

# **DSO Screening Exam: 2017 In-Class Exam**

Gregory Faletto

**Exercise 1 (Probability/Analysis).** (a)

$$\begin{aligned}
 \dot{Z}(t) &= \frac{\partial}{\partial t} \sum_{i=1}^k \gamma_i(t) X_i = \sum_{i=1}^k \dot{\gamma}_i(t) X_i \\
 \implies \mathbb{E}(\dot{Z}(t)) &= \sum_{i=1}^k \mathbb{E}(\dot{\gamma}_i(t) X_i) = 0 \\
 \implies \text{Cov}(Z(t), \dot{Z}(t)) &= \mathbb{E}[(Z(t) - \mathbb{E}[Z(t)])(\dot{Z}(t) - \mathbb{E}[\dot{Z}(t)])] = \mathbb{E}[Z(t)\dot{Z}(t)] \\
 &= \mathbb{E}\left[\left(\sum_{i=1}^k \gamma_i(t) X_i\right) \left(\sum_{i=1}^k \dot{\gamma}_i(t) X_i\right)\right] = \sum_{i=1}^k \mathbb{E}(\gamma_i(t) \dot{\gamma}_i(t) X_i^2) + 0 = \sum_{i=1}^k \gamma_i(t) \dot{\gamma}_i(t)
 \end{aligned}$$

since  $\mathbb{E}(X_i^2) = 1$ . But

$$\sum_{i=1}^k \gamma_i \dot{\gamma}_i(t) = 0 \tag{1}$$

because

$$\sum_{i=1}^k \gamma_i^2(t) = 1 \iff \frac{\partial}{\partial t} \left( \sum_{i=1}^k \gamma_i^2(t) \right) = 0 \iff 2 \sum_{i=1}^k \gamma_i(t) \dot{\gamma}_i(t) = 0,$$

so the conclusion follows.

(b)

$$\begin{aligned}
 Z(t) &= \sum_{i=1}^k \gamma_i(t) X_i \implies Z(t) \sim \mathcal{N}\left(0, \sum_{i=1}^k \gamma_i^2(t)\right) = \mathcal{N}(0, 1) \\
 \ddot{Z}(t) &= \sum_{i=1}^k \ddot{\gamma}_i(t) X_i \implies \ddot{Z}(t) \sim \mathcal{N}\left(0, \sum_{i=1}^k \ddot{\gamma}_i^2(t)\right)
 \end{aligned}$$

Also, we have

$$\begin{aligned}
 \text{Cov}(Z(t), \ddot{Z}(t)) &= \mathbb{E}[(Z(t) - \mathbb{E}[Z(t)])(\ddot{Z}(t) - \mathbb{E}[\ddot{Z}(t)])] = \mathbb{E}[Z(t)\ddot{Z}(t)] \\
 &= \mathbb{E}\left[\left(\sum_{i=1}^k \gamma_i(t) X_i\right) \left(\sum_{i=1}^k \ddot{\gamma}_i(t) X_i\right)\right] = \sum_{i=1}^k \mathbb{E}(\gamma_i(t) \ddot{\gamma}_i(t) X_i^2) + 0 = \sum_{i=1}^k \gamma_i(t) \ddot{\gamma}_i(t) = - \sum_{i=1}^k \dot{\gamma}_i^2(t)
 \end{aligned}$$

because differentiating (1) yields

$$\frac{\partial}{\partial t} \sum_{i=1}^k \gamma_i \dot{\gamma}_i(t) = 0 \iff \sum_{i=1}^k (\dot{\gamma}_i^2(t) + \gamma_i(t) \ddot{\gamma}_i(t)) = 0 \iff \sum_{i=1}^k \gamma_i(t) \ddot{\gamma}_i(t) = - \sum_{i=1}^k \dot{\gamma}_i^2(t).$$

Since both distributions are Gaussian, we have

$$\mathbb{E}(Z(t) \mid \ddot{Z}(t)) = \mathbb{E}[Z(t)] + \frac{\text{Cov}[Z(t), \ddot{Z}(t)]}{\text{Var}[\ddot{Z}(t)]} (\ddot{Z}(t) - \mathbb{E}[\ddot{Z}(t)]) = \left[ \sum_{i=1}^k \ddot{\gamma}_i^2(t) \right]^{-1} \cdot \left[ - \sum_{i=1}^k \dot{\gamma}_i^2(t) \right] \ddot{Z}(t).$$

**Exercise 2 (Mathematical statistics; hypothesis test. so don't worry about).** (a)

(b)

**Exercise 3 (Probability; Wen says we should do this).** (a) In the  $n = 1$  case we have

$$I(m, \lambda) = \int_{\mathbb{R}} \frac{\exp(\lambda \theta_1 m_1)}{(1 + \exp(\theta_1))^\lambda} d\theta_1 = \int_{\mathbb{R}} \frac{\exp(\theta_1)^{\lambda m_1}}{(1 + \exp(\theta_1))^\lambda} d\theta_1$$

$$\text{Let } u_1 = \exp(\theta_1) \implies du_1 = \exp(\theta_1) d\theta_1 \iff d\theta_1 = \frac{1}{u_1} du_1.$$

$$\implies I(m, \lambda) = \int_0^\infty \frac{u_1^{\lambda m_1}}{(u_1 + 1)^\lambda} \frac{du_1}{u_1} = \int_0^\infty u_1^{\lambda m_1 - 1} (u_1 + 1)^{-\lambda} du_1$$

$$\text{Let } x = (u_1 + 1)^{-1} \iff u_1 = x^{-1} - 1 \implies dx = -(u_1 + 1)^{-2} du_1 \iff du_1 = -x^{-2} dx.$$

$$\begin{aligned} \implies I(m, \lambda) &= \int_{x=1}^0 \left( \frac{1}{x} - 1 \right)^{\lambda m_1 - 1} x^\lambda (-x^{-2}) dx = \int_0^1 (1 - x)^{\lambda m_1 - 1} x^{\lambda - 2 - \lambda m_1 + 1} dx \\ &= \int_0^1 x^{\lambda(1 - m_1) - 1} (1 - x)^{\lambda m_1 - 1} dx = \frac{\Gamma(\lambda(1 - m_1))\Gamma(\lambda m_1)}{\Gamma(\lambda[1 - m_1] + \lambda m_1)} = \frac{\Gamma(\lambda m_0)\Gamma(\lambda m_1)}{\Gamma(\lambda)} \end{aligned}$$

where we used the definition of the Beta function:

$$B(z_1, z_2) := \int_0^1 t^{z_1 - 1} (1 - t)^{z_2 - 1} dt = \frac{\Gamma(z_1)\Gamma(z_2)}{\Gamma(z_1 + z_2)}.$$

Now we do an induction on  $n$ . Suppose for some  $n \in \mathbb{N}$ , we have

$$\int_{\mathbb{R}^n} \exp\left(\lambda \left[ \sum_{i=1}^n \theta_i x_i \right]\right) \left(1 + \sum_{i=1}^n \exp(\theta_i)\right)^{-\lambda} d\theta = \frac{\Gamma(\lambda m_0)\Gamma(\lambda m_1) \cdots \Gamma(\lambda m_n)}{\Gamma(\lambda)}. \quad (2)$$

Then we will evaluate

$$\int_{\mathbb{R}^{n+1}} \exp\left(\lambda \left[ \sum_{i=1}^{n+1} \theta_i m_i \right]\right) \left(1 + \sum_{i=1}^{n+1} \exp(\theta_i)\right)^{-\lambda} d\theta = \int_{\mathbb{R}^{n+1}} \prod_{i=1}^{n+1} \exp(\theta_i)^{\lambda m_i} \left(1 + \sum_{i=1}^{n+1} \exp(\theta_i)\right)^{-\lambda} d\theta$$

$$\text{Let } u_i = \exp(\theta_i) \implies du_i = \exp(\theta_i) d\theta_i \iff d\theta_i = \frac{1}{u_i} du_i \iff d\theta = \prod_{i=1}^{n+1} \frac{1}{u_i} du.$$

$$= \int_{\mathbb{R}_+^{n+1}} \prod_{i=1}^{n+1} u_i^{\lambda m_i} \left(1 + \sum_{i=1}^{n+1} u_i\right)^{-\lambda} \prod_{i=1}^{n+1} \frac{1}{u_i} du = \int_{\mathbb{R}_+^{n+1}} \prod_{i=1}^{n+1} u_i^{\lambda m_i - 1} \left(1 + \sum_{i=1}^{n+1} u_i\right)^{-\lambda} du$$

$$\text{Let } x_i = \left(u_i + \frac{1}{n+1}\right)^{-1} \iff u_i = \frac{1}{x_i} - \frac{1}{n+1} \implies dx_i = -\left(u_i + \frac{1}{n+1}\right)^{-2} du_i \iff du_i = -x_i^{-2} dx_i.$$

$$\begin{aligned} &= \int_{x_{n+1}=1}^0 \int_{x_n=1}^0 \cdots \int_{x_1=1}^0 \prod_{i=1}^{n+1} \left(\frac{1}{x_i} - \frac{1}{n+1}\right)^{\lambda m_i - 1} \left(1 + \sum_{i=1}^{n+1} \left[\frac{1}{x_i} - \frac{1}{n+1}\right]\right)^{-\lambda} \prod_{i=1}^{n+1} -x_i^{-2} dx_1 \cdots dx_n dx_{n+1} \\ &= \int_{x_{n+1}=0}^1 \int_{x_n=0}^1 \cdots \int_{x_1=0}^1 \prod_{i=1}^{n+1} \left(\frac{1}{x_i} - \frac{1}{n+1}\right)^{\lambda m_i - 1} \left(\sum_{i=1}^{n+1} \frac{1}{x_i}\right)^{-\lambda} \prod_{i=1}^{n+1} x_i^{-2} dx_1 \cdots dx_n dx_{n+1} \end{aligned}$$

$$\begin{aligned}
&= \int_{x_{n+1}=0}^1 \int_{x_n=0}^1 \cdots \int_{x_1=0}^1 \prod_{i=1}^n \left( \frac{1}{x_i} - \frac{1}{n+1} \right)^{\lambda m_i - 1} \left( \sum_{i=1}^{n+1} \frac{1}{x_i} \right)^{-\lambda} \prod_{i=1}^n x_i^{-2} dx_1 \cdots dx_n \\
&\quad \cdot \left( \frac{1}{x_{n+1}} - \frac{1}{n+1} \right)^{\lambda m_{n+1} - 1} x_{n+1}^{-2} dx_{n+1}
\end{aligned} \tag{3}$$

We can rewrite the left side of the inductive hypothesis (2) using analagous substitutions as

$$\begin{aligned}
&\int_{\mathbb{R}^n} \exp \left( \lambda \left[ \sum_{i=1}^n \theta_i x_i \right] \right) \left( 1 + \sum_{i=1}^n \exp(\theta_i) \right)^{-\lambda} d\theta = \int_{\mathbb{R}_+^n} \prod_{i=1}^n u_i^{\lambda m_i} \left( 1 + \sum_{i=1}^n u_i \right)^{-\lambda} \prod_{i=1}^n u_i^{-1} du \\
&= \int_{x_n=1}^0 \cdots \int_{x_1=1}^0 \prod_{i=1}^n \left( \frac{1}{x_i} - \frac{1}{n} \right)^{\lambda m_i - 1} \left( 1 + \sum_{i=1}^n \left[ \frac{1}{x_i} - \frac{1}{n} \right] \right)^{-\lambda} \prod_{i=1}^n -x_i^{-2} dx_1 \cdots dx_n \\
&= \int_{x_n=0}^1 \cdots \int_{x_1=0}^1 \prod_{i=1}^n \left( \frac{1}{x_i} - \frac{1}{n} \right)^{\lambda m_i - 1} \left( \sum_{i=1}^n \frac{1}{x_i} \right)^{-\lambda} \prod_{i=1}^n x_i^{-2} dx_1 \cdots dx_n \\
&\quad \vdots \\
&\dots = \frac{\Gamma(\lambda m_0) \Gamma(\lambda m_1) \cdots \Gamma(\lambda m_{n+1})}{\Gamma(\lambda)}.
\end{aligned}$$

(b) **not sure about the last step on this problem**

$$\begin{aligned}
\lim_{\lambda \rightarrow 0} \lambda^n I(m, \lambda) &= \lim_{\lambda \rightarrow 0} \lambda^n \cdot \frac{\Gamma(\lambda m_0) \Gamma(\lambda m_1) \cdots \Gamma(\lambda m_n)}{\Gamma(\lambda)} = \lim_{\lambda \rightarrow 0} \frac{\Gamma(\lambda m_0)}{\Gamma(\lambda)} \cdot [\lambda \Gamma(\lambda m_1)] \cdots [\lambda \Gamma(\lambda m_n)] \\
&= \lim_{\lambda \rightarrow 0} \frac{\Gamma(\lambda m_0)}{\Gamma(\lambda)} \cdot \frac{\Gamma(\lambda m_1)}{\Gamma(\lambda)} \cdots \frac{\Gamma(\lambda m_n)}{\Gamma(\lambda)} \cdot [\lambda \Gamma(\lambda)] \cdots [\lambda \Gamma(\lambda)] \\
&= \frac{\lim_{\lambda \rightarrow 0} \lambda \Gamma(\lambda m_0)}{\lim_{\lambda \rightarrow 0} \lambda \Gamma(\lambda)} \cdot \frac{\lim_{\lambda \rightarrow 0} \lambda \Gamma(\lambda m_1)}{\lim_{\lambda \rightarrow 0} \lambda \Gamma(\lambda)} \cdots \frac{\lim_{\lambda \rightarrow 0} \lambda \Gamma(\lambda m_n)}{\lim_{\lambda \rightarrow 0} \lambda \Gamma(\lambda)} = 1.
\end{aligned}$$

(c) **again, not sure about last part** Stirling's formula:

$$\begin{aligned}
\lambda! &\sim \lambda^\lambda e^{-\lambda} \sqrt{2\pi\lambda} = \lambda^{\lambda+1/2} e^{-\lambda} \sqrt{2\pi} \iff \lim_{\lambda \rightarrow \infty} \frac{\Gamma(\lambda+1) e^\lambda}{\lambda^{\lambda+1/2} \sqrt{2\pi}} = 1 \iff \lim_{\lambda \rightarrow \infty} \frac{\lambda \Gamma(\lambda) e^\lambda}{\lambda^{\lambda+1/2} \sqrt{2\pi}} = 1 \\
&\iff \lim_{\lambda \rightarrow \infty} \frac{\lambda^{1/2} \Gamma(\lambda) e^\lambda}{\lambda^\lambda} = \sqrt{2\pi} \iff \lim_{\lambda \rightarrow \infty} \frac{\lambda^{n/2} (\Gamma(\lambda))^n e^{n\lambda}}{\lambda^{n\lambda}} = (2\pi)^{n/2} \\
\lim_{\lambda \rightarrow \infty} \lambda^{n/2} I(m, \lambda) &= \lim_{\lambda \rightarrow \infty} \lambda^{n/2} \cdot \frac{\Gamma(\lambda m_0) \Gamma(\lambda m_1) \cdots \Gamma(\lambda m_n)}{\Gamma(\lambda)} \\
&= \lim_{\lambda \rightarrow \infty} \frac{\lambda^{n/2} (\Gamma(\lambda))^n e^{n\lambda}}{\lambda^{n\lambda}} \frac{\Gamma(\lambda m_0) \Gamma(\lambda m_1) \cdots \Gamma(\lambda m_n)}{(\Gamma(\lambda))^{n+1}} \frac{\lambda^{n\lambda}}{e^{n\lambda}} \\
&= \lim_{\lambda \rightarrow \infty} \left( \frac{\lambda^{n/2} (\Gamma(\lambda))^n e^{n\lambda}}{\lambda^{n\lambda}} \right) \lim_{\lambda \rightarrow \infty} \left( \frac{\Gamma(\lambda m_0) \Gamma(\lambda m_1) \cdots \Gamma(\lambda m_n)}{(\Gamma(\lambda))^{n+1}} \right) \lim_{\lambda \rightarrow \infty} \left( \frac{\lambda}{e} \right)^{n\lambda} \\
&= (2\pi)^{n/2} \cdot 1 \cdot \infty = \infty.
\end{aligned}$$

(d)

**Definition 1 (Exponential Family).** Let  $n, k$  be positive integers and let  $\mu$  be a measure on  $\mathbb{R}^n$ . Let  $t_1, \dots, t_k : \mathbb{R}^n \rightarrow \mathbb{R}$ . Let  $h : \mathbb{R}^n \rightarrow [0, \infty]$ , and assume  $h$  is not identically zero. For any  $w = (w_1, \dots, w_k) \in \mathbb{R}^k$ , define

$$a(w) := \log \left[ \int_{\mathbb{R}^n} h(x) \exp \left( \sum_{i=1}^k w_i t_i(x) \right) d\mu(x) \right], \quad \forall x \in \mathbb{R}^n$$

The set  $\{w \in \mathbb{R}^k\}$  is called the **natural parameter space**. On this set, the function

$$f_w(x) := h(x) \exp \left( \sum_{i=1}^k w_i t_i(x) - a(w) \right), \quad \forall x \in \mathbb{R}^n$$

satisfies  $\int_{\mathbb{R}^n} f_w(x) d\mu(x) = 1$  (by the definition of  $a(w)$ ). The set of functions (which can be interpreted as probability density functions, or as probability mass functions according to  $\mu$ )  $\{f_w : w \in \Theta : a(w) < \infty\}$  is called a  **$k$ -parameter exponential family in canonical form**.

Let  $X_{p_0, \dots, p_j} : \Omega \rightarrow \mathbb{R}^n$  be the random variable described by this distribution, parameterized by  $p_0, \dots, p_j$  and with  $\delta_{e_0}, \dots, \delta_{e_n}$  fixed (otherwise this is not an exponential family, for reasons described below). Then we have

$$f_{p_0, \dots, p_j}(x) = \mathbb{P}(X_{p_0, \dots, p_j} = x) = \begin{cases} p_j / \sum_{i=0}^n p_i & x = \delta_{e_j} \\ 0 & \text{otherwise} \end{cases} = I(x = \delta_{e_k}) p_k, k \in \{0, 1, \dots, n\}.$$

Let

$$h(x) := \begin{cases} 1 & x \in \{\delta_{e_0}, \dots, \delta_{e_n}\} \\ 0 & \text{otherwise.} \end{cases}$$

Let  $k = n + 1$ , and let  $t_i(x) := I(x = \delta_{e_{i-1}})$  (an indicator function that equals 1 if  $x = \delta_{e_{i-1}}$  and 0 otherwise). Let  $w_i := \log(p_{i-1})$ , and let  $a(w) := 0$ . Then

$$\begin{aligned} h(x) \exp \left( \sum_{i=1}^k w_i t_i(x) - a(w) \right) &= I(x \in \{\delta_{e_0}, \dots, \delta_{e_n}\}) \exp \left( \sum_{i=1}^k I(x = \delta_{e_{i-1}}) \log(p_{i-1}) \right) \\ &= I(x = \delta_{e_{i-1}}) p_{i-1}, i \in \{1, \dots, n+1\} = I(x = \delta_{e_k}) p_k, k \in \{0, 1, \dots, n\} = f_{p_0, \dots, p_j}(x), \end{aligned}$$

so  $F$  is an exponential family. (Where the mass function equals 0 depends on  $\delta_{e_0}, \dots, \delta_{e_n}$ . So if these are not fixed and are instead parameters of the distribution, then where the mass function equals 0 depends on the parameters of the distribution, which is not allowed in an exponential family;  $h(x)$  determines where the mass function equals 0, and it may only depend on  $x$ , not the parameters of the distribution.)

**stuff with conjugate prior has to do with Bayesian statistics; not our concern.**

**Exercise 4 (Convex Optimization).** Like theorem 2 in convex optimization lecture notes 3; probably don't need to worry about since not covered in Boyd. Don't prioritize.

**Exercise 5 (High-dimensional statistics).** (a) If  $p > n$ , then even if  $\mathbf{X}$  is full rank, it has a nullspace with nonzero entries. That is, the columns of  $\mathbf{X}$  must be linearly dependent, so there are infinitely many non-zero vectors  $\mathbf{v}$  such that  $\mathbf{X}\mathbf{v} = \mathbf{0}$ . Therefore for any solution  $\hat{\beta}$  satisfying

$$\hat{\beta} \in \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

and any vector  $\mathbf{v}$  in the nullspace of  $\mathbf{X}$ ,  $\hat{\beta} + \mathbf{v}$  is also a solution since

$$\min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 = \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 = \|\mathbf{y} - \mathbf{X}\hat{\beta} - \mathbf{0}\|_2^2 = \|\mathbf{y} - \mathbf{X}(\hat{\beta} + \mathbf{v})\|_2^2$$

$$\iff \hat{\beta} + \mathbf{v} \in \arg \min_{\beta \in \mathbb{R}^p} \|\mathbf{y} - \mathbf{X}\beta\|_2^2.$$

(b) If we assume that  $\|\beta_0\|_0 = s \leq n$ , then it may be that there is only one  $\mathbf{v}$  in the nullspace of  $\mathbf{X}$  such that  $\|\hat{\beta}\|_0 = s$ . Intuitively, this makes sense because as long as the features with zero coefficients in  $\beta_0$  are sufficiently uncorrelated with the features with nonzero coefficients and the true coefficients are not too small, it is very unlikely that it is possible to construct another  $s$ -sparse vector with equally low empirical risk by replacing one of the true features with only one of the other features (or some combination thereof). (We would also hope that  $\text{rank}(\mathbf{X}) > s$ . Otherwise, it could either be the case that features in the true model are linearly dependent, in which case the true model is not identifiable, or some of the noise features are linearly dependent with some of the true features, which could also preclude identifiability.)

(c) **Solution in Section 1.9.1 of Linear Regression notes.**

We will follow the analysis from [Zhao and Yu \[2006\]](#). The lasso problem is convex but not necessarily strictly convex if  $p > n$ . That is, there is some flat region, so the minimizer may not be unique. Consider the KKT conditions from convex optimization:

$$g(\beta) = \arg \min \left\{ \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right\} = \arg \min \{f_1(\beta) + f_2(\beta)\}$$

Then  $\hat{\beta}$  is a lasso solution if and only if 0 is in the subdifferential of  $g(\hat{\beta})$ . Note that

$$\partial g(\hat{\beta}) = \nabla f_1 + \partial f_2 = \frac{1}{n} \mathbf{X}^T (\mathbf{X}\hat{\beta} - \mathbf{y}) + \lambda \begin{bmatrix} \vdots \\ \partial |\hat{\beta}_j| \\ \vdots \end{bmatrix} = \frac{1}{n} \mathbf{X}^T (\mathbf{X}\hat{\beta} - \mathbf{y}) + \lambda \begin{bmatrix} \vdots \\ \begin{cases} \text{sgn}(\hat{\beta}_j) & \hat{\beta}_j \neq 0 \\ [-1, 1] & \hat{\beta}_j = 0 \end{cases} \\ \vdots \end{bmatrix}$$

Now assume  $\text{supp}(\hat{\beta}) = \text{supp}(\beta_0)$  (that is, assume lasso recovers the correct support). Suppose the first  $s$  features are nonzero and consider one of them (so we know that we should have  $\hat{\beta}_j \neq 0$ ):

$$0 \in \partial g(\hat{\beta}) \implies 0 \in \partial_j g(\hat{\beta}) = \left[ \frac{1}{n} \mathbf{X}^T (\mathbf{X}\hat{\beta} - \mathbf{y}) \right]_j + \lambda \text{sgn}(\hat{\beta}_j)$$

Therefore

$$\frac{1}{n} \mathbf{X}_A^T (\mathbf{X}\hat{\beta} - \mathbf{y}) + \lambda \text{sgn}(\hat{\beta}_j) = 0 \tag{4}$$

where  $X_A$  is a submatrix of  $\mathbf{X}$  containing the columns corresponding to the features in the true support, is our first condition. Next, consider what happens for  $j > s$  (features not in the true support). We have

$$\begin{aligned} 0 \in \partial g(\hat{\beta}) &\implies 0 \in \partial_j g(\hat{\beta}) = \left[ \frac{1}{n} X^T (X\hat{\beta} - \mathbf{y}) \right] + \lambda[-1, 1] \\ &\implies \left\| \frac{1}{n} X_{A^c}^T (X\hat{\beta} - \mathbf{y}) \right\|_\infty \leq \lambda \end{aligned} \quad (5)$$

where  $X_{A^c}$  is a submatrix of  $\mathbf{X}$  containing the columns corresponding to the features not in the true support, is our boundary condition. Recall the true model

$$y = X\beta_0 + \varepsilon$$

and consider the case  $X = [\mathbf{X}_1 \quad \mathbf{X}_2]$  where  $\mathbf{X}_1$  are the features in the true model and  $\mathbf{X}_2$  are noise features; that is,  $\beta_0 = \begin{bmatrix} \beta_1 \\ 0 \end{bmatrix}$ . Then we are assuming

$$\hat{\beta}_{\text{lasso}} = \begin{bmatrix} \hat{\beta}_1 \\ 0 \end{bmatrix}.$$

We have from (4)

$$\begin{aligned} 0 &= \frac{1}{n} \mathbf{X}_1^T (X\hat{\beta} - \mathbf{y}) + \lambda \text{sgn}(\hat{\beta}_1) = \frac{1}{n} \mathbf{X}_1^T (\mathbf{X}_1 \hat{\beta}_1 - \mathbf{X}_1 \beta_1 - \varepsilon) + \lambda \text{sgn}(\hat{\beta}_1) \\ &\iff \frac{1}{n} \mathbf{X}_1^T \mathbf{X}_1 (\hat{\beta}_1 - \beta_1) = \frac{1}{n} \mathbf{X}_1^T \varepsilon - \lambda \text{sgn}(\hat{\beta}_1) \end{aligned}$$

Let's assume that  $\text{sgn}(\hat{\beta}) = \text{sgn}(\beta_0)$  (sign consistency).

$$\iff \frac{1}{n} \mathbf{X}_1^T \mathbf{X}_1 (\hat{\beta}_1 - \beta_1) = \frac{1}{n} \mathbf{X}_1^T \varepsilon - \lambda \text{sgn}(\beta_1)$$

which is linear in  $\hat{\beta}$ . Solving, we have

$$\iff \hat{\beta}_1 - \beta_1 = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} (\mathbf{X}_1^T \varepsilon - n\lambda \text{sgn}(\beta_1)) \iff \hat{\beta}_1 = \beta_1 + (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \varepsilon - \lambda \text{sgn}(\beta_1)) \quad (6)$$

Looking at the second (boundary) condition (5), we have

$$\left\| \frac{1}{n} \mathbf{X}_2^T (X\hat{\beta} - \mathbf{y}) \right\|_\infty \leq \lambda. \quad (7)$$

Consider that

$$X\hat{\beta} - \mathbf{y} = \mathbf{X}_1 \hat{\beta}_1 - \mathbf{X}_1 \beta_1 - \varepsilon = \mathbf{X}_1 (\hat{\beta}_1 - \beta_1) - \varepsilon$$

Substituting in the result from (6) yields

$$X\hat{\beta} - \mathbf{y} = \mathbf{X}_1 [(n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \varepsilon - \lambda \text{sgn}(\beta_1))] - \varepsilon$$

which when we plug into (7) yields

$$\begin{aligned} & \left\| \frac{1}{n} \mathbf{X}_2^T \left[ \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1)) - \boldsymbol{\varepsilon} \right] \right\|_{\infty} \leq \lambda. \\ \iff & \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1)) - \frac{1}{n} \mathbf{X}_2^T \boldsymbol{\varepsilon} \right\|_{\infty} \leq \lambda. \end{aligned}$$

Using the Triangle Inequality, we have

$$\begin{aligned} & \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1)) - \frac{1}{n} \mathbf{X}_2^T \boldsymbol{\varepsilon} \right\|_{\infty} \\ & \leq \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1)) \right\|_{\infty} + \left\| \frac{1}{n} \mathbf{X}_2^T \boldsymbol{\varepsilon} \right\|_{\infty} \\ & \leq \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} \right\|_{\infty} \cdot \left\| n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1) \right\|_{\infty} + \left\| \frac{1}{n} \mathbf{X}_2^T \boldsymbol{\varepsilon} \right\|_{\infty} \quad (8) \\ & \vdots \end{aligned}$$

- (d) **Solution in Section 1.9.1 of Linear Regression notes.** Theorem 4 from [Zhao and Yu \[2006\]](#) states that if  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$  for some  $\sigma^2 > 0$  and some other regularity conditions hold, then the Strong Irrepresentable Condition implies the lasso has sign consistency (which implies model selection consistency) with high probability for  $\lambda$  satisfying

$$\lambda(n) \propto n^{(1+c_4)/2}$$

where  $c_4$  is a constant between 0 and 1 and smaller than  $c_2 - c_1$ , where  $c_1$  is a constant governing how large the true model can be with  $0 \leq c_1$  and  $c_2$  is a constant governing the minimum size of the coefficients in the true model with  $c_1 < c_2 \leq 1$ .

Recall that the Strong Irrepresentable Condition is the following: if as in part 5(c)  $\mathbf{X}_1 \in \mathbb{R}^{n \times q}$  is a matrix containing only the  $q$  columns of  $\mathbf{X}$  corresponding to features in the true model and  $\mathbf{X}_2 \in \mathbb{R}^{n \times p-q}$  contains only the  $p - q$  columns of  $\mathbf{X}$  corresponding to irrelevant features, the Strong Irrepresentable Condition is

$$\begin{aligned} & \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 \left( \frac{1}{n} \mathbf{X}_1^T \mathbf{X}_1 \right)^{-1} \operatorname{sgn}(\boldsymbol{\beta}_{(1)}^n) \right\|_{\infty} \leq 1 - \eta \\ \iff & \left\| \left( \mathbf{X}_1^T \mathbf{X}_1 \right)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \right\|_{\infty} \leq 1 - \eta \end{aligned}$$

for some  $\boldsymbol{\eta} \in \mathbb{R}^n$  with  $\boldsymbol{\eta} \succeq \mathbf{0}$ , where  $\boldsymbol{\beta}_{(1)}^n$  are the coefficients of  $\mathbf{X}_1$  in the true model. That is, a regression of the variables in  $\mathbf{X}_2$  against the variables in  $\mathbf{X}_1$  may not have any coefficients that are larger in absolute value than  $1 - \eta$  for some  $\eta > 0$ . This implies that the features in  $\mathbf{X}_2$  are not too strongly correlated with the features in  $\mathbf{X}_1$  (the features of  $\mathbf{X}_2$  are “irrepresentable” by the features in  $\mathbf{X}_1$ ).

(Roughly speaking, the required regularity conditions are that the number of features in the design matrix  $p_n$  is not excessively large, the eigenvalues of  $\mathbf{X}_1$  are not too small, the coefficients of the relevant features in the true model are not too small, and the model is sufficiently sparse. The conditions of this theorem allow  $p_n$  to grow asymptotically with  $n$ ; the requirements for lower bounded eigenvalues and coefficients on the true coefficients are consequences of allowing  $p_n$  to grow with  $n$ .)



## References

P. Zhao and B. Yu. On Model Selection Consistency of Lasso. *Journal of Machine Learning Research*, 7: 2541–2563, 2006.