Getting the most bang for your buck: Predicting and maximizing the number of new genomic

variants in future experiments

Lorenzo Masoero¹, Federico Camerlenghi², Stefano Favaro³, Tamara Broderick¹ \blacktriangle ¹MIT, CSAIL,²University of Milano,³University of Torino

- © Modern ML tackles increasingly complex real-world problems
- 8 Often need extremely large amount of data
- 😕 Data-gathering can be expensive, requires careful planning
- Need quantitative framework to know
- 1. How much can we *expect to learn* from new data
- 2. How to collect data to *learn the most* possible

We focus on **genomics** \rightarrow *study rare functional variants* [rare diseases, precision medicine] We provide:

- 1. Fast, accurate prediction of # new genomic variants to be seen in future samples
- 2. Optimal experimental design: how to design follow-up studies in genomics to learn the most

Data sketch & problem description



Goal: design future experiments to find rare functional variants

- Problem: hard, and costly: present in few individuals, need a lot of data; data-gathering is expensive 😕
- Available methods: predict # new variants as we get more data, but...
 - 8 Can be inaccurate (Ionita-Laza et al. [2009], Zou et al. [2016])
 - 8 No guidance on optimal sampling under a fixed budget (Gravel [2014]) 8 Can't accommodate changes in experimental conditions
- Our method: BNP approach for optimal experimental design
 - © Fast, accurate predictions for # new variants
 - © Useful guidance on design of future experiments
 - C Able to accommodate changes in experimental setup: useful for reliable optimal experimental design

Bayesian nonparametric model

- BNP models grow in complexity as more data is collected
- We use the Indian buffet process [IBP]: prior for binary matrices



- $\alpha > 0$: mass parameter: scales total # variants observed
- $\sigma \in [0,1)$: *discount parameter*: controls power law behavior

- $c > \sigma$: concentration parameter: modulates frequencies of widespread counts Derive $U_N^{(M)}$: # new variants in M additional samples given N initial $X_{1,N}$

$$\begin{array}{c|c} U_N^{(M)} & | & X_{1:N} \sim Pois(\underbrace{\alpha \sum_{m=1...M} \frac{(c+\sigma)_{N+m-1}}{(c+1)_{N+m-1}}}_{\text{given N initial samples}}) \\ & \text{ additional samples} & \text{ in M additional samples} & \text{ in M additional samples} \end{array}$$

Use posterior mean $P_N^{(M)}$ as pointwise predictor via **empirical Bayes**:

$$\alpha^*, \sigma^*, c^* = argmin_{\alpha, \sigma, c} \left[\sum_{m=1}^{N-n} \{ P_n^{(m)}(\alpha, c, \sigma) - \left(U_n^{(m)} \mid X_{1:n} \right) \} \right]$$

In practice: $n := \lfloor \frac{2N}{3} \rfloor$

Experiments: prediction and optimal design

- Sanity check: test predictor $P_N^{(M)}$ on synthetic and real data
- Performs comparably to best alternative when experimental conditions don't change between pilot and follow-up studies



new variants depends on samples' quality when experimental conditions change:

- competing methods fail to capture change in samples' quality between pilot and follow up
- BNP approach can successfully adapt and produce useful predictions



Given a pilot study, design follow-up to maximize amount of new information (# new variants) under a fixed budget: Under a fixed budget, solve quality/quantity tradeoff:

- **Quality**: sequencing depth λ -
- **Quantity**: # new samples M

Solve constrained optimization using BNP predictor



