

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

Gregory Faletto

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

$$\dot{Z}(t) = \frac{\partial}{\partial t} \sum_{i=1}^k \gamma_i(t) X_i = \sum_{i=1}^k \frac{\partial \gamma_i(t)}{\partial t} X_i$$

$$\implies \mathbb{E}(\dot{Z}(t)) = \sum_{i=1}^k \mathbb{E}\left(\frac{\partial \gamma_i(t)}{\partial t} X_i\right) = 0$$

$$\implies \text{Cov}(Z(t), \dot{Z}(t)) = \mathbb{E}[(Z(t) - \mathbb{E}[Z(t)])(\dot{Z}(t) - \mathbb{E}[\dot{Z}(t)])] = \mathbb{E}[Z(t)\dot{Z}(t)]$$

$$= \mathbb{E}\left[\left(\sum_{i=1}^k \gamma_i(t) X_i\right) \left(\sum_{i=1}^k \frac{\partial \gamma_i(t)}{\partial t} X_i\right)\right] = \sum_{i=1}^k \mathbb{E}\left(\gamma_i(t) \frac{\partial \gamma_i(t)}{\partial t} X_i^2\right) + 0 = \sum_{i=1}^k \gamma_i(t) \frac{\partial \gamma_i(t)}{\partial t}$$

$$\vdots$$

$$\mathbb{E}\left[\sum_i \gamma_i \frac{d\gamma_i}{dt}\right] = 0$$

because

(b)

$$\ddot{Z}(t) = \sum_{i=1}^k \frac{\partial^2 \gamma_i(t)}{\partial t^2} X_i$$

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

(b)

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

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

$$\text{Let } u_1 = 1 + \exp(\theta_1) \implies \partial u_1 = \exp(\theta_1) \partial \theta_1.$$

$$= \int_{\mathbb{R}} \frac{(u_1 - 1)^{\lambda m_1 - 1}}{u_1^\lambda} du_1$$

$$= \dots = \frac{\Gamma(\lambda(1 - m_1))\Gamma(\lambda m_1)}{\Gamma(\lambda)} = \frac{\Gamma(\lambda m_0)\Gamma(\lambda m_1)}{\Gamma(\lambda)}$$

(b)

**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

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

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

$$\begin{aligned} 0 \in \partial g(\hat{\beta}) &\implies 0 \in \partial_j g(\hat{\beta}) = \left[ \frac{1}{n} X^T (X\hat{\beta} - \mathbf{y}) \right] + \lambda[-1, 1] \\ &\implies \left\| \frac{1}{n} X_{A^c}^T (X\hat{\beta} - \mathbf{y}) \right\|_{\infty} \leq \lambda \end{aligned} \quad (2)$$

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 (1)

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

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

$$\iff \hat{\beta}_1 - \beta_1 = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} (\mathbf{X}_1^T \varepsilon - n\lambda \text{sgn}(\beta_1)) \iff \hat{\beta}_1 = \beta_1 + (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \varepsilon - \lambda \text{sgn}(\beta_1)) \quad (3)$$

Looking at the second (boundary) condition (2), we have

$$\left\| \frac{1}{n} \mathbf{X}_2^T (X\hat{\beta} - \mathbf{y}) \right\|_{\infty} \leq \lambda. \quad (4)$$

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 (3) 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 (4) 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

$$\begin{aligned} & \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1)) - \frac{1}{n} \mathbf{X}_2^T \boldsymbol{\varepsilon} \right\|_{\infty} \\ & \leq \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} (n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1)) \right\|_{\infty} + \left\| \frac{1}{n} \mathbf{X}_2^T \boldsymbol{\varepsilon} \right\|_{\infty} \\ & \leq \left\| \frac{1}{n} \mathbf{X}_2^T \mathbf{X}_1 (n^{-1} \mathbf{X}_1^T \mathbf{X}_1)^{-1} \right\|_{\infty} \cdot \left\| n^{-1} \mathbf{X}_1^T \boldsymbol{\varepsilon} - \lambda \operatorname{sgn}(\boldsymbol{\beta}_1) \right\|_{\infty} + \left\| \frac{1}{n} \mathbf{X}_2^T \boldsymbol{\varepsilon} \right\|_{\infty} \quad (5) \\ & \vdots \end{aligned}$$

- (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.