# Smaller $p$-values in genomics studies using distilled auxiliary information

Jordan Bryan joint work with Peter D. Hoff

Duke University, Department of Statistical Science

# Table of contents

# Multiple hypothesis testing in genomics

**Functional genomics** is the field that seeks to catalogue the function of genes and their protein products

Data generation in functional genomics has exploded in recent years, thanks to advent of RNAi and CRISPR, improvements in DNA and RNA sequencing.

*'Omics' data can provide information on the size and composition of biological entities and thus determine the boundaries of the problem at hand. Biologists can then proceed to investigate function using classical hypothesis-driven experiments. It is still unclear whether even this marriage of the two methods will deliver a complete understanding of biology, but it arguably has a better chance than either method on its own.*

*- Nature Methods editorial [1]*

*'Omics' data should provide more than an ad-hoc mechanism for generating hypotheses. The quantitative evidence contained in the correlative structure of large omics datasets should directly aid the validation of biological hypotheses.*

An experimental biologist wants to test the existence of a collection of biological effects

$$\Theta = \{\theta_1, \ldots, \theta_m\}$$

A statistician has access to a massive genomics data corpus

$$T = \{T_1, \ldots, T_K\}$$

Can the statistician use $T$ to help the biologist resolve $\Theta$?

# Frequentist, assisted by Bayes
## $p$-values

1. Identify a biological effect of interest: $\theta_j$
2. Conduct an experiment to obtain samples $Y_{1j}, \ldots, Y_{n_j j} \overset{iid}{\sim} N(\theta_j, \sigma^2)$
3. Find estimator $\hat{\sigma}^2$ such that $\nu \hat{\sigma}^2 / \sigma^2 \sim \chi_\nu^2$
4. Calculate $T_j = \bar{Y}_j / \sqrt{\hat{\sigma}^2 / n_j}$
5. Decide whether $\theta_j = 0$ or $\theta_j \neq 0$ based on $T_j$

Supposing $|T_j|$ is the test statistic, Hoff (2019) [2] notes that the formula for the classical *p*-value for the two-sided test may be written as

$$p_j = 1 - \left| F_\nu(T_j) - F_\nu(-T_j) \right| \tag{1}$$

where $F_\nu$ is the cumulative distribution function of the $t_\nu$ distribution

- Recall that $p_j \sim U(0, 1)$ under the null hypothesis $\theta_j = 0$
- Type I error can be controlled at level $\alpha$ by rejecting when $p_j < \alpha$

For any $b_j$ statistically independent of $T_j$, Hoff (2019) [2] also notes that another quantity is uniformly distributed under $\theta_j = 0$, namely:

$$p_j^{\text{FAB}} = 1 - \left| F_\nu \left( T_j + b_j \right) - F_\nu \left( -T_j \right) \right| \tag{2}$$

### Important

- Valid frequentist $p$-value
- Corresponding test is **more powerful** than the classical test if $b_j$ and $\theta_j$ have the same sign
- Approaches $p$-value from a one-sided (oracle) test as $b_j \to \pm\infty$ for $\text{sign}(\theta_j) = \pm 1$

## Bayes optimal choice of $b_j$

If prior knowledge about $\theta_j$ represented by

$$\theta_j \sim N(m_j, s_j^2) \tag{3}$$

then FAB $p$-value corresponding to Bayes-optimal level-$\alpha$ test has

$$b_j^{OPT} = 2m_j\sigma/s_j^2 \tag{4}$$

If $\sigma^2$ must be estimated from the data, can use plug-in estimator:

$$b_j := 2m_j\tilde{\sigma}/s_j^2 \tag{5}$$

Q: How do we get information $m_j, s_j^2$ that is independent of $T_j$?

A: Model relationships among biological effects $\{\theta_j : j = 1, \ldots, m\}$

Suppose we collect auxiliary information about the biological effects $\{\theta_j : j = 1, \ldots, m\}$ into the rows of a matrix $\mathsf{X}$. Then let

$$
\begin{aligned}
\bar{\mathsf{Y}}|\boldsymbol{\theta} &\sim N_m\left(\boldsymbol{\theta}, \operatorname{diag}(\sigma^2/n_j)\right) \\
\boldsymbol{\theta}|\boldsymbol{\beta} &\sim N_m(\mathsf{X}\boldsymbol{\beta}, \tau^2 \mathsf{I}_m) \\
\boldsymbol{\beta} &\sim N_p(\mathsf{0}, \psi^2 \mathsf{I}_p)
\end{aligned}
\tag{6}
$$

Under this combined sampling and *linking* model, $\bar{\mathsf{Y}}_{-j}$ gives us indirect information about $T_j$ via estimators like

$$
\tilde{m}_j = \operatorname{E}(\theta_j|\bar{\mathsf{Y}}_{-j}) \quad \tilde{s}_j^2 = \operatorname{Var}(\theta_j|\bar{\mathsf{Y}}_{-j})
\tag{7}
$$

The combined sampling and linking model states that $\bar{Y}$ is marginally normally distributed with mean $0$ and covariance

$$\psi^2 XX^T + \tau^2 I_m + \sigma^2 \begin{bmatrix} 1/n_1 & & \\ & \ddots & \\ & & 1/n_m \end{bmatrix} \tag{8}$$

The variation in $\bar{Y}$ is decomposed into

1. Variation along the principal directions of $X$
2. Isotropic variation (i.e. variation not explained by $X$)
3. Measurement error

Conditional on $\sigma^2, \tau^2, \psi^2$, closed form solutions for $\tilde{m}_j, \tilde{s}_j^2$ exist

## Empirical Bayes for variance components

1. $\tilde{\sigma}^2$ can be computed from replicate measurements
2. Given $\tilde{\sigma}^2$, can compute $\tilde{\tau}^2$ and $\tilde{\psi}^2$ by maximizing the marginal likelihood function $p(\bar{\mathbf{Y}}|\tau^2, \psi^2, \tilde{\sigma}^2)$

$$\arg\max_{\tau^2, \psi^2} \left\{ -\sum_{i=1}^m \left[ \log(\psi^2 \lambda_i + \tilde{\sigma}^2/\bar{n} + \tau^2) + \frac{||Q_i^T Y||_2^2}{\psi^2 \lambda_i + \tilde{\sigma}^2/\bar{n} + \tau^2} \right] \right\} \quad (9)$$

Efficient representation in terms of eigenvalues, eigenvectors of $\mathbf{XX}^T$ by making approximation $\tilde{\sigma}^2/n_j \approx \tilde{\sigma}^2/\bar{n}, \forall j$

# Distilled auxiliary information from genomics data

Ideally we would like to have auxiliary information $X$ that is relevant to a wide range of genomics contexts:

- Differential expression analysis
- CRISPR modifier screens
- Drug discovery screens

Experimental techniques and technologies differ, but *genes* and *cancer cell lines* recur

If we had auxiliary features for **genes** and **cancer cell lines**, we could use the FAB framework to boost power of hypothesis tests for any experiment in which these entities appear

## Where to find auxiliary features?

- Recall massive genomics data $T = \{T_1, \ldots, T_K\}$
- Let each $T_k$ be a matrix of measurements for experimental modality $k$; cancer cell lines on the rows, genes on the columns
- Can decompose tensor $T$ into gene, cancer cell line, and experimental modality constituents

**Figure 1:** Low-rank tensor factorization for covariate distillation. Model explored in different forms in [4], [3], [5]

The probability model for tensor entries:

$$T_{lgk}|\mathbf{R}_k, \mathbf{U}_l, \mathbf{V}_g, \tau_k^2 \sim N\left(\mu_k + \mathbf{U}_l^T \mathbf{R}_k \mathbf{V}_g, \tau_k^2\right) \tag{10}$$

**Convenient extensions**

- Probit likelihood for binary data (e.g. mutations)
- "Tobit" likelihood for positive, continuous data (e.g. gene expression)
- MARginalize over missing data within sampling steps

1. Obtain estimators of gene and cancer cell line auxiliary features via tensor probability model and Gibbs sampling: $\hat{U} = \mathrm{E}(U|T)$, $\hat{V} = \mathrm{E}(V|T)$

2. Construct auxiliary information matrix: $X = \hat{V} \otimes \hat{U}$

3. Obtain estimators $\tilde{m}_j = \mathrm{E}\left(\theta_j|\bar{Y}_{-j}\right)$, $\tilde{s}_j^2 = \mathrm{Var}(\theta_j|\bar{Y}_{-j})$ for each biological effect $\theta_j$ using sampling, linking model

4. Set $\tilde{b}_j := 2\tilde{m}_j\tilde{\sigma}/\tilde{s}_j^2$

5. Calculate FAB $p$-value $p_j^{\mathsf{FAB}} = 1 - \left| F_\nu\left(T_j + \tilde{b}_j\right) - F_\nu\left(-T_j\right) \right|$

# Validation by simulation and application to selected studies

Simulate 10,000 null datasets ($\boldsymbol{\theta} = \mathbf{0}$, $m = 250$, $n_j = 6 \; \forall j$)



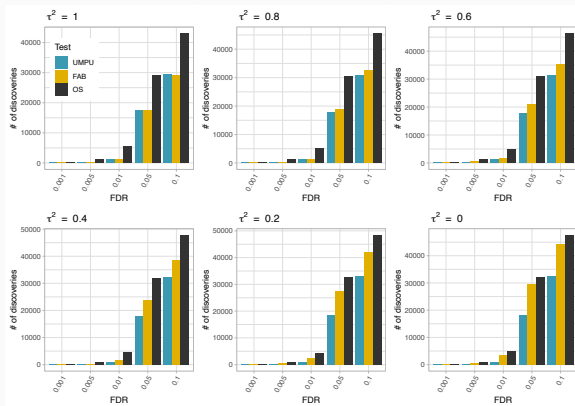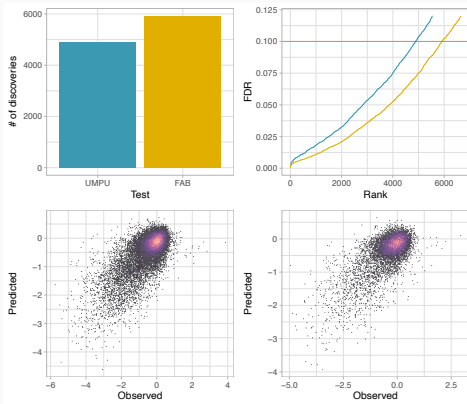Figure 2: Empirical distribution of UMPU and FAB $p$-values under null simulation. Both achieve target FDR up to Monte Carlo error.

**Figure 3:** FAB procedure cleanly interpolates between two-sided test and one-sided oracle test as $\psi^2/\tau^2 \to \infty$.

# Auxiliary information from the Cancer Dependency Map

The tensor **T** has a publicly available incarnation: *depmap.org*

1. RNAseq
2. CRISPR KO
3. RNAi KD
4. Mutation calls

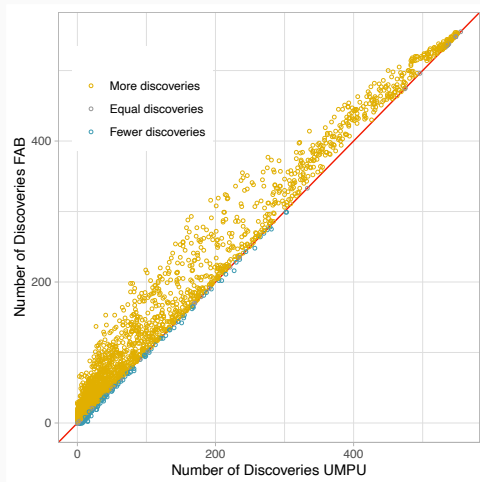500 to 1000+ cancer cell lines, 15,000 to 30,000+ genes in each dataset

**Figure 4:** FAB procedure leads to more discoveries on new experiments using modalities contained in **T**

**Figure 5:** Relationships among cancer cell lines lead to more discoveries for 67% of tested compounds. At least as many discoveries as classical procedure for 87% of tested compounds.
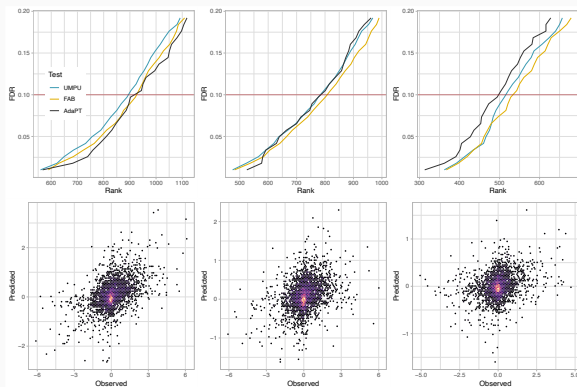
Figure 6: Gene-gene relationships distilled from cancer genomics data contain information relevant to non-cancer study.

## Take home points

1. FAB procedure offers increase in statistical power, which can accumulate to hundreds of additional discoveries in the presence of many hypotheses. Strict type I error and FDR control.

2. Cell line and gene representations derived from tensor model can be used with FAB procedure to improve statistical power in many genomics contexts.

3. Little downside to proposed FAB procedure. Reverts to classical hypothesis testing when $\psi^2 \ll \tau^2$.

Code available at *https://github.com/j-g-b/BTF*

📄 Defining the scientific method.
*Nature Methods*, 6(4):237–237, Apr. 2009.

📄 P. D. Hoff.
Smaller $p$-values via indirect information.
*arXiv:1907.12589 [stat]*, July 2019.
arXiv: 1907.12589.

📄 R. Salakhutdinov and A. Mnih.
Bayesian probabilistic matrix factorization using Markov chain Monte Carlo.
In *Proceedings of the 25th international conference on Machine learning - ICML '08*, pages 880–887, Helsinki, Finland, 2008. ACM Press.

I. Sutskever, J. B. Tenenbaum, and R. R. Salakhutdinov.
**Modelling Relational Data using Bayesian Clustered Tensor Factorization.**
In Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 1821–1828. Curran Associates, Inc., 2009.

J. Ye.
**Generalized Low Rank Approximations of Matrices.**
*Machine Learning*, 61(1-3):167–191, Nov. 2005.

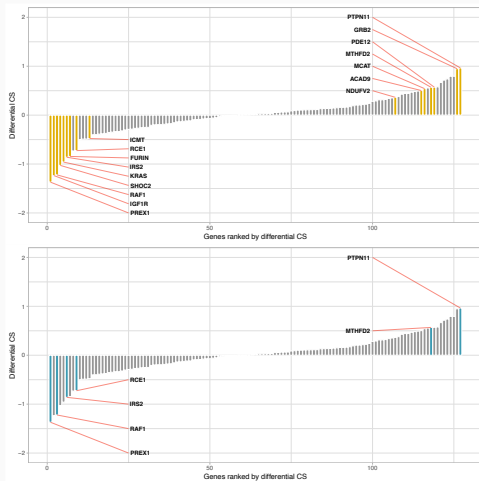Compare the vectorized form of the tensor probability model for the $k^{\text{th}}$ modality

$$\text{vec}(\mathsf{T}_k)\,|\mathsf{R}_k, \mathsf{V}, \mathsf{U}, \tau_k^2 \sim N_{n_V \times n_U}\left((\mathsf{V} \otimes \mathsf{U})\text{vec}(\mathsf{R}_k)\,, \tau_k^2 \mathsf{I}_{n_V \times n_U}\right) \qquad (11)$$

to the linking model for multiple hypotheses

$$\boldsymbol{\theta}|\boldsymbol{\beta}, \mathsf{X}, \tau^2 \sim N_m(\mathsf{X}\boldsymbol{\beta}, \tau^2 \mathsf{I}_m) \qquad (12)$$

Conditional on $\mathsf{X} = \mathsf{V} \otimes \mathsf{U}$, the linking model is equivalent to the tensor model applied to a *new* experimental modality $k'$

Figure 7: FAB procedure can be effective even when $m$ is small compared to $d_V \times d_U$