

Lecture 1: Introduction to the mean field asymptotics

Lecturer: Song Mei

Scriber: Song Mei

Proof reader: Song Mei

In this course, we study the computational and statistical aspects of statistical models in the high dimensional asymptotic limit (the mean-field asymptotics). We will introduce heuristic tools in physics including the replica method and the cavity method. These tools can be made rigorous using approaches including the Gaussian comparison inequality, the leave-one-out analysis and approximate message passing algorithms. Applications of these methods include the spiked matrix model, the LASSO problem, and the double-descent phenomenon.

1 Motivating example: The LASSO problem

We will get a flavor of the difference between the non-asymptotic theory and the asymptotic theory using the example of LASSO.

Let $\mathbf{x}_0 \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{n \times d}$, $\mathbf{w} \in \mathbb{R}^n$ and $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \mathbf{w} \in \mathbb{R}^n$. We consider the case $d \gg n$ but hope that \mathbf{x}_0 is sparse in some sense. To recover \mathbf{x}_0 given \mathbf{A} and \mathbf{y} , we solve the following LASSO problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \frac{\lambda}{n} \|\mathbf{x}\|_1. \quad (1)$$

Note that different papers use different normalization of the LASSO problem. Here the normalization I used is such that the presentation is simpler. When you read a paper on LASSO, you should first look at their normalization and then interpret the results.

1.1 Non-asymptotic theory of LASSO

A line of papers studied the LASSO risk in the non-asymptotic regime. The following result is due to [NRWY12]. Theorem 2 is a fully deterministic statement: the result is satisfied by any deterministic \mathbf{A} , \mathbf{x}_0 , \mathbf{w} , and \mathbf{y} .

Definition 1 (Restricted strong convexity). *We say a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ satisfy restricted strong convexity property, if there exists universal constants c_1 and c_2 , such that for any $\mathbf{x} \in \mathbb{R}^d$, we have*

$$\frac{\|\mathbf{A}\mathbf{x}\|_2^2}{n} \geq c_1 \|\mathbf{x}\|_2^2 - c_2 \frac{\log d}{n} \|\mathbf{x}\|_1^2. \quad (2)$$

Theorem 2 ([NRWY12]). *For any $\mathbf{A} \in \mathbb{R}^{n \times d}$ satisfying the RSC property (2) with constant c_1 and c_2 , there exists universal constants $c < \infty$ such that as long as $\lambda \geq 2\|\mathbf{A}^\top \mathbf{w}\|_\infty$, for any $\mathbf{x}_0 \in \mathbb{R}^d$ and $S \subseteq [d]$ with $|S| \leq n/(c \log d)$, the LASSO estimator (1) satisfies*

$$\|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2 \leq c \frac{\lambda^2 |S|}{n^2} + c \frac{\lambda}{n} \|\mathbf{x}_{0,S^c}\|_1 + c \frac{\log d}{n} \|\mathbf{x}_{0,S^c}\|_1^2.$$

Theorem 2 does not tell us whether there exists a matrix that satisfy the RSC property. The following proposition tells us that, for Gaussian random matrix \mathbf{A} , RSC property holds with high probability.

Proposition 3. *For $\mathbf{A} \in \mathbb{R}^{n \times d}$ with $A_{ij} \sim \mathcal{N}(0, 1)$, Eq. (2) is satisfied for some constant c_1 and c_2 with high probability as $n \rightarrow \infty$.*

In the following, we will make simpler assumptions to understand Theorem 2.

Corollary 4. Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ with $A_{ij} \sim \mathcal{N}(0, 1/\|\mathbf{x}_0\|_2^2)$. Let $\mathbf{x}_0 \in \mathbb{R}^d$ be k -sparse with the support of \mathbf{x}_0 given by S . Let \mathbf{w} be σ^2 -sub-Gaussian. Then for any $\delta > 0$, there exists constant $C(\delta)$ such that, as long as we take $n \geq C(\delta)k \log d$ and $\lambda \geq C(\delta) \cdot \sigma \sqrt{n \log d}$, then with probability at least $1 - \delta$, the LASSO estimator (1) satisfies

$$\|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2 / \|\mathbf{x}_0\|_2^2 \leq \frac{C(\delta)\sigma^2 k \log d}{n}.$$

The corollary tells us that, to well-estimate a k -sparse ground truth vector, it is enough to have sample size $n \gg k \log d$.

1.2 High dimensional asymptotics of LASSO

Note that the non-asymptotic theory of LASSO does not allow us to consider the proportional regime $n \propto k \propto d$. In many cases, however, this proportional regime is very interesting. It would be desirable to establish a theory to characterize the performance of LASSO in this regime.

Theorem 5 ([BM11]). We consider the asymptotic limit when $n/d \rightarrow \delta \in (0, \infty)$ as $d \rightarrow \infty$. Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ with $A_{ij} \sim \mathcal{N}(0, 1/n)$. Let $\mathbf{x}_0 \in \mathbb{R}^d$ with $\mathbf{x}_{0,i} \sim_{iid} \mathbb{P}_0$. Let $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Let $\hat{\mathbf{x}}$ be the LASSO estimator (1). Then we have

$$\lim_{d, n \rightarrow \infty} \frac{1}{d} \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2 = \mathbb{E}_{(X_0, Z) \sim \mathbb{P}_0 \times \mathcal{N}(0, 1)} [(\eta(X_0 + \tau_* Z; \theta_*) - X_0)^2],$$

where $\eta(x) = \text{sign}(x) \cdot (|x| - 1)_+$ is the soft thresholding function and $\tau_* = \tau_*(\alpha_*)$ is the largest solution of

$$\tau^2 = \sigma^2 + \delta^{-1} \mathbb{E}_{(X_0, Z) \sim \mathbb{P}_0 \times \mathcal{N}(0, 1)} \{[\eta(X_0 + \tau Z; \alpha_* \tau) - X_0]^2\},$$

and α_* is the unique non-negative solution of

$$\lambda = \alpha \tau_*(\alpha) \cdot \left[1 - \delta^{-1} \mathbb{E}[\eta'(X_0 + \tau_*(\alpha) Z; \alpha \tau_*(\alpha))] \right].$$

Moreover, for any Lipschitz function ψ , we have almost surely

$$\lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^d \psi(\hat{x}_i, x_{0,i}) = \mathbb{E}_{(X_0, Z) \sim \mathbb{P}_0 \times \mathcal{N}(0, 1)} [\psi(\eta(X_0 + \tau_* Z; \alpha_* \tau_*), X_0)].$$

We can plot the limiting risk versus the regularization parameter λ , which is given in Figure 1.2. This curve gives the precise U-shaped curve for the Bias and Variance tradeoff of LASSO estimator. Note that this U-shaped curve cannot be completely captured by the non-asymptotic theory, since the non-asymptotic theory doesn't give lower and upper bounds that match up to $1 + o(1)$. The sharp characterization of the risk is an advantage of the high dimensional asymptotic theory.

1.3 Comparison of non-asymptotic theory and high dimensional asymptotics

Here we present a table that compares the non-asymptotic theory versus the asymptotic theory.

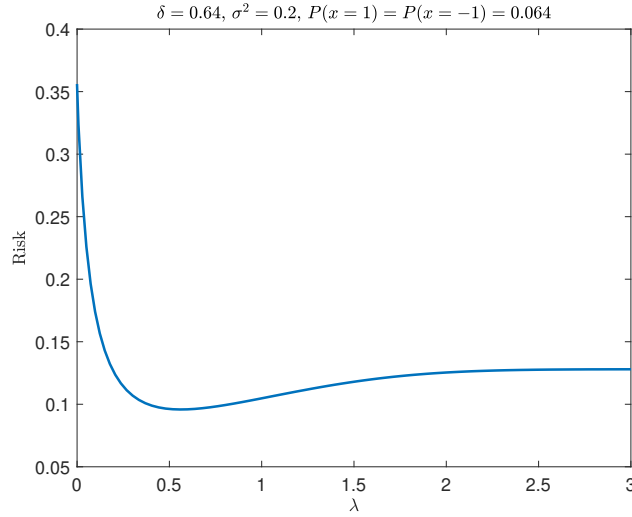


Figure 1: The risk of the LASSO estimator

	Non-asymptotics theory	High dimensional asymptotics
Typical regime	(Relatively) Strong signal-to-noise ratio ($n \gg k \log d$)	Constant signal-to-noise ratio ($n \propto d \propto k$)
Advantages	Less model assumptions. Result holds for any finite parameter size.	Precise asymptotic formula: upper and lower bounds match sharply.
Limitations	A gap of upper and lower bounds up to constant or logarithmic factors.	More detailed model assumptions. (Sometimes) hard to control how large should the parameter be so that the asymptotic regime kick in.
When useful?	Characterize the behavior of a model or an algorithm with general assumptions.	Identify the exact location of phase transition.
Examples	Statistical learning theory: bounding excessive risk by uniform convergence. Analyzing the non-convex landscape of empirical risk minimization.	The phase transition phenomenon for compressed sensing. Understanding the double-descent phenomenon. The optimal loss function in machine learning

2 Methods from statistical physics

The focus of this course is to analyze statistical models through the high dimensional asymptotic viewpoint. In many cases, we are interested in deriving the asymptotic formula instead of proving the formula rigorously, and statistical physics tools can be used to predict these formula. These predicted formula can be simply verified through experiments, but they are difficult to be made rigorous. Many predictions have been made rigorous in some way, and typically proving these formula is much more complicated than deriving them. Furthermore, there are still some predictions that are hard to be proved rigorously.

In this course, we will introduce two useful tools in physics: the “replica method” and the “cavity method”. We will show how these tools can be used to predict interesting behaviors of statistical models and algorithms. Simple models will be used as examples in class: the spiked GOE matrix and the LASSO problem. We will revisit these models several times. First we will show how the replica method can be used to predict the behavior of these models. Then we will show how these predictions can be proved using rigorous tools. These rigorous tools include the Gaussian comparison theorem, the Stieltjes transforms, and

approximate message passing (AMP) algorithms, etc.

References

- [BM11] Mohsen Bayati and Andrea Montanari, *The lasso risk for gaussian matrices*, IEEE Transactions on Information Theory **58** (2011), no. 4, 1997–2017.
- [NRWY12] Sahand N Negahban, Pradeep Ravikumar, Martin J Wainwright, and Bin Yu, *A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers*, Statistical science **27** (2012), no. 4, 538–557.