncertainty among model components using the delta method.
- We present simulation studies assessing performance.
- We apply GBMs to copy ratio estimation and RNA-seq analysis.

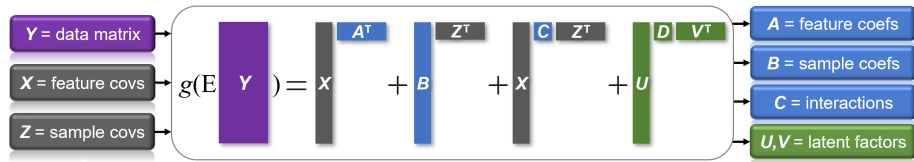
Outline

- 1 Generalized bilinear models (GBMs)
- 2 Previous work
- 3 Estimation
- 4 Inference (uncertainty quantification)
- 5 Applications
 - Copy ratio estimation in cancer genomics
 - RNA-seq gene expression analysis

Outline

- 1 Generalized bilinear models (GBMs)
- 2 Previous work
- 3 Estimation
- 4 Inference (uncertainty quantification)
- 5 Applications
 - Copy ratio estimation in cancer genomics
 - RNA-seq gene expression analysis

Generalized bilinear models (GBMs)



- Suppose the data matrix $\mathbf{Y} = (Y_{ij}) \in \mathbb{R}^{I \times J}$ satisfies

$$g(E(\mathbf{Y})) = \mathbf{X} \mathbf{A}^T + \mathbf{B} \mathbf{Z}^T + \mathbf{X} \mathbf{C} \mathbf{Z}^T + \mathbf{U} \mathbf{D} \mathbf{V}^T$$

where the link function g is applied element-wise.

- We refer to this as a *generalized bilinear model* (Choulakian, 1996).
- The “bilinear” part $\mathbf{U} \mathbf{D} \mathbf{V}^T$ is a low-rank matrix that captures latent effects due, for example, to unobserved covariates such as batch.

Identifiability and interpretation of GBM parameters

- It is important that the parameters are uniquely determined by the data distribution. We prove that identifiability holds under certain constraints.
- For interpretability, we also assume that in X and Z , the first column is all ones and the rest of the columns have mean zero.
- Then, the parameters for entry (i, j) can be interpreted as follows:

Overall intercept
Sample offset
Feature offset
Effect of k th feature covariate on sample j
Effect of ℓ th sample covariate on feature i
Interaction
Latent factors

$$g(\mathbb{E}(Y_{ij})) = c_{11} + a_{j1} + b_{i1} + \sum_{k=2}^K (c_{k1} + a_{jk})x_{ik} + \sum_{\ell=2}^L (c_{1\ell} + b_{i\ell})z_{j\ell} + \sum_{k=2}^K \sum_{\ell=2}^L c_{k\ell}x_{ik}z_{j\ell} + \sum_{m=1}^M u_{im}d_{mm}v_{jm}.$$

Outcome distributions

- We consider discrete exponential dispersion families (EDFs).

- Specifically, we suppose $Y_{ij} \sim f(y \mid \theta_{ij}, r_{ij})$ where

$$f(y \mid \theta, r) = \exp(\theta y - r\kappa(\theta))h(y, r).$$

- For any discrete EDF,

$$\begin{aligned}\mu &= E(Y) = r\kappa'(\theta) \\ \sigma^2 &= \text{Var}(Y) = r\kappa''(\theta).\end{aligned}$$

- For sequencing data, we focus on negative binomial outcomes, which is a special case of discrete EDF.
- We parametrize the dispersions as $1/r_{ij} = \exp(s_i + t_j + \omega)$.

Outline

- 1 Generalized bilinear models (GBMs)
- 2 Previous work
- 3 Estimation
- 4 Inference (uncertainty quantification)
- 5 Applications
 - Copy ratio estimation in cancer genomics
 - RNA-seq gene expression analysis

Previous work

- There is an extensive literature on models involving an unknown low-rank matrix UDV^T .
- We settle for covering only the most directly related previous work.

Previous work: Normal bilinear models without covariates

- Consider the following special case:

$$Y_{ij} = c + a_i + b_j + \sum_{m=1}^M u_{im} d_m v_{jm} + \varepsilon_{ij}$$

where $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma_{ij}^2)$.

- Principal components analysis (PCA) is equivalent to maximum likelihood estimation in this model with $\sigma_{ij}^2 = \sigma^2$.
- Estimation for this model:
Gollob (1968), Mandel (1969), Gabriel (1978), Gabriel and Zamir (1979).
- Hypothesis testing for which factors to include:
Gollob (1968), Mandel (1969), Freeman (1973), Gauch (1988, 2006).
- Confidence regions for parameters:
Goodman and Haberman (1990), Chadoeuf and Denis (1991), Dorkenoo and Mathieu (1993), Denis and Gower (1996).

Previous work: Normal bilinear models with covariates

- Consider the following special case:

$$\mathbf{Y} = \mathbf{X}\mathbf{A}^T + \mathbf{B}\mathbf{Z}^T + \mathbf{X}\mathbf{C}\mathbf{Z}^T + \mathbf{U}\mathbf{D}\mathbf{V}^T + \boldsymbol{\varepsilon}$$

where $\boldsymbol{\varepsilon}$ is a matrix of residuals with $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma_{ij}^2)$.

- Work on this model was inspired by Tukey (1962), who suggested combining regression with factor analysis.
- Estimation for this model, assuming $\sigma_{ij}^2 = \sigma^2$:
Gabriel (1978), Takane and Shibayama (1991).
- Hypothesis testing and confidence regions, assuming $\sigma_{ij}^2 = \sigma^2$:
Perry and Pillai (2013) show how to perform inference for univariate linear projections of \mathbf{A} and \mathbf{B} .

Previous work: Going beyond normal outcomes

- In many applications, it is unreasonable to assume normal outcomes.
- A classical approach is to transform the data and then assume a normal outcome model.
- However, there is unlikely to be a transformation that simultaneously achieves (a) approximate normality, (b) common variance, and (c) additive effects.
- More principled approach: Extend the generalized linear model (GLM) framework to handle latent factors, as suggested by Gower (1989).

Previous work: GBMs with covariates

- Consider the general case:

$$g(\mathbb{E}(\mathbf{Y})) = \mathbf{X}\mathbf{A}^T + \mathbf{B}\mathbf{Z}^T + \mathbf{X}\mathbf{C}\mathbf{Z}^T + \mathbf{U}\mathbf{D}\mathbf{V}^T.$$

- Previous authors have considered models of this form:

Choulakian (1996), Gabriel (1998), de Falguerolles (2000), Townes (2019).

- Townes (2019) develops a fast estimation algorithm using diagonal approximations to Fisher scoring updates for ℓ_2 -penalized estimation.
- Limitations of previous work:
 - ▶ uncertainty quantification is not addressed,
 - ▶ a single common dispersion parameter is assumed, and
 - ▶ identifiability constraints are not explicitly enforced during estimation.

Outline

- 1 Generalized bilinear models (GBMs)
- 2 Previous work
- 3 Estimation**
- 4 Inference (uncertainty quantification)
- 5 Applications
 - Copy ratio estimation in cancer genomics
 - RNA-seq gene expression analysis

Estimation algorithm

- We provide an algorithm for *maximum a posteriori* GBM estimation that extends previous work by:
 - ▶ estimating row- and column-specific dispersion parameters,
 - ▶ improving numerical stability, and
 - ▶ explicitly enforcing identifiability constraints during estimation.
- Basic idea: Iteratively cycle through the components of the model, updating each in turn using an optimization-projection step.
- “Optimization-projection” = unconstrained optimization step and a likelihood-preserving projection onto the constrained parameter space.

Estimation: Solutions to challenges

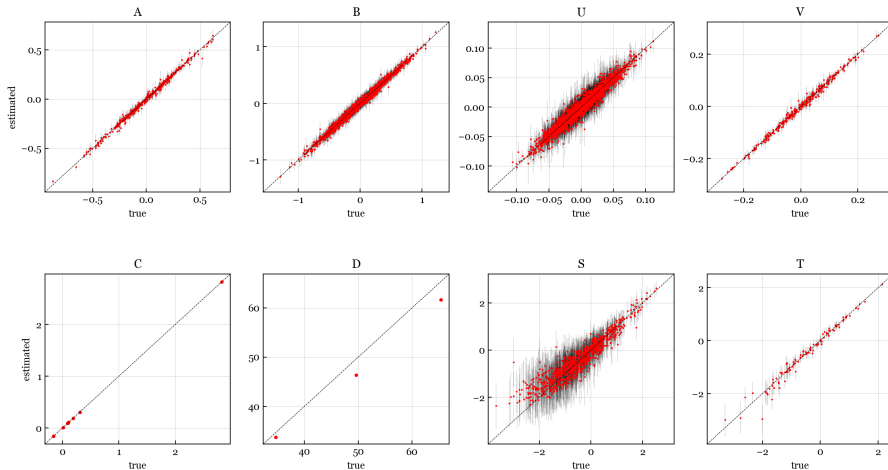
- Some key aspects of our approach:
 - ▶ Exploit the GBM structure to derive fast Fisher scoring updates.
 - ▶ Initialize using least squares for A , B , and C , with $UDV^T = 0$.
 - ▶ Use bounded, regularized Fisher scoring steps for numerical stability.
 - ▶ Derive likelihood-preserving projections to enforce constraints.
 - ▶ Relax dependencies and constraints by optimizing UD and VD rather than U and V .

Estimation: Simulation study

- We assess estimation performance in simulations with known true parameters.
- In each simulation run:
 - ▶ covariates are generated using a copula model with Normal, Gamma, or Binary marginals,
 - ▶ true parameters are generated using a Normal or Gamma scheme, and
 - ▶ outcomes are generated using the log link and a NB (negative binomial), LNP (log-normal Poisson), Poisson, or Geometric distribution.
- We abbreviate each combination of outcome/covariate/parameter scheme, e.g., NB/Binary/Normal.

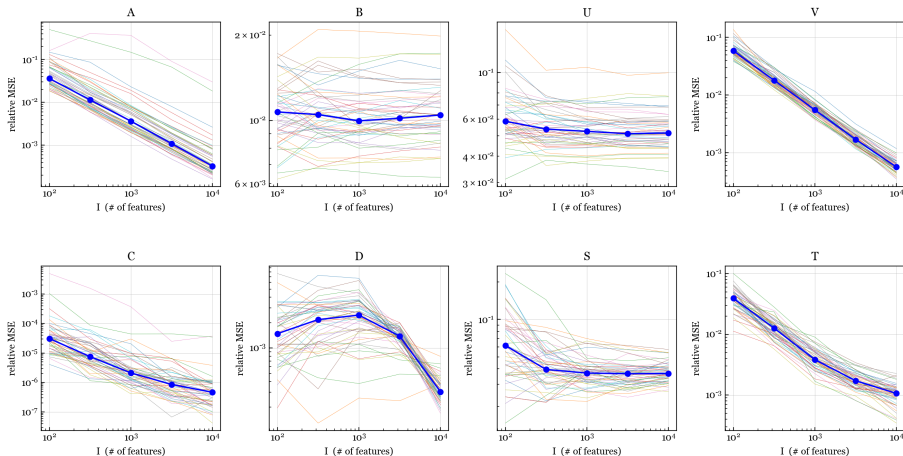
Estimation: Typical example

Scatterplots of estimated versus true parameters for a typical simulated data matrix
(NB/Normal/Normal, 1000 rows, 100 cols, 4 feature covs, 2 samples covs, and 3 factors)

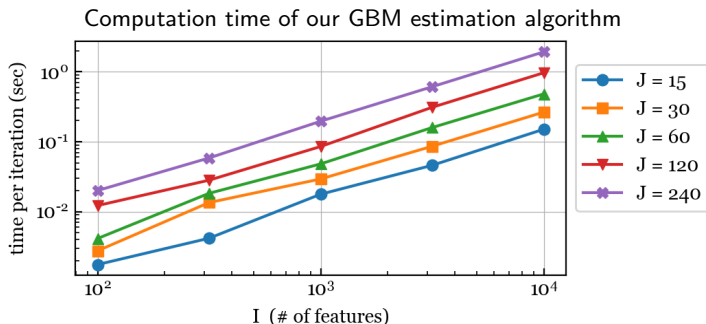


Estimation: Error tends to zero with increasing data

Relative mean-squared error between estimated and true parameter values
(50 runs of NB/Normal/Normal, 100 cols, 4 feature covs, 2 samples covs, and 3 factors)



Estimation: Computation time is linear in size of matrix



- Computation time grows linearly with I (# rows) and J (# cols).
- Each dot is the average over 10 runs of the NB/Normal/Normal scheme with 4 feature covs, 2 samples covs, and 3 factors.
- The empirical results agree with the theory.

Outline

- 1 Generalized bilinear models (GBMs)
- 2 Previous work
- 3 Estimation
- 4 Inference (uncertainty quantification)
- 5 Applications
 - Copy ratio estimation in cancer genomics
 - RNA-seq gene expression analysis

Inference (uncertainty quantification)

- Most latent factor methods do not fully account for uncertainty.
- To remove batch effects in gene expression, several methods estimate UDV^T and then treat V as known, handling uncertainty only in U .
Leek and Storey (2007, 2008), Sun et al. (2012), Risso et al. (2014)
- CNV detection methods often fit UDV^T and just subtract it off.
Fromer et al. (2012), Krumm et al. (2012), Jiang et al. (2015)
- Bayesian inference provides full uncertainty quantification, but MCMC is slow in large parameter spaces with strong dependencies.
- Variational Bayes is faster, but relies on factorized approximations that tend to underestimate uncertainty.

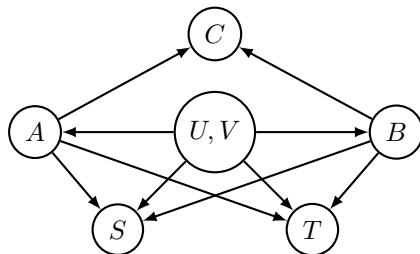
Stegle et al. (2010), Buettner et al. (2017), Babadi et al. (2018)

Inference: Novel method – “delta propagation”

- We provide a fast, accurate method for GBM uncertainty quantification.
- In particular, we introduce *delta propagation*, a general technique for propagating uncertainty among model components using the delta method.
- Delta propagation can be done analytically using closed-form expressions involving the gradient and the Fisher information.

Inference: Outline of GBM inference algorithm

Diagram of uncertainty propagation scheme for GBM inference



Outline:

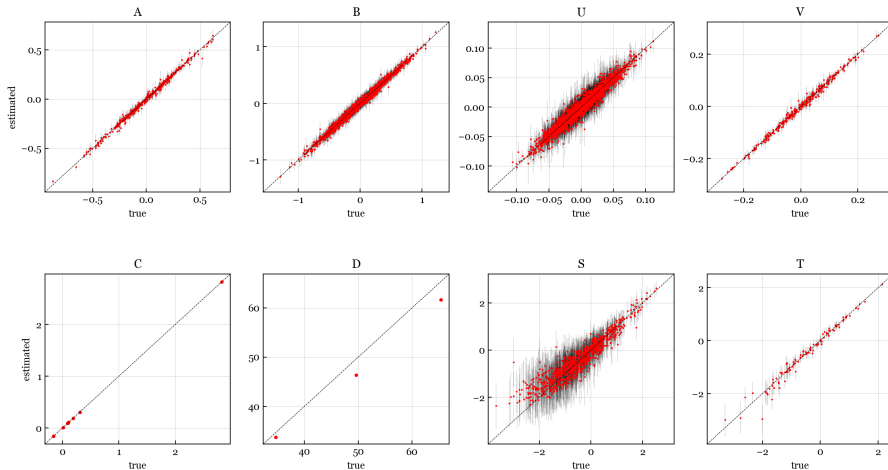
- 1 Compute conditional uncertainty for each parameter matrix/vector.
- 2 Compute joint uncertainty in (U, V) accounting for constraints.
- 3 Propagate uncertainty between components using delta propagation.
- 4 Compute approximate standard errors.

Inference: Simulation study

- To assess the accuracy of standard errors produced by our algorithm, we consider the coverage of Wald-type confidence intervals.
- Ideally, a 95% confidence interval would contain the true parameter 95% of the time.
- However, even when the model is correct, this is not guaranteed since intervals are usually based on an approximation to the distribution of an estimator.
- We generate covariates, true parameters, and outcomes using the same simulation scheme as before.

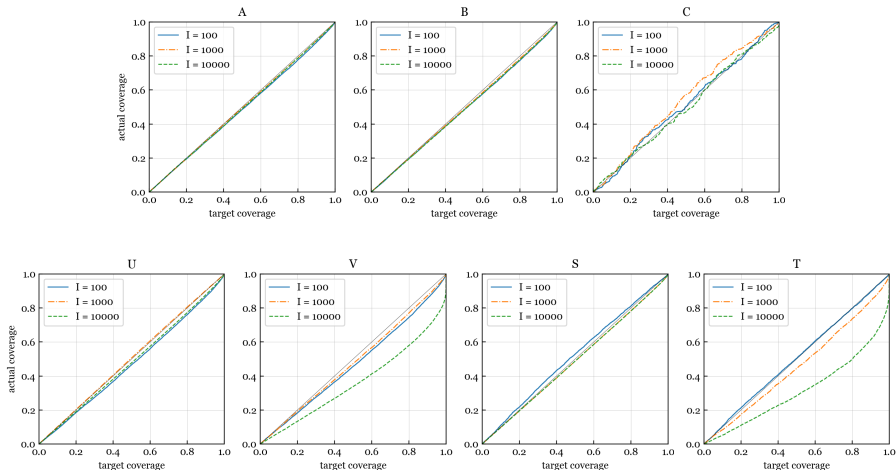
Inference: Typical example

Scatterplots of estimated versus true parameters for a typical simulated data matrix
(NB/Normal/Normal, 1000 rows, 100 cols, 4 feature covs, 2 samples covs, and 3 factors)

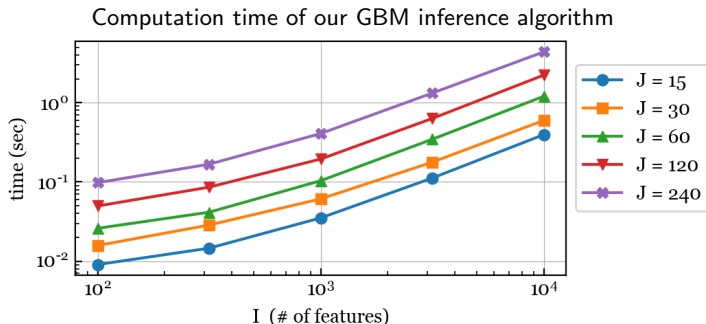


Inference: Coverage is good for most params of interest

Coverage of confidence intervals for the entries of each parameter matrix/vector
(50 runs of NB/Normal/Normal, 100 cols, 4 feature covs, 2 samples covs, and 3 factors)



Inference: Empirical assessment of computation time



- Theory indicates that computation time is linear in I (# rows) and quadratic in J (# cols).
- Thus, as I increases, the curves should become linear in I .
- Each dot is the average over 10 runs of the NB/Normal/Normal scheme with 4 feature covs, 2 samples covs, and 3 factors.

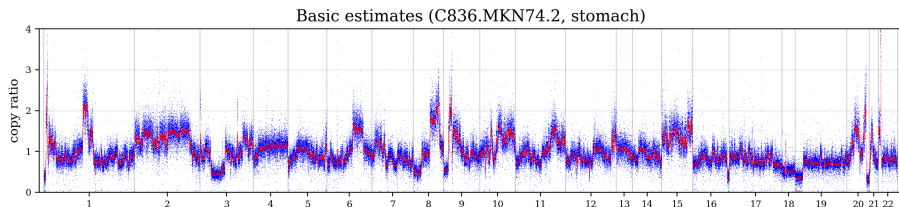
Outline

- 1 Generalized bilinear models (GBMs)
- 2 Previous work
- 3 Estimation
- 4 Inference (uncertainty quantification)
- 5 Applications
 - Copy ratio estimation in cancer genomics
 - RNA-seq gene expression analysis

Application: Copy ratio estimation in cancer genomics

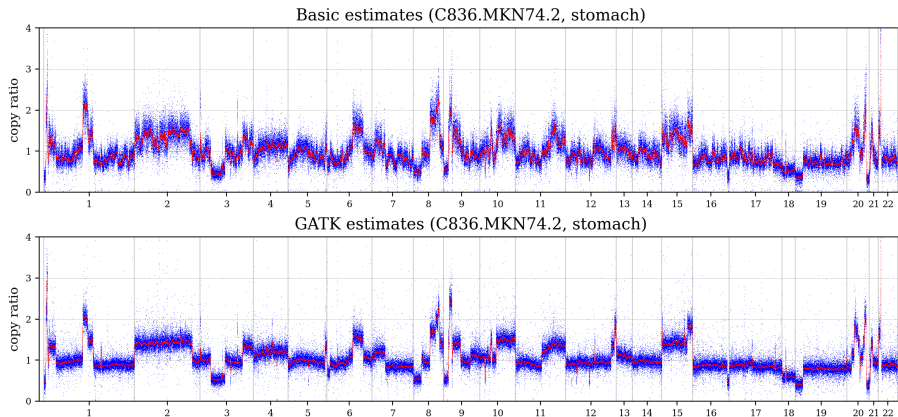
- We apply the GBM to estimate copy ratios for sequencing data.
- Copy ratio estimation is an essential step in detecting somatic copy number alterations (SCNAs), that is, duplications or deletions of segments of the genome.
- The input data is a matrix of counts where entry (i, j) is the number of reads from sample j that map to target region i of the genome.
- Goal: Estimate the copy ratio of each region, that is, the relative concentration of copies of that region in the original DNA sample.
- We illustrate on the 326 whole-exome sequencing samples from the Cancer Cell Line Encyclopedia (CCLE) (Ghandi et al., 2019).

Copy ratio estimation: Example from CCLE data



- x-axis = genomic position, blue = CR estimate, red = moving avg.
- As a baseline, we show basic row- and column- normalized estimates.
- Specifically, $\rho_{ij}^{\text{basic}} = \tilde{Y}_{ij} / (\alpha_i \beta_j)$ where $\tilde{Y}_{ij} = Y_{ij} + 0.125$, $\alpha_i = \frac{1}{J} \sum_{j=1}^J \tilde{Y}_{ij}$, and $\beta_j = \frac{1}{I} \sum_{i=1}^I \tilde{Y}_{ij} / \alpha_i$.
- These basic estimates are very noisy and are contaminated by significant technical biases.

Copy ratio estimation: Example from CCLE data



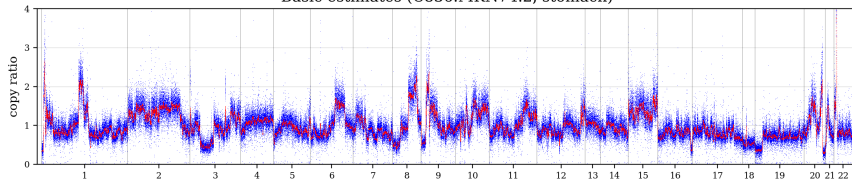
- GATK results on an illustrative sample, using a PoN with 5 factors.
- The GATK estimates are less noisy and are more locally constant.

Copy ratio estimation with the GBM

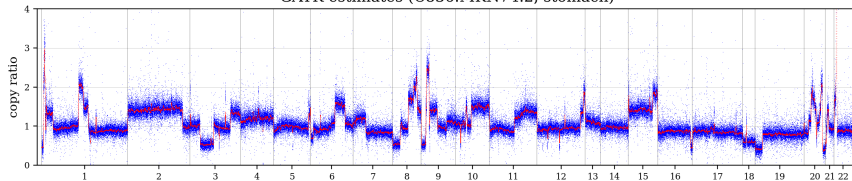
- For comparison, we run a negative binomial GBM on the adjusted (pseudo-normal) training samples to estimate latent factors U .
- Then we run a GBM on the test samples, using a feature covariate matrix X that includes this estimated U matrix.
- We use $\log(\text{length}_i)$, gc_i , and $(\text{gc}_i - \overline{\text{gc}})^2$ as region covariates, no sample covariates, and 5 latent factors.
- Model dimensions:
 - ▶ On training set: $I = 180,495$, $J = 163$, $K = 4$, $L = 1$, and $M = 5$.
 - ▶ On the test set: $I = 180,495$, $J = 163$, $K = 9$, $L = 1$, and $M = 0$.

Copy ratio estimation: Example from CCLE data

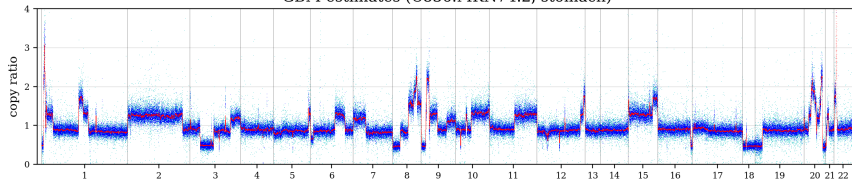
Basic estimates (C836.MKN74.2, stomach)



GATK estimates (C836.MKN74.2, stomach)

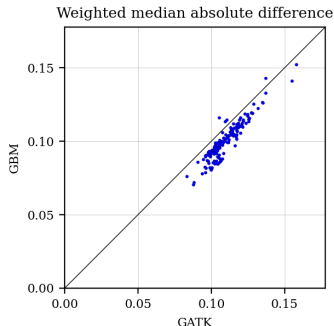
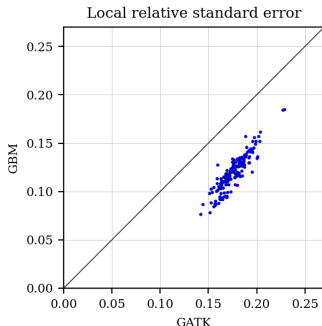


GBM estimates (C836.MKN74.2, stomach)



(In the GBM plot, estimates with low relative precision are plotted in cyan.)

Copy ratio estimation: Performance on CCLE test set

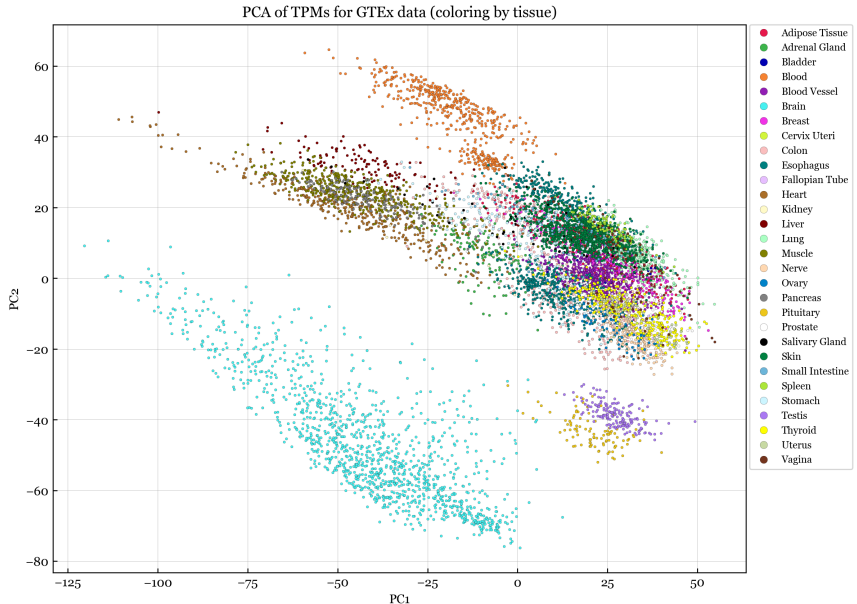


- We compare the GBM and GATK using two performance metrics:
 - ▶ Local RSE quantifies the variability of log CR estimates around a weighted moving average, accounting for the precision of each estimate.
 - ▶ Weighted MAD quantifies the typical magnitude of the slope of a weighted moving average.
- The performance gains appear to be due to using (a) model-based uncertainty and (b) a robust probabilistic model for count data.

RNA-seq: Analyzing GTEx data for aging-related genes

- We consider RNA-seq data from the Genotype-Tissue Expression (GTEx) project (Melé et al., 2015).
- 8,551 samples from 30 tissues in the human body, from 544 subjects.
- We apply the GBM to find genes whose expression changes with age, adjusting for technical biases.

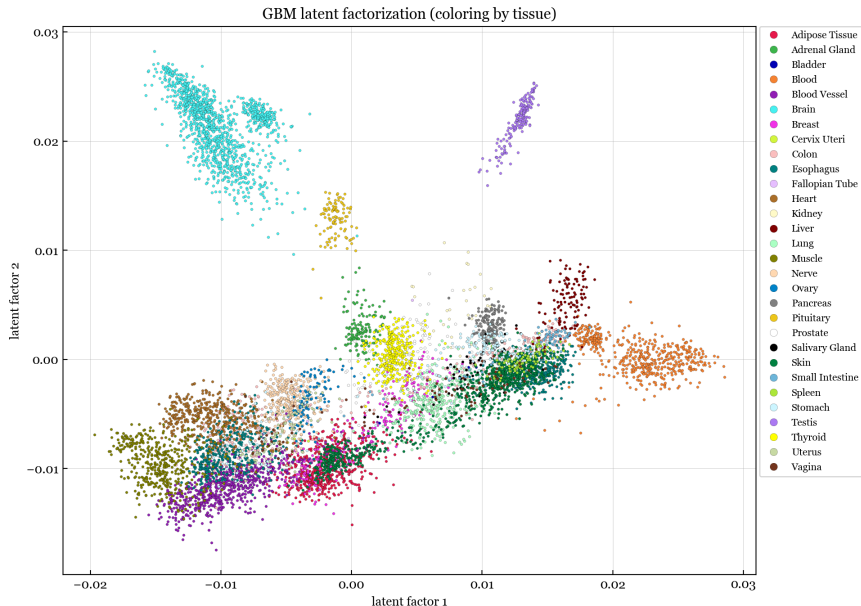
RNA-seq: PCA of GTEx data using log-transformed TPMs



RNA-seq: Visualizing GTEx data using a GBM

- Similar to PCA, we can use the GBM to visualize high-dimensional data by plotting the V matrix.
- First, we take a random subset of 5,000 genes and fit a negative binomial GBM with:
 - ▶ two latent factors,
 - ▶ no sample covariates, and
 - ▶ $\log(\text{length}_i)$, gc_i , and $(\text{gc}_i - \overline{\text{gc}})^2$ as gene covariates.
- Model dimensions: $I = 5,000$, $J = 8,551$, $K = 4$, $L = 1$, and $M = 2$.

RNA-seq: Visualizing GTEx data using a GBM



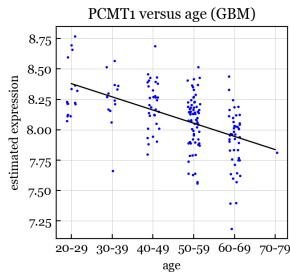
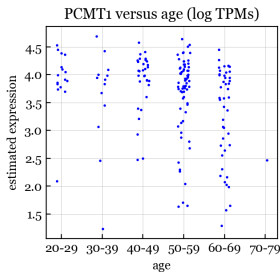
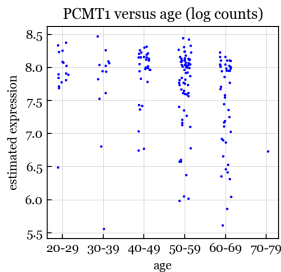
RNA-seq: Analyzing GTEx data for aging-related genes

- To find aging-related genes, we add subject age as a sample covariate.
- For illustration, we present results for the “Heart - Left Ventricle” subtissue (Heart-LV).
- We ran the GBM on the 176 Heart-LV samples in the test set, using:
 - ▶ the 19,853 genes with nonzero median across these samples,
 - ▶ gene covariates: $\log(\text{length}_i)$, gc_i , and $(gc_i - \overline{gc})^2$,
 - ▶ sample covariates: `smexncrt` (exonic rate) and age (subject age),
 - ▶ 3 latent factors.
- This choice of subtissue and model was based on the exploratory phase.

RNA-seq: Analyzing GTEx data for aging-related genes

- In this GBM, each gene has a coefficient describing how its expression changes with age.
- Using our GBM inference algorithm, we compute a p-value for each gene to test whether this coefficient is nonzero.
- 2,444 genes were significantly associated with age in Heart-LV, controlling Type I error at 0.05 using Bonferroni.
- For comparison, simple linear regression on the log-transformed TPMs yields only 1 significant gene.
- This indicates that the GBM has much greater power than a simple standard approach.

RNA-seq: Expression of the top aging-related gene



- The top GBM hit for Heart-LV is PCMT1 ($p\text{-value} = 1.1 \times 10^{-47}$).
- PCMT1 is involved in the repair and degradation of damaged proteins, and is a well-known aging gene (Tacutu et al., 2018).
- GBM-estimated expression of PCMT1 exhibits a clear downward linear trend with age.
- The log TPMs for PCMT1 are noisier and the trend is much less clear.

RNA-seq: Top age-related GO terms (Biological Process)

- To test for enrichment of Gene Ontology (GO) term gene sets, we run DAVID on the top 1000 GBM hits for Heart-LV.
- These results are highly consistent with known aging biology (López-Otín et al., 2013).

GO term ID	Description	Count	p-value	Benjamini
GO:0098609	cell-cell adhesion	48	5.1e-12	1.5e-08
GO:0006418	tRNA aminoacylation for protein translation	16	1.4e-09	2.0e-06
GO:0006099	tricarboxylic acid cycle	12	3.7e-07	3.6e-04
GO:1904871	positive regulation of protein localization to Cajal body	7	1.1e-06	6.1e-04
GO:1904851	positive regulation of establishment of protein localization to telomere	7	1.1e-06	6.1e-04
GO:0006607	NLS-bearing protein import into nucleus	10	1.3e-06	6.2e-04
GO:0006914	autophagy	22	1.8e-05	7.6e-03
GO:0016192	vesicle-mediated transport	24	2.6e-05	8.3e-03
GO:0006511	ubiquitin-dependent protein catabolic process	24	2.6e-05	8.3e-03
GO:0006888	ER to Golgi vesicle-mediated transport	24	3.5e-05	1.0e-02
GO:0006886	intracellular protein transport	31	4.3e-05	1.1e-02
GO:1904874	positive regulation of telomerase RNA localization to Cajal body	7	8.3e-05	2.0e-02
GO:0006090	pyruvate metabolic process	8	9.6e-05	2.1e-02
GO:0070125	mitochondrial translational elongation	16	1.1e-04	2.2e-02
GO:0006446	regulation of translational initiation	10	1.5e-04	2.8e-02
GO:0043039	tRNA aminoacylation	5	1.6e-04	3.0e-02
GO:0018107	peptidyl-threonine phosphorylation	10	2.9e-04	4.9e-02
GO:0000462	maturation of SSU-rRNA from tricistronic rRNA transcript	9	3.3e-04	5.4e-02
GO:0006610	ribosomal protein import into nucleus	5	3.7e-04	5.6e-02
GO:0016236	macroautophagy	14	4.0e-04	5.9e-02

RNA-seq: Top age-related GO terms (Cellular Component)

GO term ID	Description	Count	p-value	Benjamini
GO:0016020	membrane	220	9.8e-21	3.7e-18
GO:0005739	mitochondrion	157	1.2e-20	3.7e-18
GO:0070062	extracellular exosome	242	4.3e-16	9.1e-14
GO:0005829	cytosol	282	1.0e-15	1.6e-13
GO:0005913	cell-cell adherens junction	57	9.5e-15	1.2e-12
GO:0005737	cytoplasm	380	2.3e-13	2.4e-11
GO:0043209	myelin sheath	36	4.7e-13	4.2e-11
GO:0005759	mitochondrial matrix	47	5.7e-09	4.5e-07
GO:0005654	nucleoplasm	217	1.1e-08	7.8e-07
GO:0000502	proteasome complex	18	1.4e-08	8.0e-07
GO:0005743	mitochondrial inner membrane	56	1.4e-08	8.0e-07
GO:0042645	mitochondrial nucleoid	14	3.5e-07	1.8e-05
GO:0014704	intercalated disc	14	8.5e-07	4.2e-05
GO:0005832	chaperonin-containing T-complex	7	2.5e-06	1.1e-04
GO:0005643	nuclear pore	16	5.2e-06	2.2e-04
GO:0043231	intracellular membrane-bounded organelle	55	2.7e-05	1.1e-03
GO:0002199	zona pellucida receptor complex	6	2.9e-05	1.1e-03
GO:0043034	costamere	8	5.4e-05	1.9e-03
GO:0043234	protein complex	42	7.8e-05	2.6e-03
GO:0045254	pyruvate dehydrogenase complex	5	1.5e-04	4.6e-03

Conclusion

- GBMs provide a flexible framework for the analysis of matrix data.
- Delta propagation is a novel general technique for uncertainty quantification.
- Our algorithms enable accurate GBM estimation and inference in modern applications.
- Possible directions for future work:
 - ▶ extend to more general bilinear model structures,
 - ▶ seek theoretical guarantees for delta propagation, and
 - ▶ try applying delta propagation to other models.
- Preprint is on arXiv: <https://arxiv.org/abs/2010.04896>

Inference in generalized bilinear models

Jeff Miller

Joint work with Scott L. Carter

Harvard T.H. Chan School of Public Health
Department of Biostatistics

IBEST-IMCI Seminar
Univ of Idaho, Apr 29, 2021

Preprint: <https://arxiv.org/abs/2010.04896>

Backup slides

Identifiability of GBMs

- For reliable results, it is important to ensure that the parameters are uniquely determined by the data distribution.
- We prove that identifiability holds under the following constraints:
 - ▶ $X^T X$ and $Z^T Z$ are invertible,
 - ▶ $X^T B = 0$, $Z^T A = 0$, $X^T U = 0$, and $Z^T V = 0$,
 - ▶ $U^T U = I$ and $V^T V = I$,
 - ▶ D is a diagonal matrix such that $d_{11} > d_{22} > \dots > d_{MM} > 0$, and
 - ▶ the first nonzero entry of each column of U is positive.
- More precisely, the function

$$\eta(A, B, C, D, U, V) = XA^T + BZ^T + XCZ^T + UDV^T$$

is one-to-one on the set of parameters satisfying these constraints.

Previous work: GBMs without covariates

- Consider the following special case:

$$g(E(Y_{ij})) = c + a_i + b_j + \sum_{m=1}^M u_{im}d_mv_{jm}.$$

- This allows non-normal outcomes, but does not include covariates.
- Estimation for this model:
Goodman (1979, 1981, 1986, 1991), Van Eeuwijk (1995).
- Hypothesis testing for which factors to include:
Van Eeuwijk (1995).

Estimation: Challenges (1/2)

- 1 Estimating the dispersions is tricky due to nonobvious biases, arithmetic underflow/overflow, and occasional lack of convergence.
- 2 Standard GLM methods are inapplicable. Even without UDV^T , vectorization of the linear terms is computationally prohibitive.
- 3 Optimizing UDV^T is challenging due to the dependencies among U , D , and V and the orthonormality constraints $U^T U = I$ and $V^T V = I$.
- 4 The singular value decomposition (SVD) doesn't help estimate UDV^T since it implicitly assumes every entry has the same variance.

Estimation: Challenges (2/2)

- 5 Computational efficiency is needed to handle large high-throughput sequencing datasets.
- 6 A good initialization procedure is crucial for numerical stability.
- 7 Even with a good initialization, optimization methods occasionally diverge. In a large GBM, there are so many parameters that even occasional divergences lead to failure with high probability.
- 8 It is not obvious how to enforce the identifiability constraints without compromising the algorithm convergence properties.

Inference: Delta propagation method

- In fixed-dimension parametric models, the asymptotic covariance of the MLE is equal to the inverse of the Fisher information matrix.
- However, inverting the full Fisher info is intractable in large GBMs.
- Inverting the Fisher info for each component (e.g., F_a^{-1} for A) is fast, but underestimates uncertainty since it treats all else as known.
 - ▶ Thus, it can be thought of as the conditional uncertainty.
- Delta propagation is a general technique for approximating the additional variance due to uncertainty in the other components.
- Basic idea: Write the estimator for each component as a function of the other components, and propagate the variance of the other components through this function using a 1st order Taylor approx.

Estimation: Theoretical computational complexity is linear

Computation time complexity of each update in the estimation algorithm

Operation	Time complexity
Computing η	$O(IJ \max\{K, L, M\})$
Updating A	$O(IJK^2)$
Updating B	$O(IJL^2)$
Updating C	$O(IJ \max\{K^2, L^2\})$
Updating D, U , and V	$O(IJM^2)$
Updating S and T	$O(IJ)$
Total per iteration	$O(IJ \max\{K^2, L^2, M^2\})$

Notation:

- ▶ I = # of rows
- ▶ J = # of columns
- ▶ K = # of feature covariates
- ▶ L = # of sample covariates
- ▶ M = # of latent factors

Inference: Theoretical computational complexity

Computation time complexity of the inference algorithm

Operation	Time complexity
Preprocessing	$O(IJ \max\{K, L, M\})$
Conditional uncertainty for each component	$O(IJ \max\{K^2, L^2, M^2\})$
Joint uncertainty in (U, V)	$O(IJ^2 M^3)$
Propagate uncertainty between components	$O(IJ \max\{K^3, L^3, M^3\})$
Compute approximate standard errors	$O(IJ)$
Total	$O(IJ \max\{K^3, L^3, JM^3\})$

Notation:

- ▶ I = # of rows
- ▶ J = # of columns
- ▶ K = # of feature covariates
- ▶ L = # of sample covariates
- ▶ M = # of latent factors

We have experimented extensively but have not found a faster alternative that provides well-calibrated standard errors.

Copy ratio estimation using a panel of normals

- Leading methods employ a panel of normals (PoN) to estimate technical biases using PCA.
- Cancer samples are then de-noised by adjusting out the top PCs that were estimated from the PoN.
- GATK's `CreateReadCountPanelOfNormals` and `DenoiseReadCounts` tools provide CR estimates using this approach.
- For reproducibility purposes, we use a pseudo-PoN from CCLE:
 - ▶ Split the 326 CCLE samples into training and testing sets of equal size.
 - ▶ On the training samples, segment the basic CR estimates and subtract off the segment means (in log space).
 - ▶ Run `CreateReadCountPanelOfNormals` on the adjusted training data.
 - ▶ Run `DenoiseReadCounts` on the test data using the resulting PoN file.

RNA-seq: Analyzing GTEx data for aging-related genes

- To find aging-related genes, we add subject age as a sample covariate.
- We analyze each subtissue separately, due to the heterogeneity of tissues/subtissues.
- We used a random subset of 108 subjects during an exploratory model-building phase.
- The remaining 436 subjects were used during a testing phase with the selected model.