

Replica methods

Homework submissions are expected to be in pdf format produced by L^AT_EX. For programming exercises, you can use your favorite programming language, but matlab or Python are recommended. Grading of programming exercises will be based on the figures produced by the simulations, but you need also submit the source code in the original format (e.g. ‘.m’ or ‘.py’).

Exercise 1: Large deviation of the overlap matrix

Let $\mathbf{x}_0 = (x_{0,1}, \dots, x_{0,d})^\top \in \mathbb{R}^d$ with $\frac{1}{d} \sum_{i=1}^d \delta_{x_{0,i}} \rightarrow \mathbb{P}_0$ and $\mathbb{E}_{X \sim \mathbb{P}_0}[X^2] = 1$ (this will be satisfied with high probability if $x_{0,i} \sim_{iid} \mathbb{P}_0$). Let $\mathbf{x}_a \sim_{iid} \text{Unif}([-1, 1]^d)$ for $1 \leq a \leq k$. Denote $\mathbf{X} = (\mathbf{x}_a)_{0 \leq a \leq k}$, and denote

$$\widehat{\boldsymbol{\mu}}(\mathbf{X}) = \left(\langle \mathbf{x}_0, \mathbf{x}_a \rangle / d \right)_{1 \leq a \leq k} \in \mathbb{R}^k, \quad \widehat{\mathbf{Q}}(\mathbf{X}) = \left(\langle \mathbf{x}_a, \mathbf{x}_b \rangle / d \right)_{1 \leq a, b \leq k} \in \mathbb{R}^{k \times k}.$$

We further define (assuming the limit exists)

$$\text{ent}(k, \boldsymbol{\mu}, \mathbf{Q}) \equiv \lim_{\varepsilon \rightarrow 0} \lim_{d \rightarrow \infty} \frac{1}{d} \log \mathbb{P} \left(\|\widehat{\boldsymbol{\mu}}(\mathbf{X}) - \boldsymbol{\mu}\|_2^2 \leq \varepsilon, \|\widehat{\mathbf{Q}}(\mathbf{X}) - \mathbf{Q}\|_F^2 \leq \varepsilon \right).$$

Question 1: Please use the field theoretic approach to calculate $\text{Ent}(k, \boldsymbol{\mu}, \mathbf{Q})$.

Hint: You can start with the following heuristics (for $\mathbf{Q} = (q_{ab})_{0 \leq a, b \leq k}$)

$$\mathbb{P}(\mathbf{Q}(\mathbf{X}) \approx \mathbf{Q}) \doteq \frac{1}{2^{dk}} \int_{[-1, 1]^{dk}} \prod_{1 \leq a \leq k} \delta(\langle \mathbf{x}_0, \mathbf{x}_a \rangle - d \cdot \mu_a) \prod_{1 \leq a \leq b \leq k} \delta(\langle \mathbf{x}_a, \mathbf{x}_b \rangle - d \cdot q_{ab}) \prod_{1 \leq a \leq k} d\mathbf{x}_a,$$

where $\int_{[-1, 1]^d} f(\mathbf{x}_a) d\mathbf{x}_a$ is integration with respect to Lebesgue measure on $[-1, 1]^d$.

Question 2: Take

$$\boldsymbol{\mu} = \begin{bmatrix} \mu \\ \dots \\ \mu \end{bmatrix}, \quad \mathbf{Q} = \begin{bmatrix} q_1 & \dots & q_0 \\ \dots & \dots & \dots \\ q_0 & \dots & q_1 \end{bmatrix}.$$

That is $\mu_a = \mu$, $q_{aa} = q_1$ for $1 \leq a \leq k$, and $Q_{ab} = q_0$ for $1 \leq a \neq b \leq k$.

Please simplify $\overline{\text{ent}}(k, \mu, q_0, q_1) \equiv \text{ent}(k, \boldsymbol{\mu}, \mathbf{Q})$.

Question 3: Give a simple expression for

$$T(\mu, q_0, q_1) \equiv \lim_{k \rightarrow \infty} \frac{1}{k} \overline{\text{ent}}(k, \mu, q_0, q_1).$$

Hint: Use the Gaussian identity formula and the replica trick $\lim_{k \rightarrow 0} k^{-1} \log \mathbb{E}[Z^k] = \mathbb{E}[\log Z]$.

Exercise 2: Replica formula for Box-LASSO

Let $\mathbf{x}_0 \in \mathbb{R}^d$ with $\frac{1}{d} \sum_{i=1}^d \delta_{x_{0,i}} \rightarrow \text{Unif}(\{1, -1\})$ (this will be satisfied with high probability if $x_{0,i} \sim_{iid} \text{Unif}(\{1, -1\})$). Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ with $A_{ij} \sim_{iid} \mathcal{N}(0, 1/n)$. Let $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ with $\varepsilon_i \sim_{iid} \mathcal{N}(0, \sigma^2)$. Assume all these random variables are independent. Let $\mathbf{y} = \mathbf{A}\mathbf{x}_0 + \boldsymbol{\varepsilon}$. Denote the Box-LASSO estimator by

$$\hat{\mathbf{x}} \equiv \arg \min_{\mathbf{x} \in [-1, 1]^d} \frac{1}{2n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2,$$

and the training error by

$$L_n \equiv \min_{\mathbf{x} \in [-1,1]^d} \frac{1}{2n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2.$$

Question 1: Calculate $\lim_{n \rightarrow \infty} \mathbb{E}[L_n]$ using the replica method. Many results in class can be reused here. You will also use the results in Exercise 1. The answer is expected to be represented as some variational formula which involves integration with respect to some Gaussians.

Question 2: For $\psi : \mathbb{R}^2 \rightarrow \mathbb{R}$, give a prediction for the expected observable

$$\lim_{d \rightarrow \infty, n/d \rightarrow \delta} \frac{1}{d} \sum_{i=1}^d \mathbb{E}[\psi(\hat{x}_i, x_{0,i})].$$

Especially, derive the asymptotics for (here we choose not to round the estimator $\hat{\mathbf{x}}$)

$$\lim_{d \rightarrow \infty, n/d \rightarrow \delta} \mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2]/d.$$

Exercise 3: Simulating the risk of Box-LASSO (Programming exercise)

Please do the following simulation: for each $d = 50, 200, 800$, and take $n = 2d$, and each $\sigma^2 \in [0, 8]$ with grid size 0.2, obtain $n_s = 10$ **independent** samples of instances $(\mathbf{A}^{(i)}, \mathbf{x}_0^{(i)}, \mathbf{y}^{(i)})$, $i \in [n_s]$. Compute the Box-LASSO estimator $\hat{\mathbf{x}}$ as in Exercise 2. Denote

$$R = \frac{1}{d} \|\hat{\mathbf{x}} - \mathbf{x}_0\|_2^2.$$

Make an error bar plot of $(\hat{\mathbb{E}}[R], \widehat{\text{std}}[R])$ versus σ^2 (σ^2 will be the x-axis), with $\hat{\mathbb{E}}[R]$ the mean variable and $\widehat{\text{std}}[R]$ the error variable. Please plot the curve for $d = 50, 200, 800$ on the same figure.

Exercise 4: Prediction versus simulations (Extra credit)

Plot your prediction on the same figure (in Exercise 2). You earn the extra credit if the prediction and the simulation match.