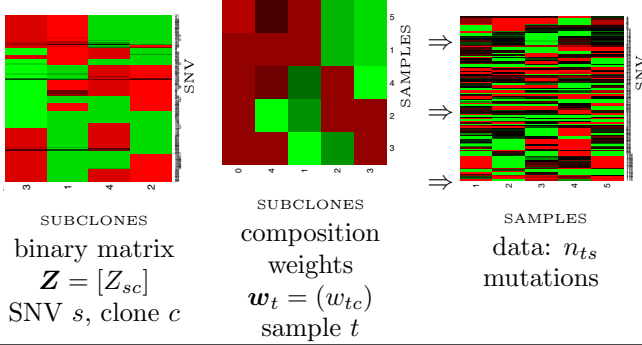


- (iii) Prior $p(\mathbf{Z})$ on $(S \times C)$ binary matrix \mathbf{Z} ,
prior $p(\mathbf{w})$ on mixture weights w_{tc} for composition
(i).

Modeling Tumor Heterogeneity

PETER MÜLLER and YANXUN XU, UT Austin

JUHEE LEE, UCSC, YUAN JI, U Chicago & NorthShore, K.
GULUKOTA, NorthShore Health System



Tumor Heterogeneity

- Mutations acquired over a tumor's life history
- Every new mutation gives rise to a new subpopulation of cells ("subclone")
- \rightarrow heterogeneous population of cells, composed of subpopulations with varying numbers of mutations.
- Tumor history imprinted in each sample as the mosaicism of mutations.

Data

SNV: point mutations, $s = 1, \dots, S$

Data: $N_{st} = \#$ reads mapped to locus of SNV s in sample t .
 $n_{st} = \#$ of these *with* SNV.

Sampling model: $n_{st} \sim \text{Bin}(N_{st}, p_{st})$

Prior: in words,

- (i) p_{st} arises as a composition of sample t as a mixture of C latent cell subclones.
- (ii) Mutation s in subclone c is either present ($Z_{sc} = 1$) or not ($Z_{sc} = 0$).
 $\mathbf{Z}_c = (Z_{sc}, s = 1, \dots, S)$ defines subclone, c .

Inference

Goal: Reconstruct cell subpopulations = estimate \mathbf{Z} and C .

Problem: Deconvolution of p_{st} as a mixture of binary indicators Z_{sc}

$$p_{st} = \sum_c w_{tc} Z_{sc} + w_{t0} p_{s0}$$

plus "background noise"

Real problem: \mathbf{Z} is latent, need to infer \mathbf{Z} from the data.

Identifiability: In principle even feasible with one sample.
Weights are identified across mutations s .

Prior

Latent cell types: $p(\mathbf{Z})$ on $(S \times C)$ binary matrix, w. random C .

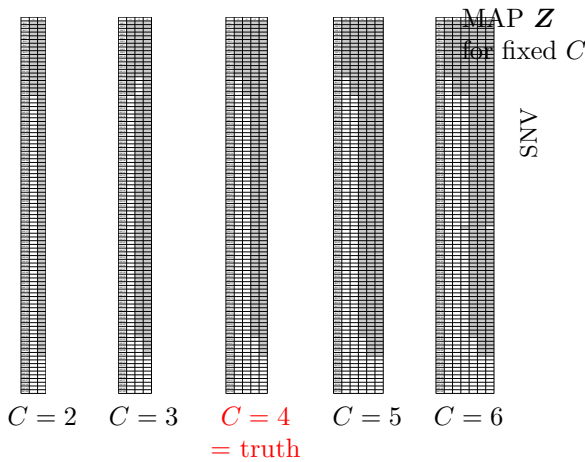
Feature allocation: Think of SNV s selecting cell types c
Features (dishes) = c ; experimental units (customers) = s

Random feature allocation: define $p(\mathbf{Z})$ as

- $p(Z_{sc} = 1 | \pi_c) = \pi_c, c = 1, \dots, C$
- $\pi_c \sim \text{Be}(\frac{\alpha}{C}, 1)$
- Drop unselected features

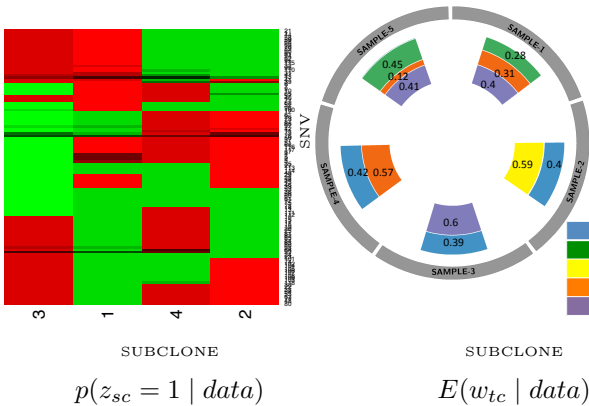
IBP as $C \rightarrow \infty$.

Composition of sample t as mix of cell types:
 $(w_{tc}, c = 1, \dots, C) \sim \text{Dir}(\cdot)$.



Results – Pancreatic Cancer

$n = 5$ samples of pancreatic cancer (PDAC, pancreatic ductal adenocarcinoma).



Computation

MCMC: random $C \rightarrow$ transdimensional MCMC, to move across C

Sloooow mixing: challenging to define good transition prob with reasonable mixing, like in any model comparison with informative data & high posterior correlation. posterior simulation for given C is straightforward.

IBF model comparison: Casella & Moreno (2006 JASA) propose propose model selection with IBF (intrinsic Bayes factor). Stylized simplified version:

FBF Model Comparison

Casella & Moreno’s IBF model comparison:

1. split data into training data $n'_{st} = \beta n_{st}; N'_{st} = \beta N_{st}$ and test data $n''_{st} = (1 - \beta)n_{st}$ and $N'' \dots$
2. informative prior $p_1(\mathbf{Z}, \mathbf{w}, p_0 | C) \equiv p_1(\mathbf{Z}, \mathbf{w}, p_0 | n', C)$
3. average over all possible splits.

We do the same, but w/o step 3 and

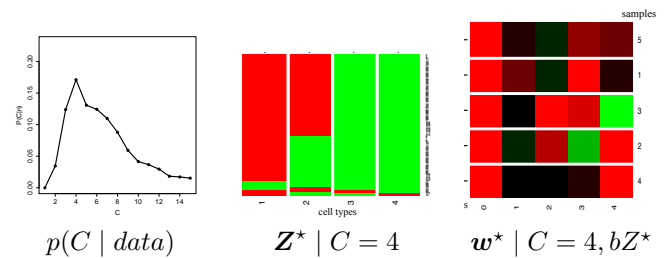
4. use MH proposal

$$q(\tilde{C}, \tilde{\mathbf{Z}}, \tilde{\mathbf{w}}, \tilde{p}_0) = q(\tilde{C} | C) p_1(\tilde{\mathbf{Z}}, \tilde{\mathbf{w}}, \tilde{p}_0 | \tilde{C})$$

Acceptance ratio simplifies and we get a well mixing MC

In the binomial likelihood, splitting the data can also be characterized as fractional BF (O’Hagan 1995).

Pancreatic Cancer



MAD Bayes for TH

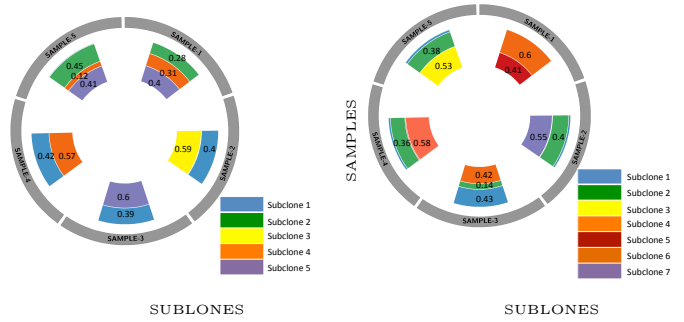
with YANXUN XU, UT Austin; YUAN YUAN, Baylor C.of Med.; YUAN JI and KAMALAKAR GULUKOTA, NorthShore Hospital.

DP mixture: Kulis & Jordan (2012) recognize log posterior \approx criterion function in k-means – voila! This is for normal sampling, asymptotically for small variance and shrinking total mass.

IBP: Broderick et al. (2013) extend a similar argument to the IBP, with normal sampling and small variance and shrinking rate of new features,

IBP with binomial sampling: same argument can be made :-)
 using increasing scaling of Bin with β and shrinking IBP par γ , using $\gamma = \exp(-\beta\lambda^2)$

Approx posterior: use k-means with different starting values to characterize posterior.



Slide 12

$S = 118$ SNV's in KEGG pathway

$S = 7000$ SNV's

Slide 15

log posterior \approx k-means like criterion:

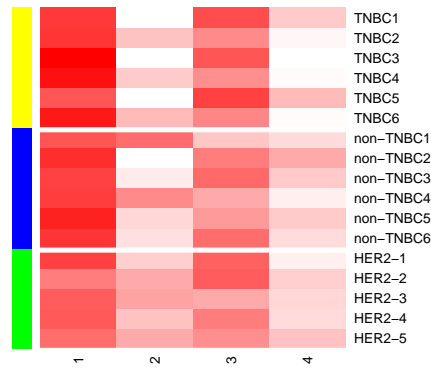
$$-\frac{1}{\beta} \log L \approx \sum_{s=1}^S \sum_{t=1}^T \{-n_{st} \log(p_{st}) - (N_{st} - n_{st}) \log(1 - p_{st})\} + C\lambda^2$$

Iterative FL-means algorithm:

- min w.r.t. \mathbf{Z} for fixed \mathbf{w}
- min w.r.t. \mathbf{w} for fixed \mathbf{Z}
- consider incrementing $\mathbf{Z} \rightarrow \tilde{\mathbf{Z}}$ with one (randomly selected) SNV and min w.r.t. \mathbf{w} . keep if it lowers the criterion

Slide 13

Results - Breast Cancer
 Horvath et al. (2013): $n = 17$ BC patients, $S = 329$ SNV's.

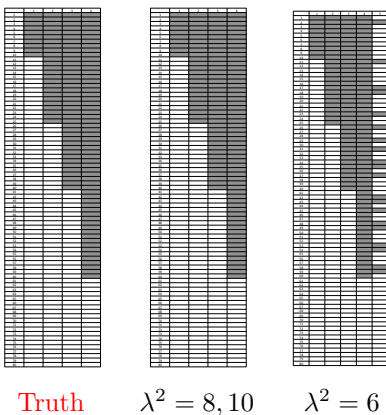


Estimated w_{tc} with $C = 4$.

Slide 16

Simulation

True and estimated \mathbf{Z}



Slide 14

Haplotypes vs. Subclones

Haplotypes: for a diploid organism, pairs of haplotypes define a unique genome.

Prior $p(\mathbf{w})$: this implies that w_{tc} should be (a mixture of) tied in pairs, i.e., $w_{tc} = w_{tc'}$, with *pairs* defining subclones.

cIBP:

Alternatively, we use a categorical generalization of the IBP, $p(\mathbf{Z})$ for a random trinary matrix with $z_{sc} \in \{0, 1, 2\}$

| | | | | | |
|----|-----|-----|-----|-----|-----|
| | 1 | 2 | 3 | 4 | 5 |
| 1 | 0.5 | 1 | 0 | 1 | 0 |
| 2 | 1 | 0.5 | 1 | 1 | 1 |
| 3 | 0.5 | 0 | 0 | 0 | 0.5 |
| 4 | 0.5 | 0 | 0.5 | 0 | 0.5 |
| 5 | 1 | 1 | 0.5 | 0.5 | 0.5 |
| 6 | 1 | 0 | 0.5 | 0 | 0 |
| 7 | 1 | 0 | 0 | 0 | 0 |
| 8 | 1 | 0.5 | 0 | 0.5 | 1 |
| 9 | 1 | 0.5 | 1 | 1 | 1 |
| 10 | 0.5 | 0 | 0 | 0 | 1 |

Slide 17

Results - Pancreatic Cancer

$n = 5$ samples of pancreatic cancer (PDAC, pancreatic ductal adenocarcinoma). Estimated w_{tc} :

Summary

TH: Model-based estimation of cell subpopulations is possible – and seems to work.

Big data: MCMC is not feasible anymore – alternative approaches remain feasible.

Limitations: and extensions

Tumor phylogenetics: Without condition on phylogenetic tree of subclones

A priori independent cell types: independent $\mathbf{z}_c = (Z_{1\dots S,c})$, with $p(\mathbf{z}_c = \mathbf{z}_{c'}) > 0$, a priori (i know – arrgh!)

Alternative dependent prior using DPP or others.

CNV: we conditioned on N_{st} .

Could use N_{st} to learn about CNV.

SOFTWARE: R code at <http://www.ma.utexas.edu/users/yxu/>