

CHROMATIN ARCHITECTURE ABERRATIONS IN PROSTATE CANCER AND LEUKEMIA

by

James Hawley

A thesis submitted in conformity with the requirements  
for the degree of Doctor of Philosophy

Graduate Department of Medical Biophysics  
University of Toronto

© Copyright 2021 by James Hawley

# Contents

<b>1 Hedging uncertainty in differential gene expression analyses with James-Stein estimators</b>	<b>1</b>
1.1 Abstract . . . . .	1
1.2 Introduction . . . . .	1
1.3 Results . . . . .	1
1.3.1 Derivation of the James-Stein fold change estimator . . . . .	2
1.3.2 Comparison between the OLS and James-Stein estimators . . . . .	4
1.3.3 Empirical analysis of the James-Stein estimator . . . . .	5
1.4 Discussion . . . . .	6
1.5 Methods . . . . .	6
<b>A Supplementary Material for Chapter 4</b>	<b>8</b>
A.1 Differential expression analysis with Sleuth . . . . .	8
A.2 Statistical moments of the ordinary least squares estimator . . . . .	10
A.3 Statistical moments of the James-Stein estimator . . . . .	11
A.3.1 Expected value of the James-Stein estimator . . . . .	11
A.3.2 Variance of the James-Stein estimator . . . . .	12
<b>Glossary</b>	<b>14</b>
<b>References</b>	<b>16</b>

# Chapter 1

# Hedging uncertainty in differential gene expression analyses with James-Stein estimators

J.R.H., and M.L. conceptualized the study. J.R.H. derived the statistical estimates and designed and conducted all the experiments. Figures were designed by J.R.H. The manuscript was written by J.H., and M.L. M.L. oversaw the study.

## 1.1 Abstract

## 1.2 Introduction

## 1.3 Results

- the two main approaches for reducing error in a model are to reduce the model variance or model bias Figure 1.1
- here we attempt to decrease mean square error (MSE) by simultaneously increasing the bias and decreasing the variance in fold change coefficient estimators
- derivation for the James-Stein (JS) estimator can be found in Section 1.3.1
- Equation for the JS estimator can be seen in Figure 1.1b.

- in theory, this may increase the error of some transcripts, but will decrease MSE for a set of transcripts in aggregate Appendix A.3

### 1.3.1 Derivation of the James-Stein fold change estimator

For a  $p$ -variate normal distribution  $Z \sim \mathcal{N}_p(\mu, \Sigma)$  where  $\mu$  is unknown and  $\Sigma$  is known, the following theorem holds [1]:

**Theorem 1** *The estimator  $\hat{\mu}^{(0)} = Z$ , for any mean  $\mu$ , does not minimize the MSE  $\mathbb{E}[(\mu - \hat{\mu})^2]$  for a single observation of  $Z$  when  $p \geq 3$  and  $\Sigma = I_p$ . Namely, the estimator  $\hat{\mu}^{(JS)} = \left(1 - \frac{b}{a + \|Z\|^2}\right) Z$  has a smaller MSE than  $\hat{\mu}^{(0)}$  for a single observation for sufficiently small  $b$  and large  $a$ .*

This result was generalized to non-singular covariance matrices that were not necessarily the identity matrix (Theorem 2 of [REF 2]):

**Theorem 2** *Let  $Z \sim \mathcal{N}_p(\mu, \Sigma)$ . Let  $\hat{\mu}^{(JS)} = \left(1 - \frac{c}{Z^T \Sigma^{-1} Z}\right) Z$ . If  $p \geq 3$ ,  $\text{Tr}(\Sigma) \geq 2\lambda_L$  (where  $\lambda_L$  is the largest eigenvalue of the covariance matrix,  $\Sigma$ ), and  $0 \leq c \leq 2\left(\frac{\text{Tr}(\Sigma)}{\lambda_L} - 2\right)$ , then  $\hat{\mu}^{(JS)}$  is the estimator for the mean,  $\mu$ , that minimizes the MSE for a single observation of  $Z$ .*

Consider the differential analysis model used in the Sleuth R package [3, 4] for a single transcript,  $s$ , with the simple experimental design of  $n_{nonmut}$  wild-type (WT) samples and 1 mutant sample:

$$D_s | Y_s \sim \mathcal{N}_N(\beta_{s,0} + \mathbb{I}_{mut} \beta_{s,1}, (\sigma_s^2 + \tau_s^2) I_N)$$

For the  $n_{nonmut}$  WT samples, this is equivalent to

$$D_s | Y_s \sim \mathcal{N}_{n_{nonmut}}(\beta_{s,0}, (\sigma_s^2 + \tau_s^2) I_{n_{nonmut}})$$

which can be fit with the same model process that Sleuth employs. For the single mutated sample, this model is

$$D_s | Y_s \sim \mathcal{N}(\beta_{s,0} + \beta_{s,1}, \max\{\hat{\sigma}_s^2, \tilde{\sigma}_s^2\} + \hat{\tau}_s^2) \quad (1.1)$$

The covariance matrix is the same as the mutated samples, but the mean  $\beta_{s,0} + \beta_{s,1}$  is unknown. By reparameterizing this model to consider every transcript in the single mutated sample, Equation (1.1) can be re-written as follows:

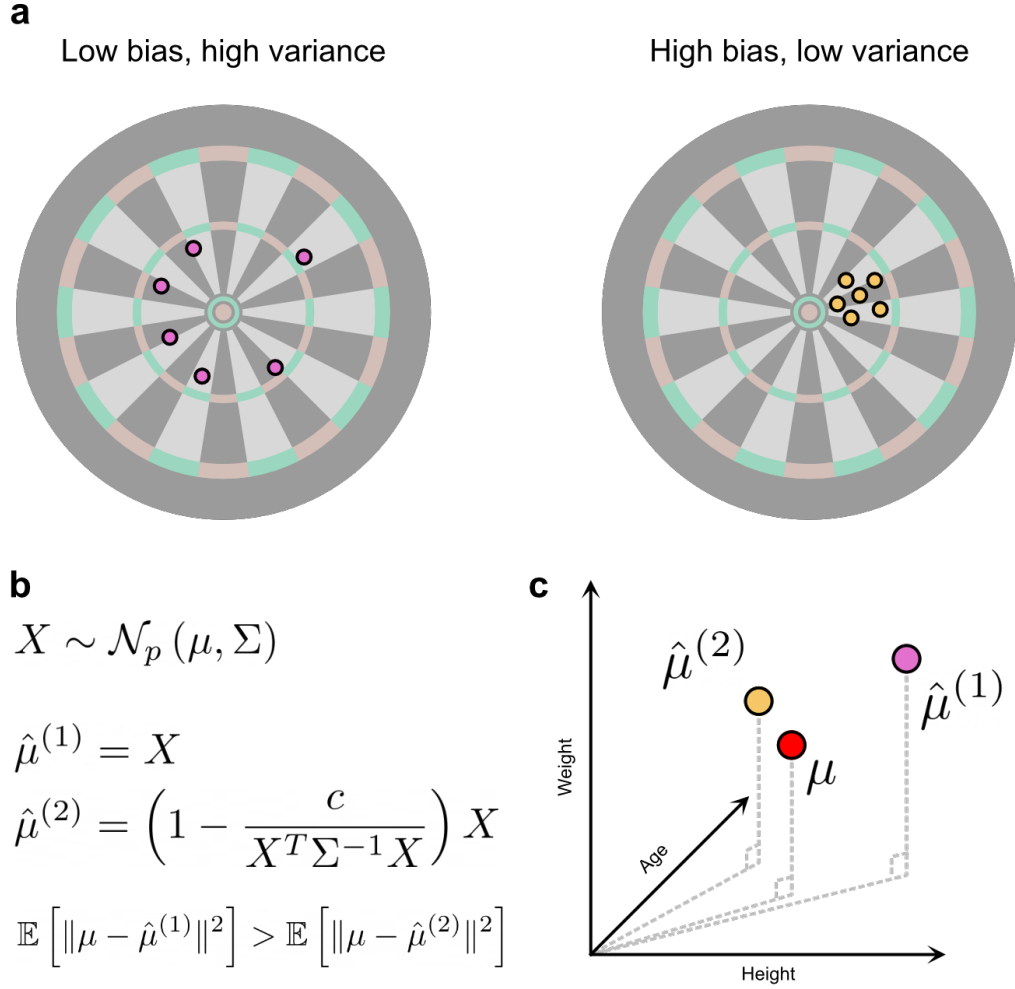


Figure 1.1: **Reducing the bias-variance tradeoff by combining information across multiple features.** **a.** Schematic of the bias-variance tradeoff for assessing model performance. Dartboard on the left shows low bias of the darts (mean is close to the bullseye) but a large variance. Dartboard on the right shows a high bias of the darts (mean is off-centre), but a small variance. **b.** For a  $p$ -variate normal distribution from which a single observation is made, the naive estimator has a higher MSE than the JS estimator, defined as  $\hat{\mu}^{(2)}$ . **c.** An analogy showing how the JS estimators work in theory. Trying to estimate the mean height, weight, and age for the entire population ( $\mu$ ) from a single person will give an estimate that is likely far from the truth ( $\hat{\mu}^{(1)}$ ). Combining information from the three variables together can produce an estimate that is closer to the truth ( $\hat{\mu}^{(2)}$ ).

$$\Delta \sim \mathcal{N}_{|S|}(\mathbf{B}_0 + \mathbf{B}_1, \Sigma) \quad (1.2)$$

where

$$\mathbf{B}_{i,s} = \beta_{s,i} \forall s \in S$$

$$\Sigma = \begin{bmatrix} \max\{\hat{\sigma}_1^2, \tilde{\sigma}_1^2\} + \hat{\tau}_1^2 & & 0 \\ & \ddots & \\ 0 & & \max\{\hat{\sigma}_{|S|}^2, \tilde{\sigma}_{|S|}^2\} + \hat{\tau}_{|S|}^2 \end{bmatrix}$$

We switch from using coefficients  $\beta_{t,i}$  to  $\mathbf{B}_{i,s}$  to avoid confusion, since  $\beta_{t,i} \in \mathbb{R}^p$  (a  $p$ -dimensional vector for each covariate in the design) whereas  $\mathbf{B}_{i,s} \in \mathbb{R}^{|S|}$  (an  $|S|$ -dimensional vector for only a single coefficient over all transcripts in  $S$ ).

Observations of a single mutated sample from this model meet the criteria for the JS estimators. A JS estimator for the unknown fold change,  $\mathbf{B}_1$ , can be constructed.

$$\hat{\mathbf{B}}_1^{(JS)} = \left( 1 - \frac{c}{(\Delta - \hat{\mathbf{B}}_0)^T \Sigma^{-1} (\Delta - \hat{\mathbf{B}}_0)} \right) (\Delta - \hat{\mathbf{B}}_0) \quad (1.3)$$

where  $\hat{\mathbf{B}}_0$  is the estimate obtained from the non-mutated samples for all transcripts  $s \in S$ .

It is straightforward to see that  $\text{Tr}(\Sigma) = \sum_{s \in S} \max\{\hat{\sigma}_s^2, \tilde{\sigma}_s^2\} + \hat{\tau}_s^2$  and that  $\lambda_L = \max_{s \in S} \{\max\{\hat{\sigma}_s^2, \tilde{\sigma}_s^2\} + \hat{\tau}_s^2\}$ .

### 1.3.2 Comparison between the OLS and James-Stein estimators

For a simple experimental design where the mutation status is the only coefficient the ordinary least squares (OLS) estimator is given by:

$$\begin{bmatrix} \hat{\beta}_{s,0}^{(OLS)} \\ \hat{\beta}_{s,1}^{(OLS)} \end{bmatrix} = \hat{\beta}_s^{(OLS)} = (X^T X)^{-1} X^T d_s = \begin{bmatrix} \bar{d}_s^{(nonmut)} \\ d_s^{(mut)} - \bar{d}_s^{(nonmut)} \end{bmatrix}$$

where

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \end{bmatrix} \in \mathbb{R}^{(N+1) \times 2}$$

is the design matrix. Looking closely at the OLS estimator for the mutation coefficient,  $\beta_{s,1}$ , it is clear that it is given by:

$$\hat{\beta}_{s,1}^{(OLS)} = d_s^{(mut)} - \hat{\beta}_{s,0}^{(OLS)} = \delta_s - \hat{\beta}_{0,s} \quad (1.4)$$

which is used directly in the definition of the JS estimator in Equation (1.3). The JS estimator for  $B_1$  can then be expressed simply as:

$$\hat{B}_1^{(JS)} = \left( 1 - \frac{c}{\left( \hat{B}_1^{(OLS)} \right)^T \Sigma^{-1} \hat{B}_1^{(OLS)}} \right) \hat{B}_1^{(OLS)} \quad (1.5)$$

From this definition, it is easy to see that the JS estimate is colinear with the OLS estimate but uniformly shrunk towards 0. For a more general experimental design, the above can be extended. Given an experimental design matrix  $X \in \mathbb{R}^{n \times p}$  where  $n > p$ ,  $\text{rank}(X) = p$  and  $\text{rank}(X^*) = p - 1$  where  $X^* \in \mathbb{R}^{(n-1) \times p}$  is the same design matrix but with one sample removed, a JS estimator for the linear coefficient uniquely specified by the one sample is given by

$$\hat{B}_i^{(JS)} = \left( 1 - \frac{c}{\left( \hat{B}_i^{(OLS)} \right)^T \Sigma^{-1} \hat{B}_i^{(OLS)}} \right) \hat{B}_i^{(OLS)}$$

It can be shown that the JS estimator is biased towards 0 with a smaller variance than the OLS estimator (see Appendix A.3).

### 1.3.3 Empirical analysis of the James-Stein estimator

- using RNA-seq data from a highly replicated yeast knockout (KO) experiment, we compared the statistical inference from differential gene expression analysis when using the OLS and JS estimators
- found that for small sample sizes ( $n = 4$ ), the JS estimators predict the largest number of true positive (TP) and false positive (FP) as well as the smallest number of true negative (TN) and false negative (FN) on average (Student's unpaired  $t$ -test,  $n = 30$ ,  $p =$ )
- for larger sample sizes ( $n \geq 6$ ), the JS estimators continue to identify more TP and FP, as well as fewer TN and FN than the OLS estimator, on average
- these effects lessen with increasing sample size, as expected

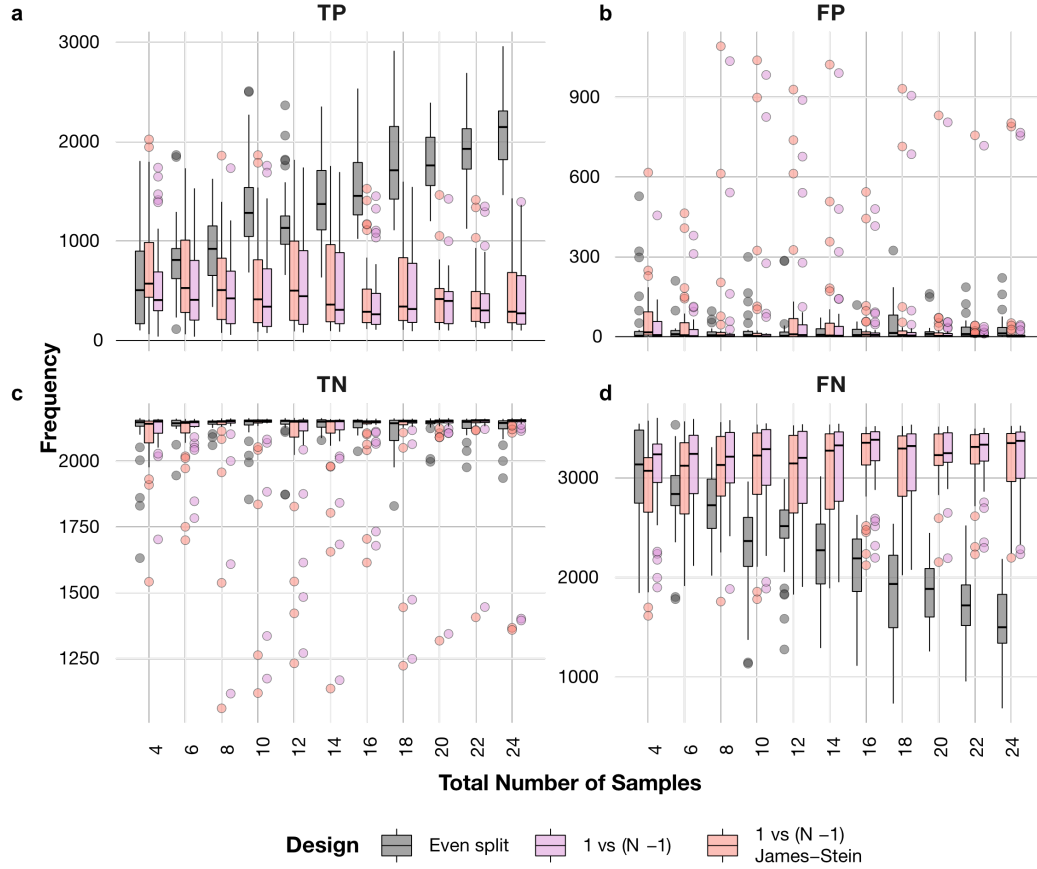


Figure 1.2: **Differential gene expression analysis of the entire yeast transcriptome with differently sized experimental designs.** Simulations ( $n = 30$ ) using randomly selected samples which were then compared to the full dataset of 48  $\Delta$  Snf2 vs 48 WT to calculate TP (a), FP (b), TN (c), and FN (d).

- both the OLS and JS methods with an unbalanced experimental design perform more poorly than the OLS method with a balanced experimental design, as expected
- to investigate where these changes in fold change estimates and  $p$ -values originates from, we investigate representative simulations
- most genes in the  $\Delta$  Snf2 model are under-expressed compared to the WT model Figure 1.3

## 1.4 Discussion

## 1.5 Methods



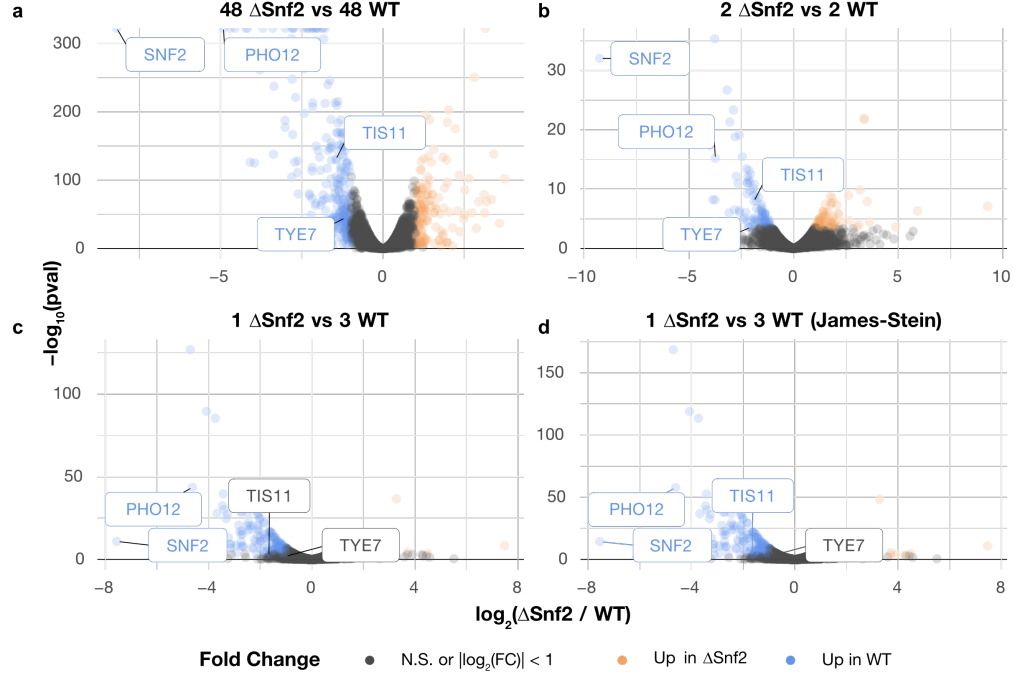


Figure 1.3: **Differential gene expression analysis of  $\Delta$  Snf2 vs WT yeast cells using different sample sizes and experimental designs.** **a.** Volcano plot of differential expression results with OLS estimates in a highly replicated experiment consisting of 48 biological replicates of each condition. **b.** The same analysis as (a) using 4 samples in total, 2  $\Delta$  Snf2 and 2 WT samples. **c.** The same analysis as (c) using 1  $\Delta$  Snf2 and 3 WT. **d.** The same analysis as (c) using the JS method instead of OLS.

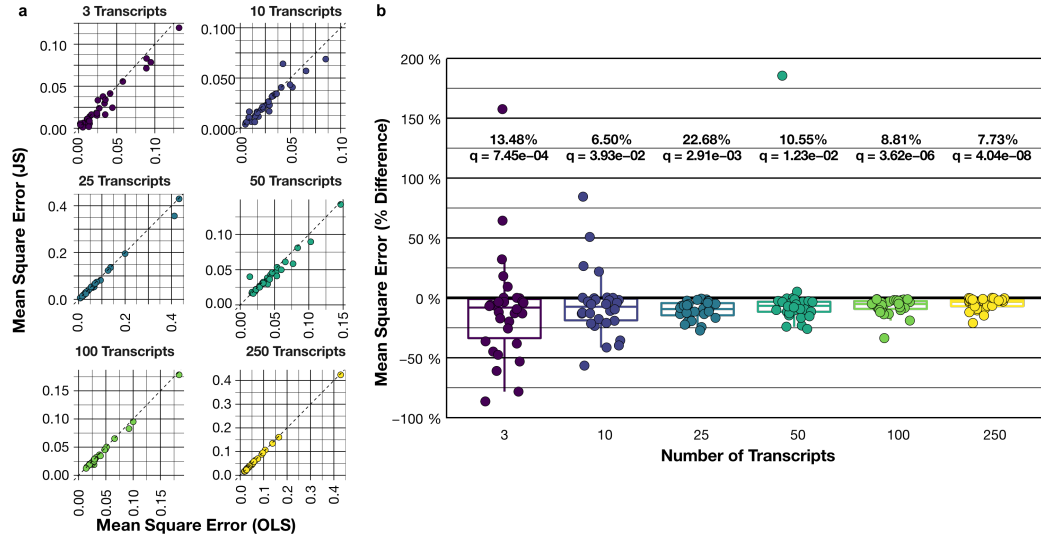


Figure 1.4: **Differential gene expression analysis focusing on a subset of transcripts, not the entire transcriptome.** All experiments use 1  $\Delta$  Snf2 vs 5 WT samples (or vice versa). **a.** Comparison of the MSE of the JS estimates ( $y$ -axis) against the OLS estimates ( $x$ -axis). The total number of transcripts in each comparison is specified above each facet. **b.** Percent different in MSE between the JS and OLS estimates. One-sided, two-sample Student's  $t$ -test,  $n = 30$ , FDR multiple test corrections.

# Appendix A

## Supplementary Material for Chapter 4

### A.1 Differential expression analysis with Sleuth

The differential expression model employed in the Sleuth (v0.30.0) [3, 4] can be described as follows. Consider a set of transcripts,  $S$ , measured in  $N$  samples with an experimental design matrix,  $X \in \mathbb{R}^{N \times p}$ , where  $p$  is the number of covariates considered. Let  $Y_{si}$  be the natural log of the abundance of transcript  $s$  in sample  $i$ . Given the design matrix

$$X = [x_1^T; x_2^T; \dots x_n^T], x_i \in \mathbb{R}^p$$

the abundance of transcripts can be modelled as a generalized linear model (GLM)

$$Y_{si} = x_i^T \beta_s + \epsilon_{si} \tag{A.1}$$

where  $\epsilon_{si} \sim \mathcal{N}(0, \sigma_s^2)$  is the biological noise of transcript  $s$  in sample  $i$  and  $\beta_s \in \mathbb{R}^p$  is the fixed effect of the covariates on the expression of transcript  $s$ .

Due to inferential noise from sequencing, each  $Y_{si}$  are not observed directly, but indirectly through the observed perturbations,  $D_{si}$ . This can be modelled as

$$D_{si}|Y_{si} = Y_{si} + \zeta_{si} \tag{A.2}$$

where  $\zeta_{si} \sim \mathcal{N}(0, \tau_s^2)$  is the inferential noise of transcript  $s$  in sample  $i$ . Both biological and inferential noise for each transcript are independent and identically distributed (IID) and independent of each other. Namely:

$$\mathbb{C}ov[\epsilon_{si}, \epsilon_{rj}] = \sigma_s^2 \delta_{i,j} \delta_{s,r}$$

$$\mathbb{C}ov[\zeta_{si}, \zeta_{rj}] = \tau_s^2 \delta_{i,j} \delta_{s,r}$$

$$\mathbb{C}ov[\epsilon_{si}, \zeta_{rj}] = 0$$

$$\forall s, r \forall i, j$$

The abundances for transcript  $s$  in all  $N$  samples can then modelled as a multivariate normal distribution

$$D_s | Y_s \sim \mathcal{N}_N(X\beta_s, (\sigma_s^2 + \tau_s^2)I_N) \quad (\text{A.3})$$

where  $I_N \in \mathbb{R}^{N \times N}$  is the identity matrix.

The goal of the differential analysis is to estimate the  $|S| \times p$  coefficients in  $B_s \forall s \in S$ , and to determine which coefficients differ significantly from 0. This is achieved through a Wald test or likelihood ratio test after estimating the inferential variance,  $\tau_s^2$ , through bootstrapping and the biological variance,  $\sigma_s^2$ , through dispersion estimation and shrinkage.

The estimator for the differential effect is the OLS estimate:

$$\hat{\beta}_s = (X^T X)^{-1} X^T d_s$$

where  $d_s$  is the observed abundances given by

$$d_{si} = \ln \left( \frac{k_{si}}{\hat{f}_i} + 0.5 \right)$$

$$\hat{f}_i = \text{median}_{s \in S^*} \frac{k_{si}}{\sqrt[N]{\prod_{j=1}^N k_{sj}}}$$

where  $k_{si}$  is the estimated read count from the Kallisto package (v0.46.1) [5] for transcript  $s$  in

sample  $i$  and  $\hat{f}_i$  is the scaling factor for sample  $i$ , calculated from the set of all transcripts that pass initial filtering,  $S^*$ .

## A.2 Statistical moments of the ordinary least squares estimator

As shown in Supplementary Note 2 of [REF 3], the estimator is unbiased, Namely

$$\mathbb{E} \left[ \hat{\beta}_s^{(OLS)} \right] = B_s \quad (\text{A.4})$$

It can also be shown that, for a covariance matrix  $\Sigma$ ,

$$\mathbb{V} \left[ \hat{\beta}_s^{(OLS)} \right] = (X^T X)^{-1} X^T \Sigma X (X^T X)^{-1}$$

In the case where  $\Sigma = (\sigma_s^2 + \tau_s^2) I_N$ , this reduces to

$$\mathbb{V} \left[ \hat{\beta}_s^{(OLS)} \right] = (\sigma_s^2 + \tau_s^2) (X^T X)^{-1}$$

Consider a simple experimental design where the only covariate of interest is the presence of a mutation. Then the design matrix, with the first column being the intercept and the second being the mutation status, looks like so:

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \end{bmatrix} \in \mathbb{R}^{(N+1) \times 2}$$

The variance of the OLS estimator is then

$$\mathbb{V} \left[ \hat{\beta}_s^{(OLS)} \right] = \frac{(\sigma_s^2 + \tau_s^2)}{n_{mut} n_{nonmut}} \begin{bmatrix} n_{mut} & -n_{mut} \\ -n_{mut} & n_{mut} + n_{nonmut} \end{bmatrix}$$

Importantly, the estimate for the coefficient measuring the effect that the presence of the mutation has variance

$$\mathbb{V} \left[ \beta_{s,mut}^{(OLS)} \right] = \frac{(\sigma_s^2 + \tau_s^2)(n_{mut} + n_{nonmut})}{n_{mut} n_{nonmut}}$$

When there is only 1 mutated sample, as per the motivation of this work, this reduces to

$$\mathbb{V} \left[ \beta_{s,mut}^{(OLS)} \right] = \frac{(\sigma_s^2 + \tau_s^2)(1 + n_{nonmut})}{n_{nonmut}} \quad (\text{A.5})$$

## A.3 Statistical moments of the James-Stein estimator

### A.3.1 Expected value of the James-Stein estimator

Due to the non-linear nature of the JS estimator, a Taylor expansion around  $B_1$  can be used to approximate the expectation. Consider:

$$\hat{B}_1^{(JS)} = \left( 1 - \frac{c}{\left( \hat{B}_1^{(OLS)} \right)^T \Sigma^{-1} \hat{B}_1^{(OLS)}} \right) \hat{B}_1^{(OLS)}$$

where

$$\begin{aligned} \hat{B}_1^{(OLS)} &\sim N_{|S|}(B_1, \Sigma) \\ \Sigma_{s,s} &= \left( \frac{n_{nonmut} + 1}{n_{nonmut}} \right) (\sigma_t^2 + \tau_t^2) \\ \Sigma_{s,t} &= 0 \forall t \neq s \end{aligned}$$

Let  $u = \Sigma^{-1/2} \hat{B}_1^{(OLS)}$ . Then

$$\begin{aligned} \mathbb{E} \left[ \hat{B}_1^{(JS)} \right] &= \mathbb{E} \left[ \hat{B}_1^{(OLS)} \right] - c \Sigma^{1/2} \mathbb{E} \left[ \frac{u}{\|u\|^2} \right] \\ &= B_1 - c \Sigma^{1/2} \mathbb{E} \left[ \frac{u}{\|u\|^2} \right] \Sigma^{1/2} \end{aligned}$$

Expanding  $\frac{u}{\|u\|^2}$  around  $a = \Sigma^{-1/2} B_1$  gives:

$$\begin{aligned} \mathbb{E} \left[ \hat{B}_1^{(JS)} \right] &= B_1 - c \Sigma^{1/2} \mathbb{E} \left[ \frac{a}{\|a\|^2} + \left( \frac{1}{\|a\|^2} - \frac{2}{\|a\|^4} a a^T \right) (u - a) + \mathcal{O}(\|u - a\|^2) \right] \\ &= \left( 1 - \frac{c}{B_1^T \Sigma^{-1} B_1} \right) B_1 + \mathcal{O}(\|u - a\|^2) \end{aligned}$$

As long as the number of transcripts being considered,  $|S|$ , is not large, and that the true coefficient of variation is not large (i.e. that  $\|u - a\|^2 \ll \|B_1\|^2$ ), the Taylor approximation is close to

$$\mathbb{E} [\hat{B}_1^{(JS)}] \approx \left(1 - \frac{c}{B_1^T \Sigma^{-1} B_1}\right) B_1 \quad (\text{A.6})$$

Thus the JS estimator is an estimate of  $B_1$  that is biased towards 0.

### A.3.2 Variance of the James-Stein estimator

The MSE of the JS estimator is related to its variance.

$$\mathbb{E} [\|\hat{B}_1^{(JS)} - B_1\|^2] = \sum_{s \in S} \mathbb{E} \left[ \left( \hat{B}_{1,s}^{(JS)} - B_{1,s} \right)^2 \right] = \sum_{s \in S} \mathbb{V} [\hat{B}_{1,s}^{(JS)}]$$

By [REF 2],  $\mathbb{E} [\|\hat{B}_1^{(JS)} - B_1\|^2] \leq \mathbb{E} [\|\hat{B}_1^{(OLS)} - B_1\|^2]$ . However, this does not imply that  $\mathbb{V} [\hat{B}_{1,s}^{(JS)}] \leq \mathbb{V} [\hat{B}_{1,s}^{(OLS)}] \forall s \in S$ . Some transcripts may have larger variances than the OLS estimator, but all transcripts in aggregate will have a smaller MSE. This is still desirable if the goal is to find if there is an effect on any transcripts in the set  $S$ , instead of a particular one within the set.

To calculate the variance for each individual transcript, a similar approach with Taylor expansions can be used, as above.

$$\begin{aligned} & \mathbb{V} [\hat{B}_1^{(JS)}] \\ & \approx \mathbb{E} \left[ \hat{B}_1^{(JS)} \left( \hat{B}_1^{(JS)} \right)^T \right] - \left( 1 - \frac{c}{B_1^T \Sigma^{-1} B_1} \right)^2 B_1 B_1^T \\ & = \Sigma^{1/2} \mathbb{E} \left[ uu^T - \frac{2c}{u^T u} uu^T + \left( \frac{c}{u^T u} \right)^2 uu^T \right] \Sigma^{1/2} - \left( 1 - \frac{c}{B_1^T \Sigma^{-1} B_1} \right)^2 B_1 B_1^T \end{aligned}$$

where, again,  $u = \Sigma^{-1/2} \hat{B}_1^{(OLS)}$ . Expanding about  $a = \Sigma^{-1/2} B_1$  yields:

$$\mathbb{V} [\hat{B}_1^{(JS)}] = \left( 1 - \frac{2c}{B_1^T \Sigma^{-1} B_1} \right) \Sigma - \frac{2c}{(B_1^T \Sigma^{-1} B_1)^2} B_1 B_1^T + \mathcal{O}(\|u - a\|^4)$$

Under similar conditions of the number of transcripts under consideration,  $|S|$ , and  $\|u - a\|^2$ , we then have that

$$\mathbb{V}[\hat{\mathbf{B}}_1^{(JS)}] \approx \left(1 - \frac{2c}{\mathbf{B}_1^T \Sigma^{-1} \mathbf{B}_1}\right) \Sigma - \frac{2c}{(\mathbf{B}_1^T \Sigma^{-1} \mathbf{B}_1)^2} \mathbf{B}_1 \mathbf{B}_1^T \quad (\text{A.7})$$

Since the diagonal elements of  $\frac{2c}{(\mathbf{B}_1^T \Sigma^{-1} \mathbf{B}_1)^2} \mathbf{B}_1 \mathbf{B}_1^T$  are all  $\geq 0$  and  $0 \leq \left(1 - \frac{2c}{\mathbf{B}_1^T \Sigma^{-1} \mathbf{B}_1}\right) \leq 1 \forall c > 0$ , the variance of the JS estimators are smaller than the OLS estimators. This smaller variance allows for more powerful statistical inferences (i.e. when using a Wald test).

Notably, the variance of the JS estimator is a function of both the mean and variance of the transcripts under consideration. This is in contrast to the OLS estimator, which is solely a function of the variance. Additionally, the off-diagonal elements of the matrix  $\mathbf{B}_1 \mathbf{B}_1^T$  imply that the JS fold change estimates are not independent of each other. This, again, contrasts with the OLS estimator, where the diagonal covariance matrix,  $\Sigma$ , implies that the fold change estimates are themselves independent of each other. The effect of this dependence on statistical inference is a function of the variance and true fold change, as can be seen from the  $\frac{2c}{(\mathbf{B}_1^T \Sigma^{-1} \mathbf{B}_1)^2}$  coefficient. While rarely true in practice, this statistical dependence can affect the results of statistical inference, in theory. For most purposes, is not expected to have a large effect on the results of statistical inference.

# Glossary

**3C** chromatin conformation capture

**AR** androgen receptor

**ChIP-seq** chromatin immunoprecipitation sequencing

**CPC-GENE** Canadian Prostate Cancer Genome Network

**crRNA** CRISPR RNA

**CRE** *cis*-regulatory element

**DEPMAP** Cancer Dependency Map

**DHS** DNase I hypersensitive sites

**FN** false negative

**FP** false positive

**FOX** forkhead box

**GLM** generalized linear model

**gRNA** guide RNA

**IID** independent and identically distributed

**JS** James-Stein

**kbp** kilobase

**KO** knockout

**MSE** mean square error



**mCRPC** metastatic castration-resistant prostate cancer

**OLS** ordinary least squares

**mRNA** messenger RNA

**PCa** prostate cancer

**RNAi** RNA interference

**RNA-seq** RNA sequencing

**shRNA** small hairpin RNA

**siRNA** small interfering RNA

**SNV** single nucleotide variants

**SV** structural variant

**TAD** topologically associated domain

**TCGA** The Cancer Genome Atlas

**TN** true negative

**TP** true positive

**TF** transcription factor

**tracrRNA** trans-activating CRISPR RNA

**UTR** untranslated region

**WGS** whole genome sequencing

**WT** wild-type

# References

1. Stein, C. *Inadmissibility of the Usual Estimator for the Mean of a Multivariate Normal Distribution* en. in *Contribution to the Theory of Statistics* **3** (University of California Press, Berkeley, California, USA, Dec. 1956), 197–206. ISBN: 978-0-520-31388-0.
2. Bock, M. E. Minimax Estimators of the Mean of a Multivariate Normal Distribution. en. *The Annals of Statistics* **3**, 209–218. ISSN: 0090-5364 (Jan. 1975).
3. Pimentel, H., Bray, N. L., Puente, S., Melsted, P. & Pachter, L. Differential Analysis of RNA-Seq Incorporating Quantification Uncertainty. en. *Nature Methods* **14**, 687–690. ISSN: 1548-7105 (July 2017).
4. Yi, L., Pimentel, H., Bray, N. L. & Pachter, L. Gene-Level Differential Analysis at Transcript-Level Resolution. *Genome Biology* **19**, 53. ISSN: 1474-760X (Apr. 2018).
5. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-Optimal Probabilistic RNA-Seq Quantification. en. *Nature Biotechnology* **34**, 525–527. ISSN: 1546-1696 (May 2016).