

A Two-stage Approach for Learning a Sparse Model with Sharp Excess Risk Analysis

Zhe Li*, Tianbao Yang*, Lijun Zhang[‡], Rong Jin[†]

*The University of Iowa, [‡]Nanjing University, [†]Alibaba Group

Main contribution

- Design a two-stage algorithm to learn a *sparse* linear model
- Reduce the order of the sample complexity from $O(1/\epsilon^2)$ to $O(1/\epsilon)$
- Reduce the constant in $O(1/\epsilon)$ for sparsity by exploring the distribution dependent sampling

Problem

- Let $x \in \mathbb{R}^d$ and $y \in \mathbb{R}$ denote an input and output pair
- Let w_* be an optimal model that minimizes the expected error

$$w_* = \arg \min_{\|w\|_1 \leq B} \frac{1}{2} E_{\mathcal{P}}[(w^T x - y)^2]$$

- **Key Problem:** w_* is not necessarily sparse
- **The goal:** to learn a *sparse* model w to achieve small excess risk

$$ER(w, w_*) = E_{\mathcal{P}}[(w^T x - y)^2] - E_{\mathcal{P}}[(w_*^T x - y)^2] \leq \epsilon$$

Three central questions

- Q: **How to learn** such a *sparse* mode with excess risk less than ϵ ?
A: Given by the following two-stage Algorithms.
- Q: What is the **sample complexity** in order to guarantee a small excess risk?
A: Given by theorem 1.
- Q: What is the **support complexity of w** to suffice for ϵ excess risk?
A: Given by theorem 2.

Related work

- Firstly, minimizes the objective function under ℓ_1 constraint and then use a randomized sparsification approach to find a sparse model [1]
- The first stage algorithm is based on EPOCH-Gradient Descent [2]

[1]: Shalev-Shartz, Shai, Srebro, Nathan, and Zhang, Trading accuracy for sparsity in optimization problems with sparsity constraints.
[2]: Hazan, Kale, Beyond the regret minimization barrier: an optimal algorithm for stochastic strongly-convex optimization.

Two-stage approach to learning a sparse model

Algorithm 1 Stochastic Optimization for Sparse Learning

- 1: **Input:** the total number of iterations T and η_1, ρ_1, T_1 .
- 2: **Initialization:** $\mathbf{w}_1^1 = 0$ and $k = 1$.
- 3: **while** $\sum_{i=1}^k T_i \leq T$ **do**
- 4: **for** $t = 1, \dots, T_k$ **do**
- 5: Obtain a sample denoted by (\mathbf{x}_t^k, y_t^k)
- 6: Compute $\mathbf{w}_{t+1}^k = \Pi_{\|\mathbf{w}\|_1 \leq B, \|\mathbf{w} - \mathbf{w}_t^k\|_2 \leq \rho_k} [\mathbf{w}_t^k - \eta_k \nabla \ell(\mathbf{w}_t^k \cdot \mathbf{x}_t^k, y_t^k)]$
- 7: **end for**
- 8: **Update** $T_{k+1} = 2T_k, \eta_{k+1} = \eta_k/2, \rho_{k+1} = \rho_k/\sqrt{2}$ and $\mathbf{w}_1^{k+1} = \sum_{t=1}^{T_k} \mathbf{w}_t^k / T_k, k = k + 1$
- 9: **end while**
- 10: **Output:** $\hat{\mathbf{w}} = \mathbf{w}_1^{k+1}$

Algorithm 2 Distribution Dependent Randomized Sparsification

- 1: **Input:** $\hat{\mathbf{w}} = (\hat{w}_1, \dots, \hat{w}_d)$ and probabilities p_1, \dots, p_d such that $\sum_{j=1}^d p_j = 1$
- 2: **Initialization:** $\tilde{\mathbf{w}}_0 = 0$.
- 3: **for** $k = 1, \dots, K$ **do**
- 4: sample $i_k \in [d]$ according to $\Pr(i_k = j) = p_j$
- 5: Compute $[\tilde{\mathbf{w}}_k]_{i_k} = [\tilde{\mathbf{w}}_{k-1}]_{i_k} + \frac{\hat{w}_{i_k}}{p_{i_k}}$
- 6: **end for**
- 7: **Output:** $\tilde{\mathbf{w}} = \frac{\tilde{\mathbf{w}}_K}{K}$

Stochastic Optimization —First Stage

Theorem 1: Assume $\|x\|_2^2 \leq R^2$, By running algorithm 1 with $\rho_1 = B, \eta_1 = 1/(2R\sqrt{T_1}), T_1 \geq (8cR + 64R\sqrt{2\log(1/\tilde{\delta})})^2$. In order to have $ER(w, w_*) \leq \epsilon$ with a high probability $1 - \delta$ over $\{(x_t^k, y_t^k)\}$, it suffice to have

$$T = \frac{cB^2T_1}{\epsilon}$$

where $\tilde{\delta} = \frac{\delta}{m}, m = \lfloor \log(cB^2/(2\epsilon) + 1) \rfloor$ and $c = \max(\kappa, 1)$.

Remark 1 (No strongly convexity assumption) The sample complexity of Algorithm 1 is $O(1/\epsilon)$ for achieving an ϵ excess risk, which is improved from $O(1/\epsilon^2)$.

Remark 2 (No sparsity assumption): The sample complexity of Algorithm 1 has a sub-linear dependence on d due to $R \leq \sqrt{d}$.

Distribution Dependent Randomized Sparsification —Second Stage

Theorem 2: Given the samples in Algorithm 1, let $p_j = \frac{\sqrt{\hat{w}_j^2 E[x_j^2]}}{\sum_{j=1}^d \sqrt{\hat{w}_j^2 E[x_j^2]}}$, $j \in [d]$ in Algorithm 2. In order to have $ER(\tilde{w}, w_*) \leq ER(\hat{w}, w_*) + \epsilon$ with a probability $1 - \delta$ over i_1, i_2, \dots, i_K , it suffices to have

$$K = \left\lceil \frac{(\sum_{i=1}^d \sqrt{\hat{w}_i^2 E[x_i^2]})^2}{\epsilon \delta} \right\rceil$$

Remark 1 (Reduced constant in $O(1/\epsilon)$ for sparsity): The value of K in theorem 2 is less than $\left\lceil \frac{\|\hat{w}\|_1^2}{\epsilon \delta} \right\rceil$, which was obtained from randomized sparsification, since $(\sum_{j=1}^d \sqrt{\hat{w}_j^2 E[x_j^2]})^2 \leq \|\hat{w}\|_1^2$. The inequality holds only when the second moments of individual features are equal.