# SLOPE MEETS AMP: DOES SLOPE OUTPERFORM LASSO?

Zhiqi Bu*, Cynthia Rush§, Jason M. Klusowski† and Weijie Su*

*University of Pennsylvania, §Columbia University, †Rutgers University

## SPARSE HIGH-DIM LINEAR REGRESSION



$y$    $X$    $\beta$    $z$

$n \times 1$    $n \times p$       $n \times 1$

$p \times 1$

- Often $p > n$
- $\beta_j \neq 0$ means the $j$th variable is relevant
- Most of entries of $\beta$ are zeros

## SLOPE PROBLEM AND PROPERTIES

Bogdan et al. (2015) proposed SLOPE problem as recovering:

$$\widehat{\beta} := \underset{b}{\operatorname{argmin}} \ \frac{1}{2}\|y - Xb\|_2^2 + J_\lambda(b)$$

where

$$J_\lambda(b) = \underbrace{\lambda_1|b|_{(1)} + \cdots + \lambda_p|b|_{(p)}}_{\text{sorted } \ell_1 \text{ norm}}$$

having $\lambda_1 \geq \ldots \geq \lambda_p \geq 0$, and $|b|_{(1)} \geq \ldots \geq |b|_{(p)}$ are the order statistics.
Clearly, when $\lambda_1 = \cdots = \lambda_p$, SLOPE reduces to LASSO.

- **Estimation:** SLOPE achieves minimax estimation properties under certain random designs without requiring knowledge of the sparsity degree of $\beta$
[Su-Candès '16; Bellec-Lucué-Tsybakov '18]

- **Testing:** SLOPE controls the false discovery rate in the case of independent predictors
[Bogdan-Berg-Sabatti-Su-Candès '15]

- **Optimization:** Since sorted $\ell_1$ norm is a norm, cost remains convex and it can be efficiently solved by using standard methods like proximal gradient descent
(see R package "SLOPE") [Bogdan-Berg-Sabatti-Su-Candès '15]

## COMPUTING SLOPE SOLUTION

Denote $\operatorname{prox}(y; \lambda) := \underset{b}{\operatorname{argmin}} \frac{1}{2}\|y - b\|_2^2 + J_\lambda(b)$.
We may solve the SLOPE problem by

- **Subgradient method:**
$$\beta^{(t+1)} = \beta^{(t)} - s_t \cdot g^{(t)}$$
where $g^{(t)}$ is a subgradient of objective function at $\beta^{(t)}$.

- **ISTA:** Iterative Shrinkage Thresholding Algorithm,
$$\beta^{(t+1)} = \operatorname{prox}\left(\beta^{(t)} + s_t X^\top(y - X\beta^{(t)}); \lambda s_t\right)$$

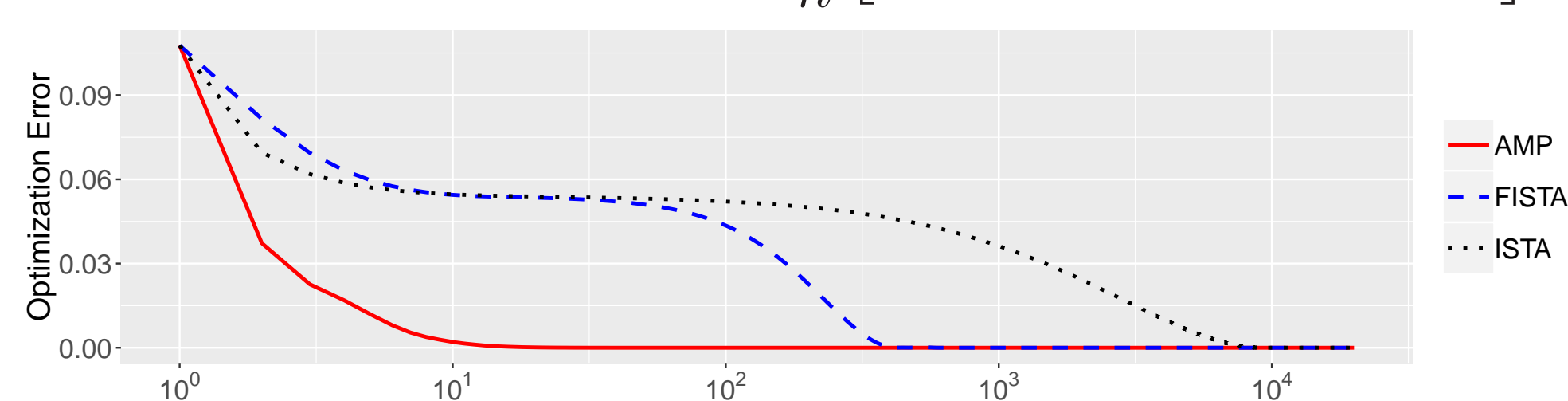- **FISTA:** Fast ISTA, with $s_{t+1} = (1 + \sqrt{1 + 4s_t^2})/2$
$$\beta^{(t+1)} = \operatorname{prox}\left(M^{(t)} + s_t X^\top(y - XM^{(t)}); \lambda s_t\right)$$
$$M^{(t+1)} = \beta^{(t)} + \frac{s_t - 1}{s_{t+1}} \cdot (\beta^{(t)} - \beta^{(t-1)})$$

- **AMP:** Approximate Message Passing [Donoho-Maleki-Montanari '10],
$$\beta^{t+1} = \operatorname{prox}(X^T r^t + \beta^t; \alpha\tau_t),$$
$$r^{t+1} = y - X\beta^{t+1} + \frac{r^t}{n}\left[\nabla \operatorname{prox}(X^\top r^t + \beta^t; \alpha\tau_t)\right].$$



## EXISTING PROCEDURE: LASSO

In 1996, Tibshirani proposed LASSO problem:

$$\min_b \ \frac{1}{2}\|y - Xb\|_2^2 + \underbrace{\lambda\|b\|_1}_{\ell_1 \text{ norm}}$$

$$= \min_b \ \frac{1}{2}\|y - Xb\|_2^2 + \lambda|b_1| + \cdots + \lambda|b_p|$$

- Useful in identifying which $\beta_j \neq 0$ (support recovery, feature selection, etc.)
- LASSO selects at most $n$ variates before it saturates

## PROXIMAL OPERATORS

- For any function $h$, $\operatorname{prox}_h$ is defined as
$$\operatorname{prox}_h(x) := \underset{b}{\operatorname{argmin}} \ \frac{1}{2}\|x - b\|^2 + h(b),$$
and $\nabla \operatorname{prox}_h$ is divergence of the proximal operator

- There exists an algorithm to compute the proximal operator when $h = J_\theta$.[Bogdan-Berg-Sabatti-Su-Candès '15]

- For SLOPE, $\nabla \operatorname{prox}_{J_{\theta_t}}(b) = \|\operatorname{prox}_{J_{\theta_t}}(b)\|_0^*$, where $\|b\|_0^*$ counts the **unique** non-zero magnitudes in $b$.

## STATE EVOLUTION

**SLOPE AMP** (Approximate Message Passing):

$$\beta^{t+1} = \operatorname{prox}(X^T r^t + \beta^t; \theta); r^{t+1} = y - X\beta^{t+1} + \frac{r^t}{n}\|\beta^{(t+1)}\|_0^*$$

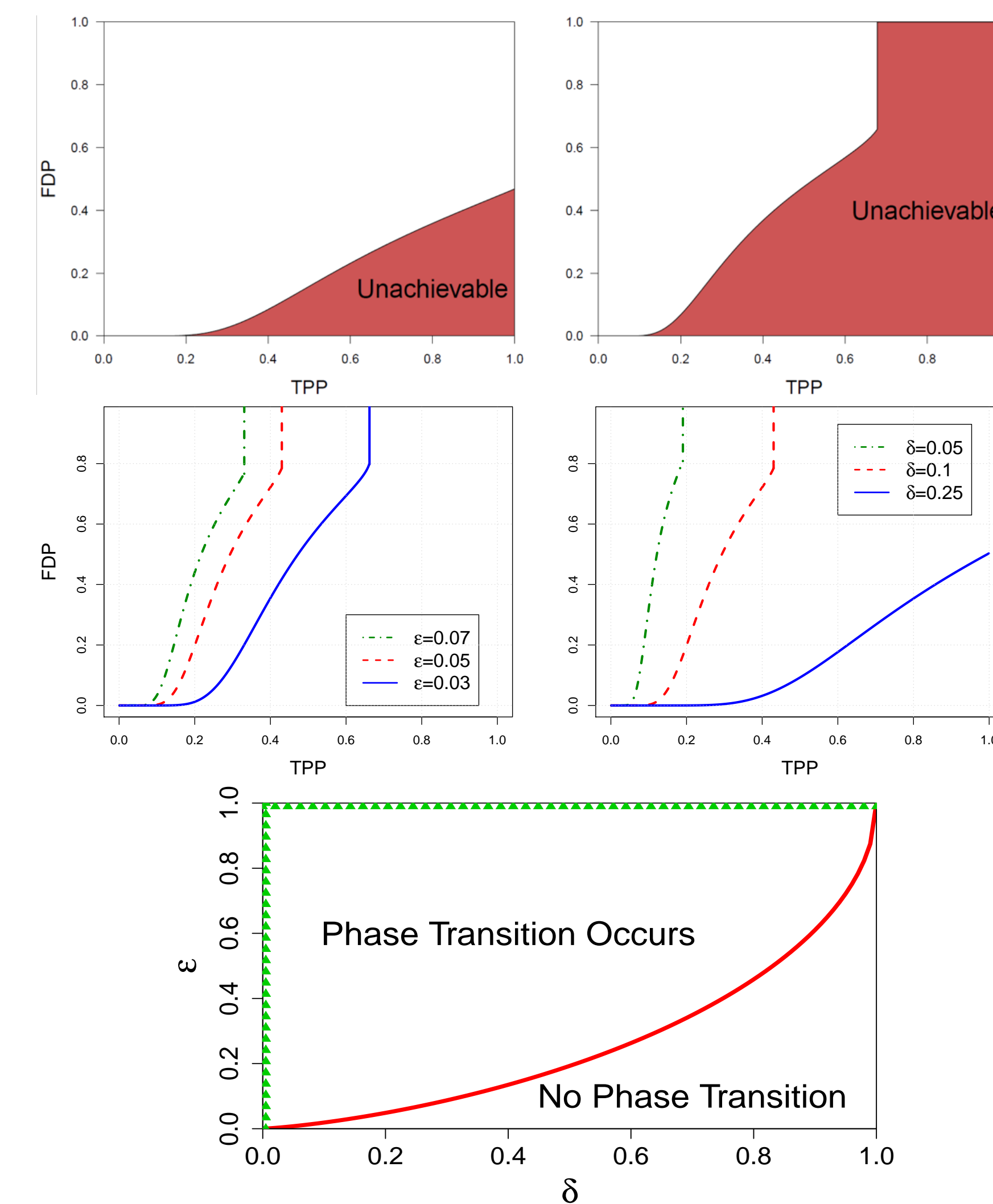> The dynamics of the AMP iterations are tracked by a recursive sequence referred to as the **state evolution**.

**State Evolution** ($n/p \to \delta$):

$$\tau^2 = F(\tau^2, \alpha\tau) := \sigma_z^2 + \lim_p \frac{1}{\delta p}\mathbb{E}\|\operatorname{prox}(B + \tau Z; \alpha\tau) - B\|^2.$$

which is solved iteratively via $\tau_{t+1}^2 = F(\tau_t^2, \alpha\tau_t)$.

## CALIBRATION

> For every iteration $t$, assign $\theta_t = \alpha\tau_t$, where $\alpha$ is a vector in the same direction as $\lambda$.

**Theorem 1** *Under conditions on $\Lambda$, the state evolution recursion with calibration defined above, has a **unique fixed point** to which the convergence monotonic in $t$, for **any** initial condition.*

**Calibration between $\lambda$ and $\alpha$:**
$$\lambda = \alpha\tau_*\left(1 - \lim_p \frac{1}{\delta p}\mathbb{E}\|\operatorname{prox}_{J_{\alpha\tau_*}}(B + \tau_* Z)\|_0^*\right).$$

## MAIN RESULTS OF SLOPE AMP

**Theorem 2** *Under some assumptions,*
$$\lim_{p \to \infty} \frac{1}{p}\|\widehat{\beta} - \beta^t\|^2 = c_t, \quad where \quad \lim_{t \to \infty} c_t = 0.$$

**Theorem 3** *Under some assumptions, for any uniformly pseudo-Lipschitz sequence of functions $\psi_p$ and for $Z \sim \mathcal{N}(0, \mathbb{I}_p)$,*
$$\lim_p \psi_p(\widehat{\beta}, \beta) = \lim_t \lim_p \mathbb{E}_Z[\psi_p(\operatorname{prox}_{J_{\alpha\tau_t}}(\beta + \tau_t Z), \beta)].$$

**Corollary 3.1** *Under some assumptions,*
$$\lim_p \frac{1}{p}\|\widehat{\beta} - \beta\|^2 = \delta(\tau_*^2 - \sigma_z^2).$$

## INFERENCE: TPP, FDP & MSE

For large enough $\epsilon := |\operatorname{supp}(\beta)|/p$ or small enough $\delta := n/p$, LASSO suffers from Donoho-Tanner phase transition: TPP is bounded away from 1 [Donoho-Tanner '09; Su-Bogdan-Candès '17 (image source)].



However, SLOPE overcomes the phase transition. Specifically we can charaterize one of the SLOPE path as a Mobius transformation: for TPP= $u$, the minimum FDP is at most $\frac{au+b}{cu+d}$ for some constants $a, b, c, d$.
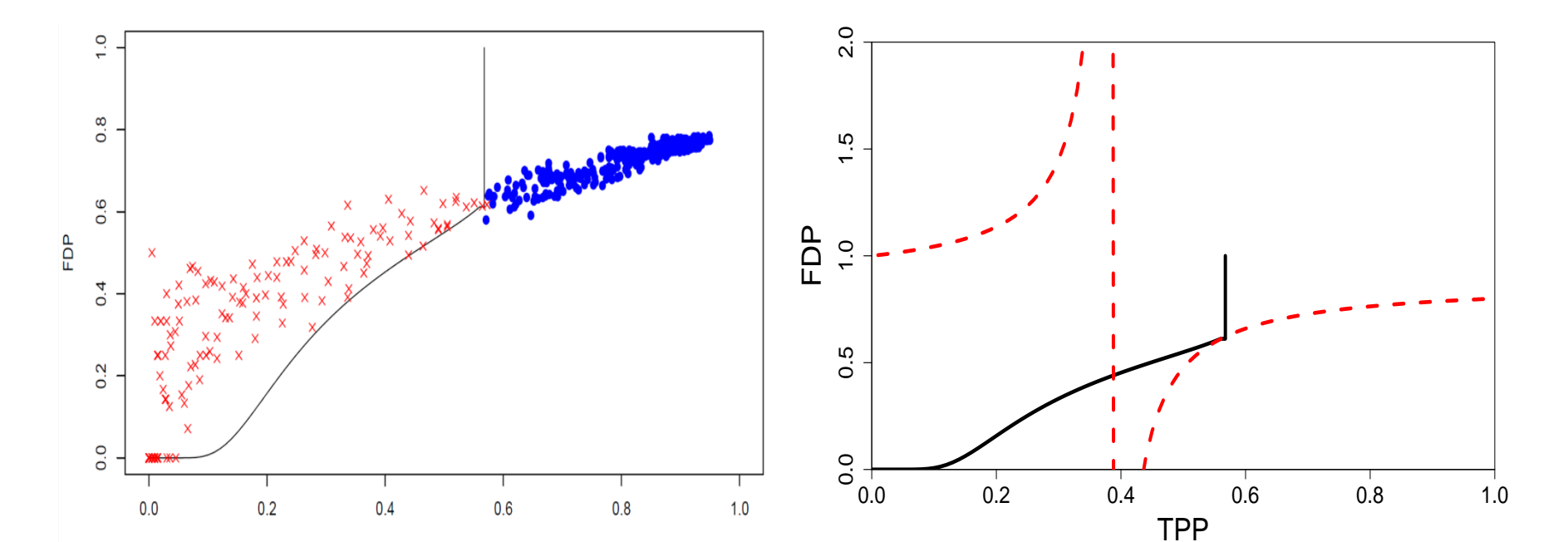


**Figure 1:** Red dot: LASSO; Blue dot: SLOPE; Black solid line: LASSO trade-off; Red dashed line: SLOPE trade-off.

In addition, fixing the signal prior and under some assumptions, we show that switching from LASSO to SLOPE gives better paths in the sense of achieving smaller FDP, larger TPP and smaller mean squared error at the same time.



## CHALLENGES

**Framework:**    $\widehat{\beta} \approx \beta^t \approx \operatorname{prox}_{J_{\alpha\tau_*}}(\beta + \tau_* Z)$

Although AMP for LASSO is well-studied [Bayati-Montanari '15], applying AMP to SLOPE is challenging because the proximal operator of sorted $\ell_1$ norm is **non-separable**.

[Berthier-Montanari-Nguyen '17] showed for general non-separable functions
$$\beta^t \approx \operatorname{prox}(\beta + \tau Z)$$

The main challenge is to show AMP iterate
$$\widehat{\beta} \approx \beta^t$$