# Undergraduate research and graduate school so far

Jason Xu, Department of Statistics

ACMS Seminar Jan 14, 2016

#### Why undergraduate research?

First of all, it's fun

- Classroom assignments can feel boring and contrived
- Work on new and relevant questions
- Low-stress way to decide whether research is for you
- Collaborate with excited and interested peers
- Travel to research programs, conferences, etc.

#### Some practical considerations

Second, it's good for you

- Closer relationship with professors
- Looks good on CV
- Attractive to employers
- Great preparation for grad school
- You often get paid

#### My undergraduate research experiences

My professor informed me about (and got me into) an NSF REU program in combinatorics and number theory

- Spent the summer in Oregon researching combinatorial game theory
- The following summer: REU in Claremont on Markov chain Monte Carlo techniques
- This one got me seriously interested in research, specifically in statistics and stochastic processes
  - I am now in grad school
  - I am in the statistics department
  - I research stochastic processes

#### Next steps

- Continued in that area as a research assistant at my home institution
- Had something to write about when applying to graduate school, fellowships, etc.
- Had stronger recommendation letters from research mentors
- Summer *after* graduation, participated in a summer program working on research problems posed by various companies (UCLA RIPS)

#### Some takeaways

- NSF REU programs are a great resource!
- UMN Duluth, Williams (SMALL), Cornell, Wisconsin-Madison (Emory?)
- Other unique programs sponsored by industry and other foundations
- Don't be afraid to apply, email professors, etc.

#### Some things I've worked on in grad school

Briefly, I've worked on inference for dynamic stochastic systems when computationally expensive techniques such as MCMC are intractable.

- Approximate Bayesian treatment of massive hidden Markov models
  - Variational inference, efficient optimization, parallel computing
  - "Machine learning", application to pattern discovery for chromatin segmentation in genomics
- Likelihood methods for partially observed Markov branching processes
  - PDEs, Fourier methods, generating functions
  - "Methodology", applications to molecular epidemiology

#### Some things I've worked on in grad school

- Accelerating generating function techniques via sparsity
  - "Computation", compressed sensing, proximal methods
- Smooth, minimally Lipschitz function interpolation
  - "Theory", analysis, empirical processes, convex optimization

Currently, developing parameter estimation techniques for new models of cell differentiation applied to recent single-cell lineage tracking of *in vivo* hematopoiesis

#### What is hematopoiesis?

Complex system in which self-renewing hematopoietic stem cells (HSCs) differentiate via a series of intermediate progenitor cell stages to produce blood cells

- Dynamics and structure are largely unknown
- Clinically important: stem cell transplantation is a mainstay of cancer therapy; all blood cell diseases are caused by malfunctions in the hematopoietic process
- Modeled using stochastic compartmental models: cells replicate and differentiate as a continuous-time *Markov branching process*

#### Proposed models of hematopoiesis



©2001 Terese Winslow. Lydia Kibiuk

#### Branching structure



#### Controversial tree structure

Science DOI: 10.1126/science.aab2116

Read Full Text to Comme

RESEARCH ARTICLE

## Distinct routes of lineage development reshape the human blood hierarchy across ontogeny

Faiyaz Notta<sup>1,1,2,\*,1</sup>, Sasan Zandi<sup>1,1,2,\*</sup>, Naoya Takayama<sup>1,2</sup>, Stephanie Dobson<sup>1,2</sup>, Olga I. Gan<sup>1</sup>, Gavin Wilson<sup>2,4</sup>, Kerstin B. Kaufmann<sup>1,2</sup>, Jessica McLeod<sup>1</sup>, Elisa Laurenti<sup>6</sup>, Cyrille F. Dunant<sup>7</sup>, John D. McPherson<sup>3,4</sup>, Lincoln D. Stein<sup>2,4</sup>, Yigal Dror<sup>5</sup>, John E. Dick<sup>1,2,±</sup>



Stem-cell scientists led by Dr. John Dick have discovered a completely new view of how human blood is made, upending dogma from the 1960s.



#### Controversial tree structure



#### Statistical simplifications



Limited to overly simple model due to statistical challenges and resolution of data, yet estimation is already difficult. Past attempts include simulation studies, normal approximation, reverse jump MCMC, method of moments.

Our dataset: rhesus macaques lineage tracking





#### Clonal Tracking of Rhesus Macaque Hematopoiesis Highlights a Distinct Lineage Origin for Natural Killer Cells

Chuanfeng Wu,<sup>1,2</sup> Brian Li,<sup>1,7</sup> Rong Lu,<sup>2,7</sup> Samson J. Koelle,<sup>1</sup> Yanqin Yang,<sup>3</sup> Alexander Jares,<sup>1</sup> Alan E. Krouse,<sup>1</sup> Mark Metzger,<sup>1</sup> Frank Liang,<sup>4</sup> Karin Loré,<sup>4</sup> Colin O. Wu,<sup>5</sup> Robert E. Donahue,<sup>1</sup> Irvin S.Y. Chen,<sup>6</sup> Irving Weissman,<sup>2</sup> and Cynthia E. Dunba<sup>-1,\*</sup>

<sup>1</sup>Hematology Branch, National Heart, Lung and Blood Institute, National Institutes of Health, Bethesda, MD 20892, USA <sup>2</sup>Institute for Stem Cell Biology and Regenerative Medicine, Stanford University School of Medicine, Palo Alto, CA 94305, USA

NIH group: cell lineage barcoding data in primates. Their in-house analysis has already produced interesting findings.

#### Our dataset: rhesus macaques lineage tracking Experimental Design



Goal: develop model framework and estimation methods

### A hidden stochastic model



#### Data and challenges:

- At each observation time, the experimental protocol yields a read count corresponding to each barcode ID present in each cell type sample
- Data: independent, identically distributed time series of observations from the hematopoietic process
- $\bullet \, \Rightarrow$  fit more complex models than previous statistical studies

#### Challenges

- Hidden process, discrete observations, experimental noise
- Huge latent populations, many interacting types

#### Richer latent branching models



Figure: Our class of models allows an arbitrary number of progenitors and mature cell types

#### Inference: our method

**Loss function estimation:** we match population correlations with empirical read count correlations across barcode lineages via minimizing the *loss function* 

$$\mathcal{L}(oldsymbol{ heta}; \mathbf{Y}) = \sum_{J} \left[ \psi_j(oldsymbol{ heta}; \mathbf{Y}) - \hat{
ho}_j(\mathbf{Y}) 
ight]^2,$$

- $\hat{\rho}_j$  denotes empirical correlation across IDs p at time  $t_j$
- ψ<sub>j</sub>(θ; Y) denotes model-based correlations that we calculate mathematically, given current parameters θ

Estimating best parameters reduces to **nonlinear least** squares optimization:  $\hat{\theta} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta; \mathbf{Y}).$ 

#### Results



Pairwise correlations, three mature cell groupings

#### Results



Pairwise correlations in zh33, one progenitor

#### Results



Pairwise correlations in zh33, 3 progenitors

### Some findings

- HSC self-renewal rate  $\hat{\lambda}$  estimated to be about once every 31 weeks, consistent with previous studies focusing on the stem cell
- Estimate of initial marking levels indicates  $\sim 15\%$  of cells are marked at the HSC level, while others are marked further down the line, and is consistent with side information from fluorescence data
- Newly estimated intermediate differentiation rates are reasonable given information about cell abundances and lineage relations from initial analyses
- Obtaining rate estimates in a generative model for the data is novel; initial distributions previously unidentifiable

## Summary

- Richer, more flexible statistical models of hematopoiesis and experimental protocol than previous studies
- First parameter estimation methods for fitting time-series data from lineage barcoding experiments
- Efficient method applies to a general class of multi-compartmental models
- Future: likelihood-based approaches, model selection strategies

Thanks for listening! More questions: jasonxu@uw.edu