

Supplementary material for feature subset selection for the multinomial logit model via mixed-integer optimization

Shunsuke Kamiya Tokyo University of Agriculture and Technology
 Ryuhei Miyashiro Tokyo University of Agriculture and Technology
 Yuichi Takano University of Tsukuba

A Key assumptions and additional lemmas

In the main body of the paper, the following conditions are assumed:

1. the loss function $\ell(y, \cdot)$ is proper convex for all $y \in [m]$, and (P) is bounded and has an interior feasible solution;
2. the Fenchel conjugate $\hat{\ell}(y, \cdot)$ is continuous on its effective domain and strictly convex for all $y \in [m]$;
3. $\text{dom } \hat{\ell}(y, \cdot)$ is nonempty, bounded and closed for all $y \in [m]$.

Before proving our theorems, we must introduce the following lemmas.

Lemma 12 (Strong duality; Theorem 9.6 in [6]). *Let $\mathbf{a}_i \in \mathbb{R}^n$ and $b_i \in \mathbb{R}$ for all $i \in [m]$. Let $f : \mathbb{R}^n \rightarrow [-\infty, +\infty]$ be the objective function of the following problem:*

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{a}_i^\top \mathbf{x} - b_i = 0, && \forall i \in [m]. \end{aligned}$$

We assume that f is proper convex and that there is a feasible solution on $\text{ri dom } f$, which denotes the relative interior of $\text{dom } f$. Moreover, if this optimization problem is bounded, the optimal values of the following two optimization problems are the same.

$$\begin{aligned} & \underset{\boldsymbol{\lambda} \in \mathbb{R}^m}{\text{maximize}} && \min \{ \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \mid \mathbf{x} \in \mathbb{R}^n \}, \\ & \underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} && \max \{ \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \mid \boldsymbol{\lambda} \in \mathbb{R}^m \}, \end{aligned}$$

where $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ is the Lagrange function

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i (\mathbf{a}_i^\top \mathbf{x} - b_i).$$

Lemma 13 (Section 3.5 in [4]). *Given $\mathcal{A} \subseteq \mathbb{R}^m$, let $f_\alpha : \mathbb{R}^n \rightarrow \mathbb{R}$ be subdifferentiable functions for all $\alpha \in \mathcal{A}$. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be the function defined by*

$$f(\mathbf{x}) = \sup_{\alpha \in \mathcal{A}} f_\alpha(\mathbf{x}).$$

Then f is subdifferentiable if \mathcal{A} is compact and the function $\alpha \mapsto f_\alpha(\mathbf{x})$ is upper semi-continuous for each \mathbf{x} . The subderivative of f is given by

$$\partial f(\mathbf{x}) = \text{conv} \left[\bigcup \{ \partial f_\alpha(\mathbf{x}) \mid \alpha \in \mathcal{A}, f_\alpha(\mathbf{x}) = f(\mathbf{x}) \} \right],$$

where conv denotes the convex hull of a set.

Lemma 14. Given $U \subseteq \mathbb{R}^p$, let $f : \mathbb{R}^n \times U \rightarrow \mathbb{R}$ be the objective function of the optimization problem

$$\underset{\mathbf{x} \in S(\mathbf{u})}{\text{maximize}} \quad f(\mathbf{x}, \mathbf{u}), \quad (14)$$

where $S : U \rightarrow \mathcal{P}(\mathbb{R}^n)$ is a constraint map. Assume that S is continuous at a point $\bar{\mathbf{u}} \in U$ and that the objective function f is continuous on $S(\bar{\mathbf{u}}) \times \{\bar{\mathbf{u}}\}$. If problem (14) has the unique optimal solution $\bar{\mathbf{x}}$ for $\mathbf{u} = \bar{\mathbf{u}}$, the following map $\Phi : U \rightarrow \mathcal{P}(\mathbb{R}^n)$ is continuous at $\bar{\mathbf{u}}$:

$$\Phi(\mathbf{u}) = \underset{\mathbf{x}}{\text{argmax}} \{f(\mathbf{x}, \mathbf{u}) \mid \mathbf{x} \in S(\mathbf{u})\}.$$

Proof. The same result is proved for a minimization problem in Theorem 3.30 in [5]. We can use this result directly. \square

Lemma 15 (Section 3.2.3 in [3]). Let $\mathcal{A} \subseteq \mathbb{R}^m$ and $f : \mathbb{R}^n \times \mathcal{A} \rightarrow [-\infty, +\infty]$. Let $g : \mathbb{R}^n \rightarrow [-\infty, +\infty]$ be the function defined by

$$g(\mathbf{x}) = \sup_{\mathbf{y} \in \mathcal{A}} f(\mathbf{x}, \mathbf{y}).$$

If $f(\cdot, \mathbf{y})$ is convex for each $\mathbf{y} \in \mathcal{A}$, g is convex.

A.1 Proof of Theorem 1

The proof proceeds along the lines of the proof of Theorem 1 in [1]. First, we derive the Lagrangian relaxation of problem (P) as follows:

$$\underset{\mathbf{w} \in \mathbb{R}^{m \times p}, \mathbf{b} \in \mathbb{R}^m, \boldsymbol{\eta} \in \mathbb{R}^{n \times m}}{\text{minimize}} \quad \underset{\boldsymbol{\alpha} \in \mathbb{R}^{n \times m}}{\text{maximize}} \quad \sum_{i=1}^n \ell(y_i, \boldsymbol{\eta}_{i\cdot}) + \frac{1}{2\gamma} \sum_{r=1}^m \|\mathbf{w}_{r\cdot}\|_2^2 + \sum_{i=1}^n \sum_{r=1}^m \alpha_{ir} (\mathbf{w}_{r\cdot}^\top \mathbf{x}_{i\cdot} + b_r - \eta_{ir}).$$

From Assumption 1, the loss function $\ell(y, \cdot)$ is proper convex for all $y \in [m]$ and problem (P) is bounded and has a interior feasible solution. Consequently, the strong duality holds by Lemma 12. We now transform this problem into

$$\underset{\boldsymbol{\alpha} \in \mathbb{R}^{n \times m}}{\text{maximize}} \quad \sum_{i=1}^n \min \{ \ell(y_i, \boldsymbol{\eta}_{i\cdot}) - \boldsymbol{\alpha}_{i\cdot}^\top \boldsymbol{\eta}_{i\cdot} \mid \boldsymbol{\eta}_{i\cdot} \in \mathbb{R}^m \} \quad (15)$$

$$+ \sum_{r=1}^m \min \{ b_r \mathbf{1}^\top \boldsymbol{\alpha}_{\cdot r} \mid b_r \in \mathbb{R} \} \quad (16)$$

$$+ \sum_{r=1}^m \min \left\{ \frac{1}{2\gamma} \|\mathbf{w}_{r\cdot}\|_2^2 + \mathbf{w}_{r\cdot}^\top \mathbf{X}^\top \boldsymbol{\alpha}_{\cdot r} \mid \mathbf{w}_{r\cdot} \in \mathbb{R}^p \right\}. \quad (17)$$

Because their decision variables are independent, these minimization problems (15)–(17) can be solved separately as follows:

- The optimal value of (15) is $-\hat{\ell}(y_i, \boldsymbol{\alpha}_{i\cdot})$ from the definition of the conjugate function.
- An equality constraint $\mathbf{1}^\top \boldsymbol{\alpha}_{\cdot r} = 0$ is obtained because problem (16) must be bounded.
- Problem (17) can be solved analytically; we have the optimal solution $\mathbf{w}_{r\cdot}^* = -\gamma \mathbf{X}^\top \boldsymbol{\alpha}_{\cdot r}$.

From these results, we have the dual problem as desired.

A.2 Proof of Lemma 2

First, for any $y \in [m]$ and $\boldsymbol{\alpha} \in \text{dom } \hat{\ell}(y, \cdot)$, the expression $\sum_{r=1}^m \sum_{j=1}^p z_j (\mathbf{x}_{\cdot j}^\top \boldsymbol{\alpha}_{\cdot r})^2$ is differentiable with respect to $\mathbf{z} \in [0, 1]^p$. In addition, $\text{dom } \hat{\ell}(y, \cdot)$ is bounded and closed from Assumption 3. The subderivative of c is therefore obtained by Lemma 13 as

$$\partial c(\mathbf{z}) = \text{conv} \left\{ \left(-\frac{\gamma}{2} \|\boldsymbol{\alpha}^\top \mathbf{x}_{\cdot j}\|_2^2 \right)_{j \in [p]} \mid \mathbf{1}^\top \boldsymbol{\alpha}_{\cdot r} = 0, \forall r \in [m], \boldsymbol{\alpha}_{i\cdot} \in \text{dom } \hat{\ell}(y_i, \cdot), \forall i \in [n], f_{\mathbf{z}}(\boldsymbol{\alpha}) = c(\mathbf{z}) \right\},$$

where $f_{\mathbf{z}}$ is the objective function (10). Moreover, $f_{\mathbf{z}}$ is strictly convex because of $\mathbf{X}\mathbf{X}^\top \succeq \mathbf{O}$ and the strict convexity of $\hat{\ell}(y, \cdot)$ from Assumption 2. That is, the map $\boldsymbol{\alpha}^*$ is a monomorphism. We thus obtain the partial derivatives of \mathbf{z} as (11).

Next we show the continuity of $\nabla c(\mathbf{z})$. The expression (11) is continuous at each $\mathbf{z} \in [0, 1]^p$ if the function $\boldsymbol{\alpha}^*$ is continuous; we therefore show the continuity of $\boldsymbol{\alpha}^*$ instead. The feasible region

$$\mathcal{A} = \left\{ \boldsymbol{\alpha} \in \mathbb{R}^{n \times m} \mid \mathbf{1}^\top \boldsymbol{\alpha}_{\cdot r} = 0, \forall r \in [m], \boldsymbol{\alpha}_i \in \text{dom } \hat{\ell}(y_i, \cdot), \forall i \in [n] \right\}$$

does not depend on \mathbf{z} , and thus the constraint map of $(D_{\mathbf{z}})$ is trivially continuous at each $\mathbf{z} \in [0, 1]^p$. For any $\mathbf{z} \in [0, 1]^p$, the objective function (10) is also continuous in $\boldsymbol{\alpha}$. From these facts and the uniqueness of $\boldsymbol{\alpha}^*(\mathbf{z})$, $\boldsymbol{\alpha}^*$ satisfies the assumptions of Lemma 14. Consequently, $\boldsymbol{\alpha}^*$ is continuous at each $\mathbf{z} \in [0, 1]^p$.

From the above discussion, c is continuously differentiable.

A.3 Proof of Lemma 3

For each $\boldsymbol{\alpha} \in \mathcal{A}$, the objective function (10) is linear in $\mathbf{z} \in [0, 1]^p$ and thus convex. Consequently, c is convex from Lemma 15.

A.4 Proof of Theorem 4

Algorithm 1 converges to an optimal solution if the following conditions are satisfied [2]:

- The optimization problem (7) is feasible and bounded.
- The objective function c is continuously differentiable and convex.

The former condition is clearly satisfied by Assumption 1, and the latter condition is also satisfied because Lemmas 2 and 3 hold under Assumptions 1–3. Because the loss function satisfies Assumptions 1–3, Algorithm 1 converges to the optimal solution in a finite number of iterations.

A.5 Proof of Proposition 5

First, we prove that the function $\ell^{\text{MNL}}(y, \cdot)$ is proper convex for any $y \in [m]$. From the definition, we have $\ell^{\text{MNL}}(y, \boldsymbol{\eta}) = -\log [\exp(\eta_y) / \sum_{s=1}^m \exp(\eta_s)]$. Because $0 < \exp(v) < +\infty$ for any $v \in \mathbb{R}$, the following inequality holds:

$$0 < \frac{\exp(\eta_y)}{\sum_{s=1}^m \exp(\eta_s)} < 1.$$

Consequently, $\text{dom } \ell^{\text{MNL}}(y, \cdot) = \mathbb{R}^m \neq \emptyset$ and $\ell^{\text{MNL}}(y, \boldsymbol{\eta}) > 0$. The convexity of $\ell^{\text{MNL}}(y, \cdot)$ is discussed in Section 3 in [7].

Next, we consider problem (P) of the MNL model. As discussed above, $0 < \ell^{\text{MNL}}(y, \boldsymbol{\eta}) < +\infty, \forall y \in [m], \boldsymbol{\eta} \in \mathbb{R}^m$. Consequently, (P) is bounded and feasible. Because $\text{dom } \ell^{\text{MNL}}(y, \cdot) = \mathbb{R}^m$ for all $y \in [m]$, (P) has an interior feasible solution.

A.6 Proof of Proposition 6

The continuity is trivial because $v \log v$ is continuous at each $v \in [0, 1]$. We also find that $\hat{\ell}^{\text{MNL}}(y, \cdot)$ is strictly convex for all $y \in [m]$ because its Hessian matrix is always positive definite; this is easily observed from the following equation:

$$\frac{\partial^2 \hat{\ell}^{\text{MNL}}}{\partial \alpha_s \partial \alpha_r}(y, \boldsymbol{\alpha}) = \begin{cases} \alpha_r^{-1} & \text{if } r \neq y \text{ and } s = r, \\ (1 + \alpha_r)^{-1} & \text{if } r = y \text{ and } s = r, \\ 0 & \text{otherwise,} \end{cases} \quad \forall \boldsymbol{\alpha} \in \mathcal{A}_y^{\text{MNL}}.$$

A.7 Proof of Proposition 7

We have $\alpha_{iy_i} = -\mathbf{1}^\top \boldsymbol{\alpha}_i^{\setminus y_i}$ from constraint (9). Consequently, the feasible region \mathcal{A}^{MNL} is bounded and closed.

A.8 Proof of Proposition 9

The second equality below is satisfied by the definition of \hat{g} :

$$\begin{aligned} g(\eta; \mathbf{p}) &= \sup\{\eta\alpha - \hat{g}(\alpha; \mathbf{p}) \mid \alpha \in [0, 1]\} \\ &= \sup\{\eta\alpha - (p_3\alpha^2 + p_2\alpha + p_1) \mid \alpha \in [0, 1]\}. \end{aligned} \quad (18)$$

We note that the objective function of problem (18) is concave from the assumption $p_3 > 0$. We differentiate the objective function with respect to α and then set it equal to zero as follows:

$$\eta - 2p_3\bar{\alpha} - p_2 = 0,$$

where $\bar{\alpha} \in \mathbb{R}$ is the stationary point. Consequently, the following holds:

$$\bar{\alpha} = \frac{\eta - p_2}{2p_3}. \quad (19)$$

Because the objective function is concave, the optimal value is given at $\alpha = 0$ and $\alpha = 1$ when $\bar{\alpha} < 0$ and $\bar{\alpha} > 1$, respectively.

These intervals can be transformed into the following intervals of η by equation (19):

$$\begin{aligned} \bar{\alpha} < 0 &\Leftrightarrow \eta < p_2, \\ \bar{\alpha} \in [0, 1] &\Leftrightarrow \eta \in [p_2, p_2 + 2p_3], \\ \bar{\alpha} > 1 &\Leftrightarrow \eta > p_2 + 2p_3. \end{aligned}$$

Consequently, we have the desired result.

A.9 Proof of Proposition 10

The second equality below is satisfied by the definition of ℓ^{Titsias} :

$$\begin{aligned} \tilde{\ell}^{\text{Titsias}}(y, \boldsymbol{\alpha}) &= \sup\{\boldsymbol{\alpha}^\top \boldsymbol{\eta} - \ell^{\text{Titsias}}(y, \boldsymbol{\eta}) \mid \boldsymbol{\eta} \in \mathbb{R}^m\} \\ &= \sup\{\boldsymbol{\alpha}^\top \boldsymbol{\eta} - \sum_{s \neq y} \log[1 + \exp(\eta_s - \eta_y)] \mid \boldsymbol{\eta} \in \mathbb{R}^m\}. \end{aligned} \quad (20)$$

Because optimization problem (20) has no constraints and has a convex objective function, we obtain an optimal solution by the gradient with respect to $\boldsymbol{\eta}$. Let us consider the following two cases.

Case 1: $r \neq y$ holds. First, we calculate the partial derivative with respect to η_r , and set it equal to zero as follows:

$$\alpha_r - \frac{\exp(\eta_r^* - \eta_y^*)}{1 + \exp(\eta_r^* - \eta_y^*)} = 0,$$

where $\boldsymbol{\eta}^* \in \mathbb{R}^m$ is an optimal solution of problem (20). That is,

$$\frac{1}{1 + \exp(\eta_y^* - \eta_r^*)} = \alpha_r. \quad (21)$$

From this equation, the following two equations are obtained:

$$\eta_y^* - \eta_r^* = \log(1 - \alpha_r) - \log \alpha_r, \quad (22)$$

$$1 + \exp(\eta_r^* - \eta_y^*) = (1 - \alpha_r)^{-1}. \quad (23)$$

Note that we assume $0 < \alpha_r < 1$ to derive these equations.

Case 2: $r = y$ holds. Similarly, we calculate the partial derivative with respect to η_r , and set it equal to zero as follows:

$$\alpha_y + \sum_{s \neq y} \left(\frac{1}{1 + \exp(\eta_y^* - \eta_s^*)} \right) = 0.$$

Consequently, the following equation is obtained from equation (21):

$$\mathbf{1}^\top \boldsymbol{\alpha} = 0. \quad (24)$$

We obtain the following derivation from equations (22) and (23):

$$\begin{aligned} \hat{\ell}^{\text{Titsias}}(y, \boldsymbol{\alpha}) &= \boldsymbol{\alpha}^\top \boldsymbol{\eta}^* - \sum_{s \neq y} \log[1 + \exp(\eta_s^* - \eta_y^*)] \\ &= \sum_{s \neq y} \alpha_s (\eta_y^* - \log(1 - \alpha_s) + \log \alpha_s) + \alpha_y \eta_y^* + \sum_{s \neq y} \log[1 - \alpha_s] \\ &= \sum_{s \neq y} (\alpha_s \log \alpha_s + (1 - \alpha_s) \log(1 - \alpha_s)) + \eta_y^* \sum_{s \in [m]} \alpha_s. \end{aligned}$$

Consequently, from equation (24), the following equation holds:

$$\hat{\ell}^{\text{Titsias}}(y, \boldsymbol{\alpha}) = \sum_{s \neq y} [\alpha_s \log \alpha_s + (1 - \alpha_s) \log(1 - \alpha_s)].$$

Because $v \log v \rightarrow 0$ when $v \rightarrow 0$, we have the desired result.

A.10 Proof of Theorem 11

The loss function $\ell^{\text{quad}} : [m] \times \mathbb{R}^m \rightarrow \mathbb{R}$ is defined by

$$\ell^{\text{quad}}(y, \boldsymbol{\eta}) = \sum_{s \neq y} g(\eta_s - \eta_y; \mathbf{p}),$$

where

$$g(\eta; \mathbf{p}) = \begin{cases} -p_1 & \text{if } \eta < p_2, \\ (\eta - p_2)^2/4p_3 - p_1 & \text{if } \eta \in [p_2, p_2 + 2p_3], \\ \eta - (p_1 + p_2 + p_3) & \text{otherwise.} \end{cases}$$

Moreover, for each $y \in [m]$, let $\hat{\ell}^{\text{quad}}(y, \boldsymbol{\alpha})$ be the conjugate defined as follows:

$$\hat{\ell}^{\text{quad}}(y, \boldsymbol{\alpha}) = \begin{cases} \sum_{s \neq y} p_3 \alpha_s^2 + p_2 \alpha_s + p_1 & \text{if } \boldsymbol{\alpha} \in \mathcal{A}_y^{\text{quad}}, \\ +\infty & \text{otherwise,} \end{cases}$$

where

$$\mathcal{A}_y^{\text{quad}} = \{\boldsymbol{\alpha} \in \mathbb{R}^m \mid \mathbf{1}^\top \boldsymbol{\alpha} = 0, \mathbf{0} \leq \boldsymbol{\alpha}^{\setminus y} \leq \mathbf{1}\}.$$

We show that these two functions and problem (P) satisfy Assumptions 1–3.

First, for each $y \in [m]$, $\ell^{\text{quad}}(y, \boldsymbol{\alpha})$ is the summation of quadratic functions; $\ell^{\text{quad}}(y, \boldsymbol{\alpha})$ is thus strictly convex in $\boldsymbol{\alpha}$ if and only if $p_3 > 0$. Consequently Assumption 2 is satisfied when $p_3 > 0$. Second, we have $\mathbf{0} \leq \boldsymbol{\alpha}^{\setminus y} \leq \mathbf{1}$ and $\alpha_y = -\mathbf{1}^\top \boldsymbol{\alpha}^{\setminus y}$ from the definition of $\mathcal{A}_y^{\text{quad}}$. Consequently, $\mathcal{A}_y^{\text{quad}}$ is bounded and closed; that is, Assumption 3 is satisfied. Finally, the function ℓ^{quad} is bounded below because $g(\eta; \mathbf{p}) \geq -p_1$ when $p_3 > 0$. Because the conjugate of a closed proper convex function is still closed proper convex, $\hat{\ell}^{\text{quad}}(y, \boldsymbol{\alpha})$ is a proper convex function. As with the proof of Theorem 5, (P) is bounded and has an interior feasible solution. Consequently, Assumption 1 is satisfied.

From the above discussion, all the assumptions of Theorem 4 are satisfied, and we have the desired result.

References

- [1] D. Bertsimas, J. Pauphilet, and B. V. Parys. Sparse classification and phase transitions: A discrete optimization perspective. arXiv preprint arXiv:1710.01352, 2017.
- [2] P. Bonami, L. Biegler, A. Conn, G. Cornuéjols, I. Grossmann, C. Laird, J. Lee, A. Lodi, F. Margot, N. Sawaya, and A. Wächter. An algorithmic framework for convex mixed integer nonlinear programs. *Discrete Optimization* 5(2):186–204, 2008.

- [3] S. Boyd and L. Vandenberghe. Convex optimization. Cambridge University Press, 2004.
- [4] S. Boyd and L. Vandenberghe. Subgradients. Lecture notes for EE364b, Stanford University, Winter Quarter 2006–2007, 2008.
- [5] M. Fukushima. Basic nonlinear optimization (in Japanese). Asakura Publishing, 2001.
- [6] T. Kanamori, T. Suzuki, I. Takeuchi, and I. Sato. Continuous optimization for machine learning (in Japanese). Kodansha Scientific, 2016.
- [7] J. D. M. Rennie. Regularized logistic regression is strictly convex. Technical report, MIT, [<http://qwone.com/~jason/writing/convexLR.pdf>], 2005. Last accessed: 2018-09-05.