

Lecture 19. Approximate message passing algorithms.

① Algorithm for Gibbs mean (Statistical estimators).

Configuration space: $\Omega \subseteq \mathbb{R}^d$. Measure $\nu_0 \in \mathcal{P}(\mathbb{R}^d)$.Random Hamiltonian: $H: \Omega \rightarrow \mathbb{R}$.Gibbs distribution: $P_\beta(d\sigma) \propto \exp\{-\beta H(\sigma)\}\nu_0(d\sigma)$.Goal: (approximately) calculate $\langle \sigma \rangle_\beta \in \mathbb{R}^d$ using efficient algorithm.

Example: Bayes linear model, Bayes estimator, and LASSO estimator.

$$x_0 \in \mathbb{R}^d, \quad x_{0i} \sim \text{iid } P_0, \quad A \in \mathbb{R}^{n \times d}$$

$$y = Ax_0 + w, \quad w_i \sim \text{iid } N(0, \sigma^2),$$

Posterior mean estimator

$$\hat{x}_{\text{Bayes}} = \int_{\mathbb{R}^d} x \times P(x|y, A) dx \in \mathbb{R}^d$$
$$P(x|A, y) \propto \exp\left\{-\frac{\|y - Ax\|_2^2}{2\sigma^2}\right\} \prod_{i=1}^d P_0(x_i)$$

LASSO estimator:

$$\hat{x}_{\text{LASSO}} = \underset{x}{\operatorname{argmin}} \left\{ \frac{1}{2n} \|y - Ax\|_2^2 + \frac{\lambda}{n} \|x\|_1 \right\} = \lim_{\beta \rightarrow \infty} \int_{\mathbb{R}^d} x \times P_\beta(x|y, A) dx$$

$$P_\beta(x|A, y) \propto \exp\left\{-\beta \left[\frac{1}{2n} \|y - Ax\|_2^2 + \frac{\lambda}{n} \|x\|_1 \right] \right\}.$$

Remark: We have seen many concentrated observables,

$$\text{e.g. } O(x) = \frac{1}{d} \sum_{i=1}^d x_i^2 \Rightarrow \langle O \rangle_\beta \approx \mathbb{E}[\langle O \rangle_\beta].$$
 Randomness come from $H, i.e. A, x_0, w$.

Randomness come from
 $H, i.e.
A, x_0, w$

So, approximating $\langle O \rangle_\beta$ can be independent of specific instance of A, x_0, w . Limiting value depends on their distribution instead of realization.However, EA of $O(x) = x_i$ typically doesn't concentrate.

$$\langle x_i \rangle_\beta \neq \mathbb{E}[\langle x_i \rangle_\beta].$$

 \Rightarrow Approximating $\langle x_i \rangle_\beta$ needs to use the specific instance H .AMP is an algorithm used to calculate $\hat{x} = \langle \sigma \rangle_\beta$.

② ISTA and FISTA for LASSO.

$$\hat{x}_{\text{LASSO}} = \arg \min_{x} \mathcal{L}(x) = \arg \min_{x} \frac{1}{2n} \|y - Ax\|_2^2 + \frac{\lambda}{n} \|x\|_1.$$

Convex optimization problem.

PGD/ISTA (Proximal gradient descent; Iterative thresholding algorithm)

$$x^{k+1} = \arg \min_x \left[\frac{1}{2s_k} \|x - (x^k - s_k \nabla f(x^k))\|_2^2 + g(x) \right] \quad \text{PGD}$$

Proximal: $x^{k+1} = \arg \min_x \left[\frac{1}{2s_k} \|x - x^{k+1}\|_2^2 + g(x) \right]$

$$\Rightarrow x^{k+1} = \eta(x^k - s_k A^T(Ax^k - y); \lambda s_k) \quad x^1 = 0.$$

$\eta(x; \theta) = (\|x\| - \theta) \cdot \mathbf{1}\{\|x\| > \theta\}.$

[Beck, Teboulle, 2008]

Thm: Suppose $f \in C^2(\mathbb{R})$, $\sup_x \|\nabla^2 f(x)\|_{\text{op}} \leq \beta$, convex.

$g \in C(\mathbb{R})$, convex.

$$\mathcal{L}(x) = f(x) + g(x). \quad \arg \min_x \mathcal{L}(x) \neq \emptyset.$$

Then $s_k = \frac{1}{\beta}$, x^k : k'th PGD iteration.

$$\Rightarrow \mathcal{L}(x^k) - \min_x \mathcal{L}(x) \leq \frac{\beta \|x^1 - x_k\|_2^2}{2k} = O\left(\frac{1}{k}\right)$$

Proof is fully deterministic, based on Jensen's inequality and simple algebra

$$\text{APGD/FISTA.} \quad \mu_1 = 0. \quad \mu_k = \frac{1 + \sqrt{1 + 4\mu_{k-1}^2}}{2}. \quad r_k = \frac{1 - \mu_k}{\mu_{k+1}}.$$

$$\begin{cases} v^{k+1} = \arg \min_x \left\{ \frac{\beta}{2} \|x - (x^k - \frac{1}{\beta} \nabla f(x^k))\|_2^2 + g(x) \right\} \\ x^{k+1} = (1 - r_k) v^{k+1} + r_k v^k. \end{cases}$$

$$x^{k+1} = (1 - r_k) \eta(x^k - \frac{1}{\beta} A^T(Ax^k - y); \frac{\lambda}{\beta}) + r_k \eta(x^{k-1} - \frac{1}{\beta} A^T(Ax^{k-1} - y); \frac{\lambda}{\beta}).$$

Thm: Under the same assumptions.

Let x^k : k'th iteration of APGD. then

$$\mathcal{L}(x^k) - \min_x \mathcal{L}(x) \leq \frac{2\beta \|x^1 - x_k\|_2^2}{k^2} = O\left(\frac{1}{k^2}\right).$$

Remark:

③ The approximate algorithm passing algorithm for LASSO.

ISTA w/ Step size 1.

$$\eta'(x; \theta) = \partial_x \eta(x; \theta).$$

$$\begin{cases} x^{k+1} = \eta(x^k + A^T z^k; \theta_k) \\ z^k = y - Ax^k \end{cases} \quad \theta_k = \lambda.$$

AMP: (Will show how this can be derived in later lectures)

$$\begin{cases} x^{k+1} = \eta(x^k + A^T z^k; \theta_k) \\ z^k = y - Ax^k + \underbrace{\omega_k z^{k-1}}_{\text{Onsager correction term. (Intuition will be given later)}}. \end{cases} \quad \theta_k, \omega_k \text{ scalars.}$$

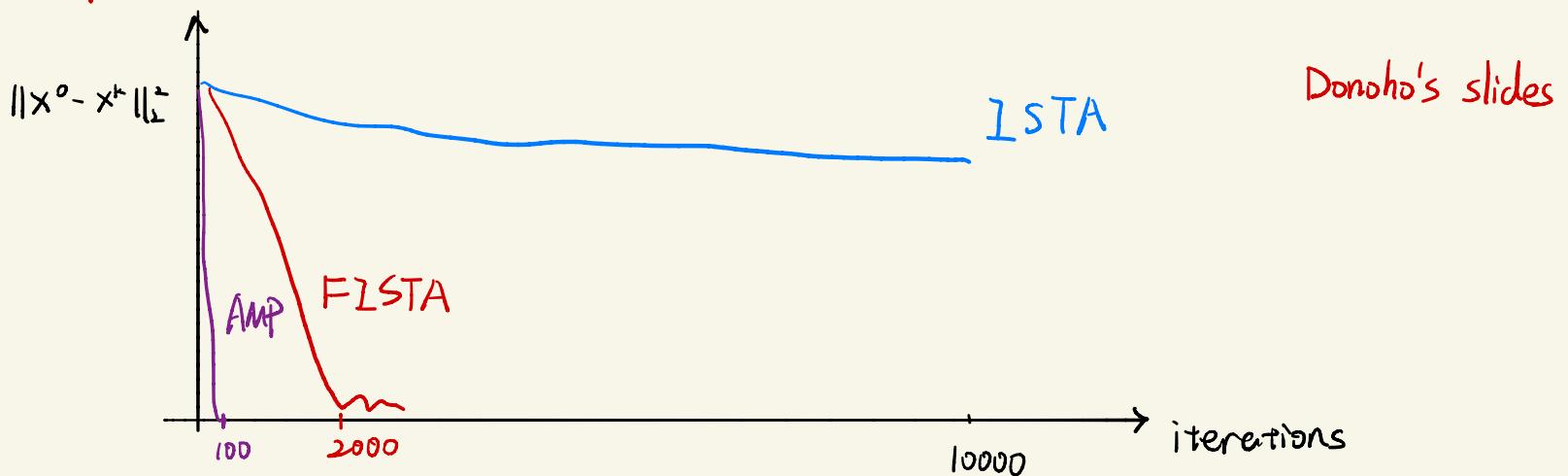
$$\omega_k = \frac{1}{S} \sum_{i=1}^d \eta'(x^{k-1} + A^T z^{k-1}; \theta_{k-1}); \quad \eta'(x^{k-1} + A^T z^{k-1}; \theta_{k-1}) \in \mathbb{R}^d$$

θ_k : needs to be properly chosen. Subtle. Now just assume arbitrary.

Remark: ⊗ AMP is not descending. Not a deterministic proof of convergence.

⊗ Convergence analysis of AMP depend on assumptions of A, y .

Performance comparison. $n=1200, d=400, k=60$ (Sparsity). $\sigma^2=0$



Assumptions:

a) $A \in \mathbb{R}^{n \times d}, \quad A_{ij} \sim \text{iid } N(0, \frac{1}{n}).$

b) $x_0 \in \mathbb{R}^d, \quad \frac{1}{d} \sum_{i=1}^d \delta_{x_0, i} \xrightarrow{\text{weak}} P_{x_0}, \quad \frac{1}{d} \sum_{i=1}^d x_0, i^2 \rightarrow \mathbb{E}_{x_0}[x_0^2]$

c) $w \in \mathbb{R}^d, \quad \frac{1}{n} \sum_{i=1}^n \delta_{w, i} \xrightarrow{\text{weak}} P_w, \quad \frac{1}{n} \sum_{i=1}^n w_i^2 \rightarrow \mathbb{E}_w[w^2] = \sigma^2.$

d) $y = Ax_0 + w \in \mathbb{R}^n.$

e) $n/d \rightarrow S.$

④ Theoretical analysis of AMP.

State evolution characterization:

$$\tau_{k+1}^2 = F(\tau_k^2, \theta_k) \quad SE$$

$$\text{where } F(\tau^2, \theta) = \sigma^2 + \frac{1}{\delta} \mathbb{E}[(\eta(x_0 + \tau G; \theta) - x_0)^2].$$

$$(x_0, G) \sim \mathbb{P}_{x_0} \times N(0, I).$$

Thm 1 ([Bayati, Montanari, 2011]).

Let assumptions. a) — e) hold. Let $\psi: \mathbb{R}^+ \rightarrow \mathbb{R}$.

be Pseudo Lipschitz. i.e. $|\psi(x) - \psi(y)| \leq K \cdot \|x - y\|_2 \cdot (1 + \|x\|_2 + \|y\|_2)$

Then almost surely.

$$\lim_{\substack{d \rightarrow \infty \\ n/d \rightarrow 0}} \frac{1}{d} \sum_{i=1}^d \psi(x_i^{k+1}, x_{0,i}) = \mathbb{E}[\psi(\eta(x_0 + \tau_k G; \theta_k), x_0)].$$

$$\text{Rmk: } S(\tau_{k+1}^2 - \sigma^2) = \lim_{\substack{d \rightarrow \infty \\ n/d \rightarrow 0}} \|x^k - x_0\|_2^2 / d \quad \psi(x, x_0) = (x - x_0)^2$$

AMP (High dim dynamics) $\xleftrightarrow{\text{Thm 1}}$ State evolution

\uparrow Prop 2.
LASSO estimator.

\downarrow Prop 1
Fixed point equation.

Let $(\tau_\infty, \alpha_\infty)$ be a proper solution of the fixed point equation

$$\begin{cases} \tau^2 = \sigma^2 + \delta^{-1} \mathbb{E} \{ [\eta(x_0 + \tau G; \alpha \tau) - x_0]^2 \} \\ \lambda = \alpha \tau \{ 1 - \delta^{-1} \mathbb{E} [\eta'(x_0 + \tau G; \alpha \tau)] \} \end{cases}$$

$$\delta \geq \delta_*(\sigma^2, \lambda)$$

Prop 1: Take $\theta_k = \alpha_\infty \tau_k$, the state evolution $\{\tau_k^2\}_{k \geq 1}$.

$$\tau_{k+1}^2 = F(\tau_k^2, \alpha_\infty \tau_k) \text{ converges to } \tau_\infty^2.$$

The convergence is exponentially fast.

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

$\{\mathbf{x}^k\}_{k \geq 1}$ AMP iterates. Then

$$\lim_{k \rightarrow \infty} \lim_{n \rightarrow \infty} \|\hat{x}(\lambda) - \mathbf{x}^k\|_2^2 / d = 0.$$

Proof strategy: \otimes Show $\exists v_k \in \partial \mathcal{C}(x^k)$ (Thm 1).

$$\lim_{k \rightarrow \infty} \lim_{d \rightarrow \infty} \|v_k\|_2^2 / d = 0.$$

\otimes Convex property of $\mathcal{C}(x)$.

Remark: $\lim_{n \rightarrow \infty} \|\hat{x}(\lambda) - \mathbf{x}^k\|_2^2 / d \leq e^{-ck}. \lim_{n \rightarrow \infty} \|\hat{x}(\lambda) - \mathbf{x}^0\|_2^2 / d. \forall k$
 $(\|\hat{x}(\lambda) - \mathbf{x}^k\|_2^2 / d \leq e^{-ck}. \|\hat{x}(\lambda) - \mathbf{x}^0\|_2^2 / d \forall k)$

Thm 2: ([Bayati, Montanari, 2011]).

Let assumptions a) – e) holds. Let $\psi: \mathbb{R}^2 \rightarrow \mathbb{R}$ be any pseudo-Lipschitz function. i.e. $|\psi(x) - \psi(y)| \leq K \cdot \|x - y\|_2 \cdot (1 + \|x\|_2 + \|y\|_2)$. Then almost surely

$$\lim_{\substack{d \rightarrow \infty \\ n/d \rightarrow s}} \frac{1}{d} \sum_{i=1}^d \psi(\hat{x}(\lambda), x_{0,i}) = \mathbb{E}[\psi(\eta(x_0 + T_\lambda G; \alpha + T_\lambda); X_0)].$$

Remark: Similar analysis and results for a wide range of problem.
 ↳ synchro. Generalized linear models.

Plans for remaining lectures: