# **Estimating the Number and Effect Sizes of Non-null Hypotheses**



# **Problem Statement**

#### Formal Statement

Let  $\nu_*$  be a distribution on  $\mathbb{R}$ . Let  $f_{\mu}$  be the known parametric distribution of test statistics with effect size  $\mu$ , e.g.  $f_{\mu} = \mathcal{N}(\mu, 1)$ . Draw  $\{\mu_i\}_{i=1}^n \sim \nu_*$  (unobserved) and  $\{X_i\}_{i=1}^n \sim f_{\mu_i}$  (observed).

### Goal

Estimate, for all  $\gamma \in \mathbb{R}$  and without overestimating,

$$\zeta_{\nu_*}(\gamma) := \mathbb{P}_{\nu_*} \left( \mu_i > \gamma \right)$$



# **Our Estimator**

Let

$$\widehat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1} \{ X_i \le t \} \text{ be the empirical CDF} \\ F_{\nu}(t) = \mathbb{P}_{\mu \sim \nu, X \sim f_{\mu}}(X \le t) \text{ be the true CDF under } \nu$$

For any  $\gamma \in \mathbb{R}$ , for a specified probability of overestimation  $\alpha$ , our estimator is given by

$$\widehat{\zeta}_n(\gamma) = \min_{\nu: ||\widehat{F}_n - F_\nu||_{\infty} \le \sqrt{\frac{\log(2/\alpha)}{2n}}} \int_{\gamma}^{\infty} \nu(x) dx.$$
(1)

Jennifer Brennan, Ramya Korlakai Vinayak, Kevin Jamieson {jrb, ramya, jamieson}@cs.washington.edu

# **Theoretical Results**

**Theorem.** For i = 1, 2, ..., n, let  $\mu_i \sim \nu_*$  and  $X_i \sim f_{\mu_i}$  where each draw is iid. Let our simultaneous estimator be given by (1). Then,

$$\mathbb{P}\left(\exists \gamma:\widehat{\zeta}_n(\gamma)>\zeta_{
u_*}(\gamma)
ight)\leq lpha.$$

Furthermore, with probability at least  $1 - \delta$ , for all  $\gamma \in \mathbb{R}$  and  $\varepsilon \in$  $(0, \zeta_{\nu_*}(\gamma)]$  we have  $\zeta_{\nu_*}(\gamma) - \hat{\zeta}_n(\gamma) \leq \varepsilon$  whenever

$$n \geq \frac{\log\left(\frac{4}{\alpha\delta}\right)}{\left(\min_{\nu:\mathbb{P}_{\nu}((\gamma,\infty))\leq\zeta_{\nu_{*}}(\gamma)-\varepsilon}||F_{\nu}-F_{\nu_{*}}||_{\infty}\right)^{2}}.$$



By the **DKW inequality**, the true CDF  $F_{\nu_*}$  is contained in the  $\ell_{\infty}$  ball around  $\widehat{F}_n$ . If we find the  $\nu$  that stays inside this ball but has the *least* mass above  $\gamma$  (the minimizer of (1), shown by  $F_{\nu_{min}}$ ), then with high probability, this is a lower bound on  $\zeta_{\nu_*}(\gamma)$ .

The sample complexity result follows from the DKW inequality and the triangle inequality. We want to bound  $||F_{\nu} - \widehat{F}_n||_{\infty}$  subject to the constraints on  $\nu$  in (1), so we control it using the DKW bound on  $||F_{\nu_*} - \widehat{F}_n||_{\infty}$ , and the constrained quantity  $||F_{\nu} - F_{\nu_*}||_{\infty}$  found in the theorem statement.

PAUL G. ALLEN SCHOOL **OF COMPUTER SCIENCE & ENGINEERING** 



## Synthetic Data

### **Gaussian test statistics**

Let  $\nu_* = (1 - \zeta_*)\delta_0 + \zeta_*\delta_{\gamma_*}$  and  $f_{\mu} = \mathcal{N}(\mu, 1)$ , as when the data are z-scores. The following figure shows the probability of detecting at least half of the discoveries, as a function of both  $\gamma_*$  and  $\zeta_*$ , for a fixed  $n = 10^4$ .



### **Poisson test statistics**

Our estimator also works on non-Gaussian data, as demonstrated here.



**Data** Two z-scores from gene knock-out experiments on each of 13,071

- Our estimator finds evidence of more discoveries than multiple testing can identify, especially at small effect sizes
- MLE suggests there may be even more discoveries – but the MLE frequently

[1] Linhui Hao et al. "Drosophila RNAi screen identifies host genes important for influenza virus replication". In: Nature 454.7206 (2008), p. 890.









