OOPS 2020
Mean field methods in high-dimensional statistics and nonconvex optimization
Lecturer: Andrea Montanari
Problem session leader: Michael Celentano
July 7, 2020

# Problem Session 1

**Problem 1: from Gordon's objective to the fixed point equations**

Recall *Gordon's min-max problem* is

$$B^*(\boldsymbol{g}, \boldsymbol{h}) := \min_{\boldsymbol{u} \in \mathbb{R}^d} \max_{\boldsymbol{v} \in \mathbb{R}^n} \left\{ \frac{1}{n} \|\boldsymbol{v}\| \langle \boldsymbol{g}, \boldsymbol{u} \rangle + \frac{1}{n} \|\boldsymbol{u}\| \langle \boldsymbol{h}, \boldsymbol{v} \rangle - \frac{\sigma}{n} \langle \boldsymbol{w}, \boldsymbol{v} \rangle - \frac{1}{2n} \|\boldsymbol{v}\|^2 + \frac{\lambda}{\sqrt{n}} \|\boldsymbol{\theta}_0 + \boldsymbol{u}\|_1 \right\}. \qquad (1)$$

In lecture, we claimed that by analyzing Gordon's objective we can show that the Lasso solution is described in terms of the solutions $\tau^*, \beta^*$ to the *fixed point equations*

$$\tau^2 = \sigma^2 + \frac{1}{\delta} \mathbb{E} \left[ (\eta(\Theta + \tau Z; \tau \lambda / \beta) - \Theta)^2 \right],$$

$$\beta = \tau \left( 1 - \frac{1}{\delta} \mathbb{E}[\eta'(\Theta + \tau Z; \tau \lambda / \beta)] \right),$$

where $\eta$ is the solution of the 1-dimensional problem

$$\eta(y; \alpha) := \arg\min_{x \in \mathbb{R}} \left\{ \frac{1}{2}(y - x)^2 + \alpha |x| \right\} = (|x| - \alpha)_+ \mathsf{sign}(x).$$

$\eta$ is commonly known as *soft-thresholding*. In this problem, we will outline how to derive the fixed point equations from Gordon's min-max problem.

(a) Define $\tilde{\boldsymbol{\theta}}_0 = \sqrt{n}\boldsymbol{\theta}_0$ and $\tilde{\boldsymbol{u}} = \sqrt{n}\boldsymbol{u}$. Prove that $B^*(\boldsymbol{g}, \boldsymbol{h})$ has the same distribution as

$$\min_{\tilde{\boldsymbol{u}}\in\mathbb{R}^d} \max_{\beta\geq 0} \left\{ \left( \left\| \frac{\|\tilde{\boldsymbol{u}}\|}{\sqrt{n}} \frac{\boldsymbol{h}}{\sqrt{n}} - \frac{\sigma\boldsymbol{w}}{\sqrt{n}} \right\| - \frac{1}{n}\langle \boldsymbol{g}, \tilde{\boldsymbol{u}}\rangle \right) \beta - \frac{1}{2}\beta^2 + \frac{\lambda}{n}\|\tilde{\boldsymbol{\theta}}_0 + \tilde{\boldsymbol{u}}\|_1 \right\}. \tag{2}$$

**Hint:** Let $\beta = \|\boldsymbol{v}\|/\sqrt{n}$

$$\frac{1}{n^{3/2}} \|v\| \langle g, \tilde{v}\rangle + \frac{1}{n^{3/2}} \|\tilde{u}\| \langle h, v\rangle - \frac{\sigma}{n} \langle w, v\rangle - \frac{1}{2n}\|v\|^2 + \frac{\lambda}{n} \|\hat{\theta}_0 + \tilde{u}\|_1$$

$$\left\langle \frac{\|\tilde{u}\|}{n^{1/2}} \frac{h}{n^{1/2}} - \frac{\sigma w}{n^{1/2}}, \frac{v}{n^{1/2}} \right\rangle$$

$$\left\| \frac{\|\tilde{u}\|}{n^{1/2}} \frac{h}{a^{1/2}} - \frac{\sigma w}{n^{1/2}} \right\| \frac{\|v\|}{n^{1/2}}$$

$$\underbrace{\qquad\qquad}_{\beta}$$

$$\sqrt{\frac{\|\tilde{u}\|^2}{n} + \sigma^2}$$

(b) Argue (heuristically) that we may approximate the optimization above by

$$\min_{\tilde{\boldsymbol{u}}\in\mathbb{R}^d}\max_{\beta\geq 0}\left\{\left(\sqrt{\frac{\|\tilde{\boldsymbol{u}}\|^2}{n}+\sigma^2}-\frac{1}{n}\langle\boldsymbol{g},\boldsymbol{u}\rangle\right)\beta-\frac{1}{2}\beta^2+\frac{\lambda}{n}\|\tilde{\boldsymbol{\theta}}_0+\tilde{\boldsymbol{u}}\|_1\right\}. \tag{3}$$

**Remark:** If we maximize over $\beta\geq 0$ explicitly, we get

$$\min_{\tilde{\boldsymbol{u}}\in\tilde{U}}\left\{\frac{1}{2}\left(\sqrt{\frac{\|\tilde{\boldsymbol{u}}\|^2}{n}+\sigma^2}-\frac{\langle\boldsymbol{g},\tilde{\boldsymbol{u}}\rangle}{n}\right)^2_+ +\frac{\lambda}{n}\|\tilde{\boldsymbol{\theta}}_0+\tilde{\boldsymbol{u}}\|_1\right\}.$$

Note that the objective on the right-hand side is locally strongly convex around any point $\tilde{\boldsymbol{u}}$ at which the first term is positive. When $n < d$, the Lasso objective is nowhere locally strongly convex. This convenient feature of the new form of Gordon's problem is very useful for its analysis. We do not explore this further here.

$$\frac{1}{2}\|y-X\theta\|^2 + \lambda\|\theta\|_1$$

$$n < d$$

$$\sqrt{x} = \min_{\tau\geq 0}\ \frac{x}{2\tau}+\frac{\tau}{2}$$

$$\max_{\beta\geq 0}\ \min_{\tilde{v}\in\mathbb{R}^d}\left\{\left(\min_{\tau\geq 0}\ \frac{\|\tilde{v}\|/n+\sigma^2}{2\tau}+\frac{\tau}{2}\right)-\frac{1}{n}\langle g,\tilde{v}\rangle\right)\beta\right.$$
$$\left.-\frac{1}{2}\beta^2+\frac{\lambda}{n}\|\hat{\theta}_0+\tilde{v}\|_1\right\}$$

$$\max_{\beta\geq 0}\ \min_{\tau\geq 0}\left\{\frac{\sigma^2\beta}{2\tau}+\frac{\tau\beta}{2}-\frac{1}{2}\beta^2+\frac{1}{n}\min_{\tilde{v}\in\mathbb{R}^d}\left\{\frac{\beta}{2\tau}\|\tilde{v}\|^2-\beta\langle g,\tilde{v}\rangle\right.\right.$$
$$\left.\left.+\frac{\lambda}{n}\|\hat{\theta}_0+\tilde{v}\|_1\right\}\right\}$$

(c) Argue that the quantity in Eq. (3) is equal to

$$\max_{\beta \geq 0} \min_{\tau \geq 0} \left\{ \frac{\sigma^2 \beta}{2\tau} + \frac{\tau\beta}{2} - \frac{\beta^2}{2} + \frac{1}{n} \min_{\tilde{u} \in \mathbb{R}^d} \left\{ \frac{\beta}{2\tau} \|\tilde{u}\|^2 - \beta \langle g, \tilde{u} \rangle + \lambda \|\tilde{\theta}_0 + \tilde{u}\|_1 \right\} \right\}. \tag{4}$$

**Hint:** Recall the identity $\sqrt{x} = \min_{\tau \geq 0} \left\{ \frac{x}{2\tau} + \frac{\tau}{2} \right\}$.   $x \geq 0$

$\tau \geq 0$

$$\frac{1}{d} \sum_{j=1}^{d} \left\{ \frac{\beta}{2\tau} \tilde{u}_j^2 - \beta g_\delta u_\delta + \lambda |\hat{\theta}_{0j} + \tilde{u}_j| \right\}$$

$$\underbrace{\hspace{8cm}}$$

$$\frac{1}{d} \sum_{j=1}^{d} \delta_{(\tilde{u}_j, g_j, \hat{\theta}_{vj})} \xrightarrow{\ w\ } \mu$$

$$\underline{\hspace{12cm}}$$

$\text{wrt } \tau \qquad -\frac{\sigma^2 \beta}{\lambda \tau^2} + \frac{\beta}{2} - \underbrace{\frac{\beta}{2n\tau^2} \mathbb{E}\left[\|\hat{u}(g)\|^2\right]}_{} = 0 \qquad \begin{array}{c} \text{soln to part} \\ \text{(e).} \end{array}$

$$-\frac{1}{\lambda \tau^2} \cdot \frac{1}{\delta} \mathbb{E}\left[ (\eta(\eta + \tau z; \lambda \tau / \beta) - \eta)^2 \right]$$

$$\tau^2 = \sigma^2 + \frac{1}{\delta} \mathbb{E}\left[ |\eta(\eta + \tau z; \lambda \tau / \beta) - \eta|^2 \right]$$

$$\underbrace{\hspace{9cm}}$$

$\underline{\text{wrt } \beta}$

$$\frac{\sigma^2}{2\tau} + \frac{\tau}{2} - \beta + \underbrace{\frac{1}{2n\tau} \mathbb{E}\left[\|\hat{u}\|^2\right] - \frac{1}{n} \mathbb{E}\left[\langle g, \hat{u}(g)\rangle\right]}_{\frac{1}{2} \frac{1}{\tau}(\tau^2 - \sigma^2) - \frac{\tau}{\delta} \mathbb{P}\left(|\eta + \tau z| \geq \lambda \tau/\beta\right)} = 0$$

$$\beta = \tau \left( 1 - \frac{1}{\delta} \mathbb{P}\left(|\eta + \tau z| \geq \lambda \tau / \beta\right) \right)$$

(d) Write
$$\hat{\boldsymbol{u}} := \arg\min_{\tilde{\boldsymbol{u}} \in \mathbb{R}^d} \left\{ \frac{\beta}{2\tau} \|\tilde{\boldsymbol{u}}\|^2 - \beta\langle \boldsymbol{g}, \tilde{\boldsymbol{u}} \rangle + \lambda\|\tilde{\boldsymbol{\theta}}_0 + \tilde{\boldsymbol{u}}\|_1 \right\}$$

in terms of the soft-thresholding operator.

(e) Compute

   (i) the derivative of $\min_{\tilde{\boldsymbol{u}} \in \mathbb{R}^d} \left\{ \frac{\beta}{2\tau} \|\tilde{\boldsymbol{u}}\|^2 - \beta\langle \boldsymbol{g}, \tilde{\boldsymbol{u}} \rangle + \lambda\|\tilde{\boldsymbol{\theta}}_0 + \tilde{\boldsymbol{u}}\|_1 \right\}$ with respect to $\tau$.

   (ii) the derivative of $\min_{\tilde{\boldsymbol{u}} \in \mathbb{R}^d} \left\{ \frac{\beta}{2\tau} \|\tilde{\boldsymbol{u}}\|^2 - \beta\langle \boldsymbol{g}, \tilde{\boldsymbol{u}} \rangle + \lambda\|\tilde{\boldsymbol{\theta}}_0 + \tilde{\boldsymbol{u}}\|_1 \right\}$ with respect to $\beta$.

$$\frac{\beta}{\tau}\left( \frac{1}{2}\|\tilde{v} - \tau g\|^2 - \frac{\cancel{\tau^2 \|g\|^2}}{2} + \frac{\lambda \tau}{\beta}\|\tilde{\theta}_0 + \tilde{v}\|_1 \right)$$

$$\theta = \tilde{\theta}_0 + \tilde{v}$$

$$\frac{1}{2}\|\theta - (\tilde{\theta}_0 + \tau g)\|^2 + \frac{\lambda\tau}{\beta}\|\theta\|_1$$

$$\underline{\hat{v} = \eta\left(\tilde{\theta}_0 + \tau g; \lambda\tau/\beta\right) - \tilde{\theta}_0}$$

(e) $\frac{d}{d\beta} \min_{x} C(x;\beta) = \partial_\beta C(\hat{x};\beta)$

$$C(\hat{x}(\beta);\beta)$$

$$\cancel{\nabla_x C(\hat{x},\beta)^T dx} + \partial_\beta C(\hat{x}(\beta);\beta) d\beta$$

Derivative w.r.t $\tau$ :  $\quad -\frac{\beta}{2\tau^2}\|\hat{v}(g)\|^2$

Derivative w.r.t $\beta$ :  $\quad \frac{1}{2\tau}\|\hat{v}(g)\|^2 - \langle g, \hat{v}(g)\rangle$

(f) Write $\frac{1}{n}\mathbb{E}[\|\widehat{u}\|^2]$ and $\frac{1}{n}\mathbb{E}[\langle g, \widehat{u}\rangle]$ as an expectation over the random variables $(\Theta, Z) \sim \widehat{\mu}_{\theta_0} \otimes \mathsf{N}(0,1)$. For the latter, rewrite it using Gaussian integration by parts.

(g) Take the derivative of the objective in Eq. (4) with respect to $\tau$ and with respect to $\beta$. Show that setting the expectations of these derivatives to 0 is equivalent to the fixed point equations.

$$\frac{1}{n}\mathbb{E}\left[\|\widehat{u}(g)\|^2\right]$$

$$= \frac{d}{n} \cdot \frac{1}{d}\mathbb{E}_g\left[\|\,\eta(\widehat{\theta}_b + \tau g;\, \lambda\tau/\beta) - \widehat{\theta}_0\|^2\,\right]$$

$$\underbrace{\hspace{5cm}}$$

$$\sum_{j=1}^{d}\left(\eta(\widehat{\theta}_{0j} + \tau g_j;\, \lambda\tau/g) - \widehat{\theta}_{0j}\right)^2$$

$$\approx \frac{d}{n}\mathbb{E}_{\widehat{\mu}}\left[\left(\eta(\Theta + \tau Z;\, \lambda\tau/\beta) - \Theta\right)^2\right]$$

$$\widehat{\mu} = \frac{1}{d}\sum_{j=1}^{d}\delta_{(\widehat{\theta}_{0j},\, g_j)}$$



$$\xrightarrow{\ \ w\ \ } \mu_\Theta \otimes \mathsf{N}(0,1)$$

$$\frac{1}{n}\mathbb{E}\left[\langle g, \eta(\widehat{\theta}_0 + \tau g;\, \lambda\tau/\beta)\rangle\right]$$

$$\frac{d}{n}\mathbb{E}_{\widehat{\mu}}\left[Z\,\eta(\Theta + \tau Z;\, \lambda\tau/\beta)\right]$$

$$\frac{d}{n}\tau\, \mathbb{E}\left[1\left\{|(h)+\tilde{z}\tilde{z}|\geq \lambda\tau/\beta\right\}\right] = \frac{d}{n}\tau\, \mathbb{P}\left(|(h)+\tilde{z}\tilde{z}|\geq \frac{\lambda\tau}{\beta}\right)$$

## Problem 2: Gordon's objective for max-margin classification

The use of Gordon's technique extends well beyond linear models. In logistic regression, we receive iid samples according to

$$y_i \sim \mathsf{Rad}(f(\langle \boldsymbol{x}_i, \boldsymbol{\theta}_0\rangle)), \quad \boldsymbol{x}_i \sim \mathsf{N}(0, \mathbf{I}_d),$$

$$f(x) = \frac{\exp(x)}{\exp(x) + \exp(-x)}.$$

For a certain $\delta^* > 0$ the following occurs: when $n/p \to \delta < \delta^*$, $n, p \to \infty$, with high probability there exists $\boldsymbol{\theta}$ such that $y_i \langle \boldsymbol{x}_i, \boldsymbol{\theta}\rangle > 0$ for all $i = 1, \ldots, n$. In such a regime, the data is *linearly separable*. The max-margin classifier is defined as

$$\widehat{\boldsymbol{\theta}} \in \arg\max_{\boldsymbol{\theta}} \left\{\min_{i \leq n} y_i \langle \boldsymbol{x}_i, \boldsymbol{\theta}\rangle : \|\boldsymbol{\theta}\| \leq 1\right\}, \tag{5}$$

and the value of the optimization problem, which we denote by $\kappa(\boldsymbol{y}, \boldsymbol{X})$, is called the *maximum margin*. To simplify notation, we assume in this problem that $\|\boldsymbol{\theta}_0\| = 1$. In this problem, we outline how to set up Gordon's problem for max-margin classification. The analysis of Gordon's objective is complicated, and we do not describe it here. See Montanari, Ruan, Sohn, Yan (2019+). "The generalization error of max-margin linear classifiers: High-dimensional asymptotics in the overparameterized regime." `arxiv:1911.01544`.

(a) Show that

$$\kappa(\boldsymbol{y}, \boldsymbol{X}) \geq \kappa \text{ if and only if } \min_{\|\boldsymbol{\theta}\| \leq 1} \|(\kappa\boldsymbol{1} - (\boldsymbol{y} \odot \boldsymbol{X}\boldsymbol{\theta}))_+\|_2 = 0.$$

Argue that

$$\min_{\|\boldsymbol{\theta}\| \leq 1} \frac{1}{\sqrt{d}} \|(\kappa\boldsymbol{1} - (\boldsymbol{y} \odot \boldsymbol{X}\boldsymbol{\theta}))_+\|_2 = \min_{\|\boldsymbol{\theta}\| \leq 1} \max_{\|\boldsymbol{\lambda}\| \leq 1, \boldsymbol{\lambda} \geq 0} \frac{1}{\sqrt{d}} \boldsymbol{\lambda}^\top (\kappa\boldsymbol{1} - \boldsymbol{y} \odot \boldsymbol{X}\boldsymbol{\theta})$$

$$= \min_{\|\boldsymbol{\theta}\| \leq 1} \max_{\|\boldsymbol{\lambda}\| \leq 1, \boldsymbol{\lambda} \odot \boldsymbol{y} \geq 0} \frac{1}{\sqrt{d}} \boldsymbol{\lambda}^\top (\kappa\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}).$$

(b) Why can't we use Gordon's inequality to compare the preceding min-max problem to

$$\min_{\|\boldsymbol{\theta}\| \leq 1} \max_{\|\boldsymbol{\lambda}\| \leq 1, \boldsymbol{y} \odot \boldsymbol{\lambda} \geq 0} \frac{1}{\sqrt{d}} (\kappa\boldsymbol{\lambda}^\top \boldsymbol{y} + \|\boldsymbol{\lambda}\| \boldsymbol{g}^\top \boldsymbol{\theta} + \|\boldsymbol{\theta}\| \boldsymbol{h}^\top \boldsymbol{\lambda})?$$

(c) Let $\tilde{\boldsymbol{x}} = \boldsymbol{X}\boldsymbol{\theta}_0$. Show that the min-max problem is equivalent to

$$\min_{\|\boldsymbol{\theta}\|\leq 1} \max_{\|\boldsymbol{\lambda}\|\leq 1, \boldsymbol{\lambda}\geq 0} \frac{1}{\sqrt{d}}\boldsymbol{\lambda}^\top(\kappa\mathbf{1} - (\boldsymbol{y}\odot\tilde{\boldsymbol{x}})\langle\boldsymbol{\theta}_0, \boldsymbol{\theta}\rangle - \boldsymbol{y}\odot\boldsymbol{X}\Pi_{\boldsymbol{\theta}_0^\perp}\boldsymbol{\theta}).$$

Here $\Pi_{\boldsymbol{\theta}_0^\perp}$ is the projection operator onto the orthogonal complement of the space spaned by $\boldsymbol{\theta}_0$. Argue that

$$\mathbb{P}\left(\min_{\|\boldsymbol{\theta}\|\leq 1} \max_{\|\boldsymbol{\lambda}\|\leq 1, \boldsymbol{\lambda}\geq 0} \frac{1}{\sqrt{d}}\boldsymbol{\lambda}^\top(\kappa\mathbf{1} - (\boldsymbol{y}\odot\tilde{\boldsymbol{x}})\langle\boldsymbol{\theta}_0, \boldsymbol{\theta}\rangle - \boldsymbol{y}\odot\boldsymbol{X}\Pi_{\boldsymbol{\theta}_0^\perp}\boldsymbol{\theta}) \leq t\right)$$

$$\leq 2\mathbb{P}\left(\min_{\|\boldsymbol{\theta}\|\leq 1} \max_{\|\boldsymbol{\lambda}\|\leq 1, \boldsymbol{\lambda}\geq 0} \frac{1}{\sqrt{d}}\left(\boldsymbol{\lambda}^\top(\kappa\mathbf{1} - (\boldsymbol{y}\odot\tilde{\boldsymbol{x}})\langle\boldsymbol{\theta}_0, \boldsymbol{\theta}\rangle) + \|\boldsymbol{\lambda}\|\boldsymbol{g}^\top\Pi_{\boldsymbol{\theta}_0^\perp}\boldsymbol{\theta} + \|\Pi_{\boldsymbol{\theta}_0^\perp}\boldsymbol{\theta}\|\boldsymbol{h}^\top\boldsymbol{\lambda}\right) \leq t\right),$$

and likewise for the comparison of the probabilities that the min-max values exceed $t$, where $\boldsymbol{g} \sim \mathsf{N}(0, \mathbf{I}_d)$ and $\boldsymbol{h} \sim \mathsf{N}(0, \mathbf{I}_n)$ independent of everything else.

(d) What is the limit in Wasserstein-2 distance of $\frac{1}{n}\sum_{i=1}^n \delta_{(y_i, \tilde{x}_i, h_i)}$?