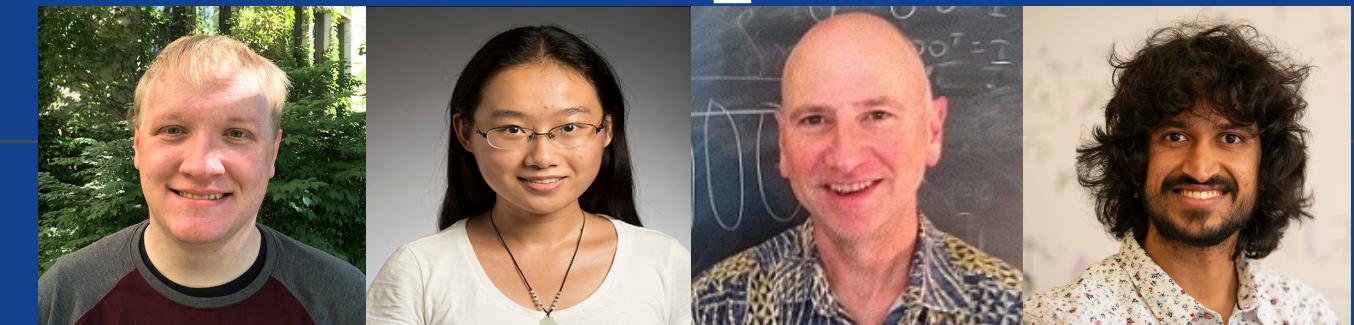


Data fission: splitting a single data point

James Leiner¹, Boyan Duan², Larry Wasserman¹, Aaditya Ramdas¹

¹Department of Statistics and Data Science, Carnegie Mellon University

²Google



Introduction

Suppose we observe data with a known distribution, up to an unknown variable of interest $\theta: X \sim P_X(\theta)$.

We explore decompositions of X into $f(X)$ and $g(X)$ such that:

1. $f(X)$ is not sufficient to reconstruct X by itself
2. There exists a function h such that $h(f(X), g(X)) = X$
3. One of the following two properties holds:
 - $f(X) \perp g(X)$ with known marginal distributions ("strong version")
 - $f(X)$ and $g(X)|f(X)$ have known and tractable distributions ("weak version")

Additive randomization

Consider observing $X \sim N(\theta, \sigma^2)$. We can then draw $Z \sim N(0, \sigma^2)$ and randomize by:

$$f(X) := X + \tau Z \quad X \sim N(\theta, \sigma^2) \quad g(X) := X - \frac{1}{\tau} Z \quad \sim N(\theta, (1 + \tau^2)\sigma^2)$$

As τ increases, more information gets allocated to $f(X)$ and less information gets allocated to $g(X)$.

Conditional randomization

Alternatively, we can use *conditional* rather than *additive* randomization and then choosing a $\tau \in (0, 1)$:

$$X \sim \text{Pois}(\theta) \quad f(X) \sim \text{Bin}(X, \tau) \quad \sim \text{Pois}(\tau\theta) \quad g(X) := X - f(X) \quad \sim \text{Pois}((1 - \tau)\theta)$$

When $\tau = \frac{1}{2}$, $f(X)$ and $g(X)$ are roughly comparable. As $\tau \rightarrow 0$, $g(X) \xrightarrow{d} X$. As $\tau \rightarrow 1$, $f(X) \xrightarrow{d} X$.

Conjugate Prior "Reversal"

Bayesian Inference: $p(\theta|X) \propto p(X|\theta) p(\theta)$

Is $X \sim P_X(\theta)$ conjugate to some other distribution P_Z ? Then draw $Z \sim P_Z(X)$

"Reverse" the prior to accomplish data fission $p(X|Z) \propto p(Z|X) p(X)$

$$f(X) := Z \quad \text{Known marginals for exponential family distributions} \quad g(X) := X \quad \text{Distribution of } g(X)|f(X) \text{ known from Bayesian inference}$$

Need more information for selection? Draw Z_1, \dots, Z_B and let $f(X) := (Z_1, \dots, Z_B)$.

Randomization with conjugacy

The exponential distribution is conjugate prior to the Poisson distribution.

$$X \sim \text{Exp}(\theta)$$

$$f(X) \sim \text{Pois}(\tau X) \quad g(X) := X$$

$$f(X) \sim \text{Geo}\left(\frac{\theta}{\theta + \tau}\right) \quad g(X)|f(X) \sim \text{Gamma}(1 + f(X), \theta + \tau)$$

When $\tau \approx \theta$, the most information is contained in $f(X)$. As $\tau \rightarrow 0$ or $\tau \rightarrow \infty$, $g(X)|f(X) \xrightarrow{d} X$.

Selective Inference

We primarily focus on the use of data fission for (potentially high dimensional) model selection and inference.

Full Dataset (X)

Randomize dataset

Use any procedure to select a model from $f(X)$

Reveal full data at inference step

Use $P(g(X)|f(X))$ for inference

Information splitting

Assuming that $P_X(\theta)$ is known up to unknown parameter of interest θ , data fissions smoothly trades off Fisher information between selection and inference by varying τ .

For "weak version"

$$\mathbb{I}_X(\theta) = \mathbb{I}_{f(X)}(\theta) + E[\mathbb{I}_{g(X)|f(X)}(\theta)]$$

For "strong version"

$$\mathbb{I}_X(\theta) = \mathbb{I}_{f(X)}(\theta) + \mathbb{I}_{g(X)}(\theta)$$

Linear Regression

We assume that y_i is the dependent variable and $x_i \in \mathbb{R}^p$ is a vector of features with corresponding design matrix $X \in \mathbb{R}^{n \times p}$.

$$Y = \mu + \epsilon \text{ with } \epsilon = (\epsilon_1, \dots, \epsilon_n)^T \sim N(0, \sigma^2 I_n) \text{ where } \sigma^2 \text{ is known and } \mu = E[Y|X] \in \mathbb{R}^n \text{ is unknown}$$

Draw $Z \sim N(0, \sigma^2)$

Use $f(Y) := Y + \tau Z$ to select a model $M \subseteq [p]$ and corresponding design matrix X_M

$$\text{Use } g(Y) \text{ to conduct inference on } \hat{\beta}(M) = \operatorname{argmin}_{\tilde{\beta}} \|g(Y) - X_M \tilde{\beta}\|^2 = (X_M^T X_M)^{-1} X_M^T g(Y)$$

Target Parameter:

$$\beta_*(M) = \operatorname{argmin}_{\tilde{\beta}} E\left[\|Y - X_M \tilde{\beta}\|^2\right] = (X_M^T X_M)^{-1} X_M^T \mu$$

Forming confidence intervals

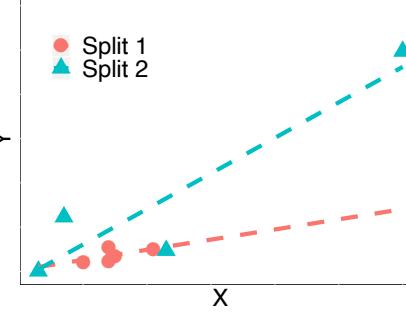
$$\hat{\beta}(M) \sim N(\beta_*(M), (1 + \tau^{-2})\sigma^2 (X_M^T X_M)^{-1})$$

Form $1 - \alpha$ confidence intervals as: $\hat{\beta}^k(M) \pm z_{\alpha/2} \sqrt{(1 + \tau^{-2})\sigma^2 (X_M^T X_M)^{-1}_{kk}}$

If there exists $\hat{\sigma} \xrightarrow{p} \sigma$, calculate this before fissioning the data and use throughout—all guarantees above still hold asymptotically.

Why not data splitting?

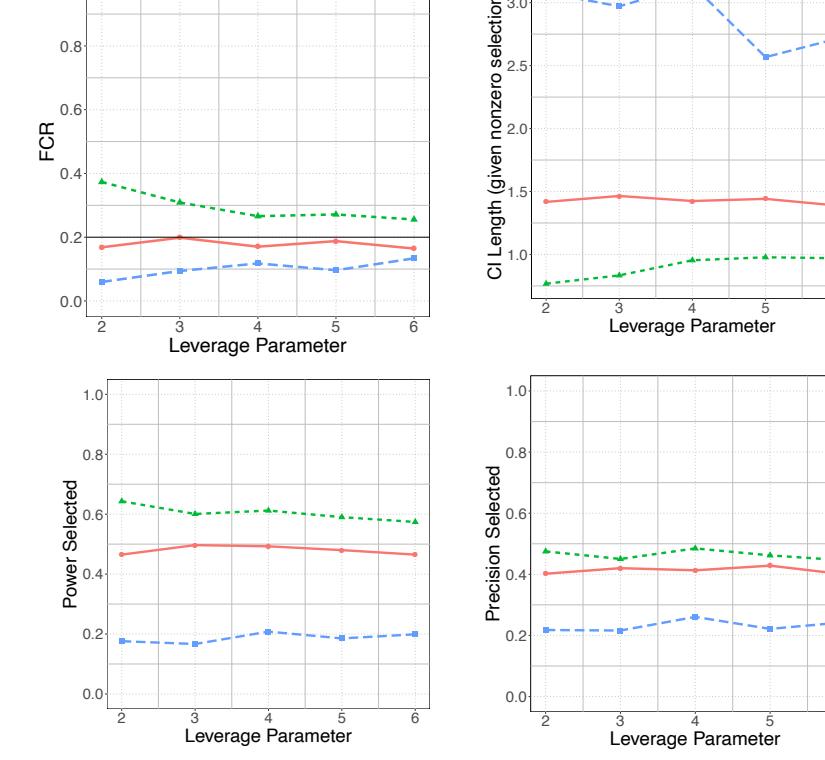
Consider a dataset with low n and a handful of high leverage points.



Data splitting cannot allocate this information to both datasets, so information tradeoff is not "smooth".

Data fission allows the information in that single point to be "split" evenly across $f(X)$ and $g(X)$.

We simulate a small dataset ($n = 20, p = 16$) with a single high leverage point and repeat the data fission procedure for 500 trials, selecting a model with LASSO. We compare data fission along with data splitting, and reusing the full dataset for both selection and inference ("double dipping").



Splitting and fission both control FCR but fission allows for tighter confidence intervals. Double dipping violates FCR control.

Data fission also enables improved power/precision at the selection stage when using LASSO for variable selection.

Interactive Hypothesis Testing

Full Dataset (X)

$f(X)$ $g(X)$

Use for CI construction on rejection set

$$p(p_i) = \min(p_i, 1 - p_i)$$

$$h(p_i) = 2\mathbb{I}\left(p_i > \frac{1}{2}\right)$$

Form rejection set adaptively using BH, AdaPT (Lei and Fithian '18), or STAR (Lei et al. '20) procedures.

Simulation results

We assume $x_i \sim N(\mu_i, 1)$ with $\mu_i = 0$ for nulls and $\mu_i = 2$ for nulls (arranged in a circle). After forming rejection set (\mathcal{R}) from fissioned data (using same Gaussian decomposition as above), we:

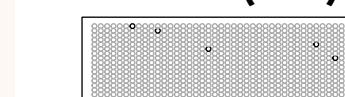
Form $1 - \alpha$ CI:

$$\frac{\sum_{i \in \mathcal{R}} g(y_i)}{|\mathcal{R}|} \pm z_{\alpha/2} \sigma \sqrt{\frac{1 + \tau^{-2}}{|\mathcal{R}|}}$$

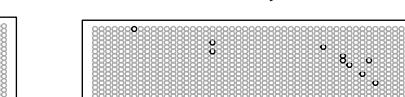
...to cover:

$$\bar{\mu} = \frac{1}{|\mathcal{R}|} \sum_{i \in \mathcal{R}} \mu_i$$

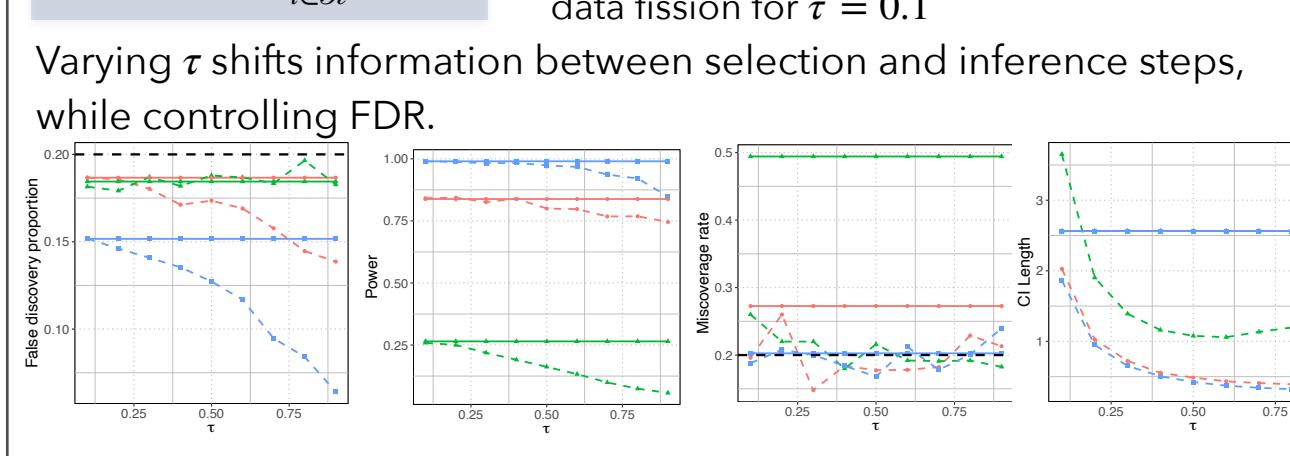
AdaPT (full):



AdaPT (fission):



Example rejection regions, with and without data fission for $\tau = 0.1$



Misspecified GLMs

We assume that y_i follows some distribution in the exponential dispersion family and attempt to model $\mu_i := E[y_i|x_i]$ through covariates $x_i \in \mathbb{R}^p$ under the assumption that $m(\mu_i) = \beta^T x_i$ for some known link function m .

Problem! Even if the distribution of y_i is known, it is unlikely that μ_i is actually a linear combination of the selected covariates for realistic selection rules.

Solution

Assumption: The analyst fissions the data such that $g(y_i) \perp g(y_k)|X, f(Y)$ for all $i \neq k$.

Use $f(Y)$ to select a model $M \subseteq [p]$ which induces a quasi-likelihood function on $g(Y)$, for some working model p :

$$L_n := \sum_{i=1}^n \log p(g(y_i)|\beta, f(Y), X_M),$$

Denote $\hat{\beta}_n(M)$ to be the empirical maximizer of L_n

Target parameter now is the KL minimizer between the working model and true distribution q :

$$\beta_n^*(M) = \operatorname{argmin}_{\beta} D_{KL}\left(\prod_{i=1}^n q(g(y_i)|X, f(Y), X_M) \parallel \prod_{i=1}^n p(g(y_i)|\beta, f(Y), X_M)\right)$$

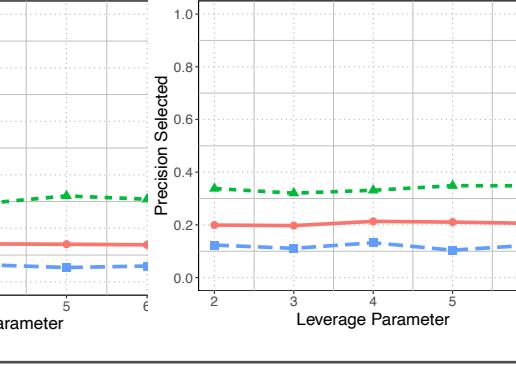
Under mild regularity conditions (Fahrmeir 1990), $\hat{\beta}_n(M) \xrightarrow{p} \beta_n^*(M)$ and asymptotically **conservative** CIs can be constructed (see arxiv preprint for details).

Fission for Poisson regression

We use the same simulation setup as linear regression, but now with a Poisson-distributed response.

Splitting and fission both control FCR but fission allows for tighter confidence intervals.

... and also enables improved power/precision at inference stage when using LASSO for variable selection



Selective Inference for Trend Filtering

We observe a time series $y_t = f_0(t) + \epsilon_t$ where $\epsilon_t \sim N(0, \sigma^2)$ and f_0 is not assumed to belong to any model class.

Trend filtering estimate \hat{f}_0 by constructing a piecewise linear function as:

$$\hat{f}_0 = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{2} \|Y - f\|_2^2 + \lambda \|\|x_{t+1} - x_t\| - (x_t - x_{t-1})\|_1 \quad \text{where } \lambda \text{ is a tuning parameter}$$

Equivalently, we can conceptualize trend filtering in two stages:

Stage 1: Knot Selection

The kink points at which \hat{f}_0 switches direction are called knots

A specific set of knots t_1, \dots, t_r implicitly defines a falling factorial basis which is a set of functions whose discrete derivatives are constant for adjacent design points up to order $k - 1$

Use $f(Y)$ to select basis A

Use $g(Y)$ to conduct inference

Forming confidence intervals

Target parameter is the

$$\mu^*(A) = A(A^T A)^{-1} A^T \mu \quad \text{where } \mu = (f_0(1), \dots, f_0(n))^T$$

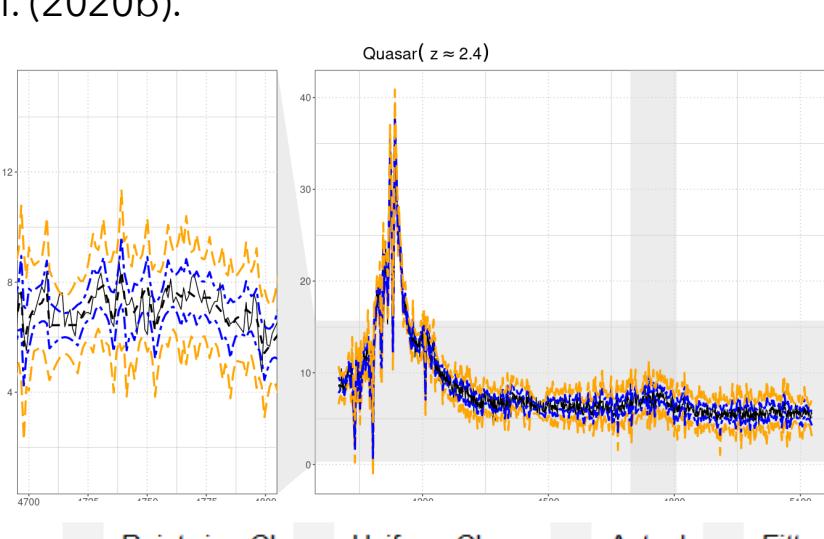
Pointwise CIs:

$$\mathbb{P}(\mu^*(A)_i \notin \text{CI}(\mu^*)_i) \leq \alpha \text{ for all } i.$$

Can be constructed exactly the same as in the linear regression example

Using Koenker (2011) for construction methodology (see arxiv preprint for details).

The above construction will control the FCR (for pointwise CIs) or simultaneous type I error rate (for uniform CIs). To test this procedure, we run it on a real data example. For an astronomical object of interest, we model the coated flux $f(\lambda)$ as a function of wavelengths λ (Politsch et al. 2020b).



Data fission appears to model the underlying trend well, while still allowing for enough information to construct tight confidence intervals