

# Outline for Week 8 Presentation on Fast Rates and Sparse Linear Prediction

Tomas Vaškevičius

## Setup

Consider a fixed-design linear regression model

$$\mathbf{y} = \mathbf{X}\mathbf{w}^* + \boldsymbol{\xi},$$

where  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $\mathbf{w}^* \in \mathbb{R}^d$ ,  $\mathbf{y} \in \mathbb{R}^n$  and  $\boldsymbol{\xi} \sim N(0, \sigma^2 I_{n \times n})$ . We assume that the following is true:

1. The  $\ell_2$  norms of columns of  $\mathbf{X}/\sqrt{n}$  are bounded by 1;
2. The “true” parameter  $\mathbf{w}^* \in \mathbb{R}^d$  is  $s$ -sparse, meaning that only  $s$  entries of  $\mathbf{w}^*$  are non-zero.

We measure an estimator’s  $\hat{\mathbf{w}} = \hat{\mathbf{w}}(\mathbf{X}, \mathbf{y})$  statistical performance via its *in-sample prediction error* defined as

$$\mathcal{E}(\hat{\mathbf{w}}) = \mathbf{E}_{\boldsymbol{\xi}} \left[ \frac{1}{n} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{X}\mathbf{w}^*\|_2^2 \right].$$

## The Problem

Under the two assumptions stated above, the  $\ell_0$  norm constrained estimator

$$\hat{\mathbf{w}}_{\ell_0} \in \operatorname{argmin}_{\{\mathbf{w} \in \mathbb{R}^d: \|\mathbf{w}\|_0 \leq s\}} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

is known to attain the minimax optimal rate (see [Raskutti, Wainwright, and Yu, 2011])

$$\mathcal{E}(\hat{\mathbf{w}}) \lesssim \frac{s\sigma^2 \log(d/s)}{n},$$

where the notation  $\lesssim$  hides multiplicative universal constants independent of problem parameters.

The computational intractability of the  $\hat{\mathbf{w}}_{\ell_0}$  estimator has motivated the development of computationally-efficient estimators, most notably the lasso, which replaces the  $\ell_0$  constraint via an  $\ell_1$  penalty term as follows:

$$\hat{\mathbf{w}}_{\ell_1}^\lambda \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1.$$

The above problem is convex and from the computational point of view it can be solved efficiently. However, the lasso is only known to satisfy the fast rate  $O(1/n)$  bound on the in-sample prediction error  $\mathcal{E}(\hat{\mathbf{w}}_{\ell_1}^\lambda)$  provided that the design matrix  $\mathbf{X}$  satisfies additional regularity assumptions. If  $\mathbf{X}$  is *only* assumed to satisfy the normalized columns assumption, several authors have proved “slow rate”  $\Omega(1/\sqrt{n})$  lower-bounds on  $\mathcal{E}(\hat{\mathbf{w}}_{\ell_1}^\lambda)$ , thus establishing *statistical* sub-optimality of the lasso in comparison with the  $\ell_0$  estimator  $\hat{\mathbf{w}}_{\ell_0}$  (see, e.g., [Candès and Plan, 2009, Foygel and Srebro, 2011, Dalalyan, Hebiri, and Lederer, 2017]). More broadly, Zhang, Wainwright, and Jordan [2014] show that the “slow rate” is intrinsic for any polynomial time algorithm that is constrained to output a *sparse* output vector  $\hat{\mathbf{w}}$ .

## Main Result

It remains unknown whether there exists a polynomial-time algorithm that attains the minimax-optimal “fast rate” on the in-sample prediction error. In the reading group meeting, we are going to review a result due to Zhang, Wainwright, and Jordan [2017], which establishes that the “slow rate” is unavoidable to a large class of M-estimators  $\hat{\mathbf{w}}_\rho^\lambda$  defined as follows:

$$\hat{\mathbf{w}}_\rho^\lambda \in \mathcal{W}_\rho^\lambda = \left\{ \mathbf{w} \in \mathbb{R}^d : \mathbf{w} \text{ is a local minimum of the function } \mathbf{w} \mapsto \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \rho(\mathbf{w}) \right\}.$$

The penalty  $\rho$  is assumed to satisfy the following conditions:

1. The function  $\rho$  is coordinate-separable: i.e., there exist univariate functions  $\rho_i$  for  $i = 1, \dots, d$  such that  $\rho(\mathbf{w}) = \sum_{i=1}^d \rho_i(w_i)$ ;
2. For each  $i$ , we have  $\rho_i(0) = 0$  and for all  $w \in \mathbb{R}$  we have  $\rho_i(w) = \rho_i(-w)$ ;
3. Each function  $\rho_i$  is non-decreasing on  $[0, \infty)$ .

We remark that  $\rho$  is not assumed to be a convex. The below result contains a slightly simplified statement of the one presented in [Zhang, Wainwright, and Jordan, 2017].

**Theorem** (Theorem 1 in [Zhang, Wainwright, and Jordan, 2017]). *For any  $d \geq n$  and large enough  $n$ , there exists a design matrix  $\mathbf{X}$  (that satisfies the column-normalization condition) such that for any penalty  $\rho$  that satisfies the three conditions above, there exists a 2-sparse vector  $\mathbf{w}^\star = \mathbf{w}^\star(\rho)$ , such that the following holds:*

$$\mathbf{E}_\xi \left[ \inf_{\lambda \geq 0} \sup_{\mathbf{w} \in \mathcal{W}_\rho^\lambda} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w}^\star\|_2^2 \right] \geq \Omega(1/\sqrt{n}).$$

## Presentation Outline

Below is the plan for the presentation:

1. First, we are going to review some basic properties of the lasso. In particular, we will review the assumption-free “slow rate” upper-bound on  $\mathcal{E}(\hat{\mathbf{w}}_{\ell_1}^\lambda)$  and the “fast rate” upper-bound that holds under the restricted eigenvalue condition.
2. We will then construct a design matrix  $\mathbf{X}$  for which the lasso (with the classical regularization parameter  $\lambda \asymp \sigma\sqrt{\log d}/\sqrt{n}$ ) incurs a suboptimal in-sample prediction error. This part of the presentation will be based on Section 2 in [Candès and Plan, 2009], which already contains the key ideas used in the proof of the Theorem of Zhang, Wainwright, and Jordan [2017].
3. Using the proof presented in [Zhang, Wainwright, and Jordan, 2017], we will extend the above sub-optimality result for the lasso for all scales of the parameter  $\lambda$ .
4. Finally, if time permits, we will discuss how the proof for the lasso extends to general penalties  $\rho$  satisfying the three conditions stated above.

## References

- Emmanuel J Candès and Yaniv Plan. Near-ideal model selection by  $\ell_1$  minimization. *The Annals of Statistics*, 37(5A):2145–2177, 2009.
- Arnak S Dalalyan, Mohamed Hebiri, and Johannes Lederer. On the prediction performance of the lasso. *Bernoulli*, 23(1):552–581, 2017.
- Rina Foygel and Nathan Srebro. Fast-rate and optimistic-rate error bounds for  $\ell_1$ -regularized regression. *arXiv preprint arXiv:1108.0373*, 2011.
- Garvesh Raskutti, Martin J Wainwright, and Bin Yu. Minimax rates of estimation for high-dimensional linear regression over  $\ell_q$ -balls. *IEEE transactions on information theory*, 57(10):6976–6994, 2011.
- Yuchen Zhang, Martin J Wainwright, and Michael I Jordan. Lower bounds on the performance of polynomial-time algorithms for sparse linear regression. In *Conference on Learning Theory*, pages 921–948, 2014.
- Yuchen Zhang, Martin J Wainwright, and Michael I Jordan. Optimal prediction for sparse linear models? lower bounds for coordinate-separable m-estimators. *Electronic Journal of Statistics*, 11(1):752–799, 2017.