Feature-Weighted Elastic Net for Competing Risk Outcomes

Lukas Burk

The elastic net objective function is given as:

$$J(\beta_0,\beta) = \frac{1}{2}||\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta||_2^2 + \lambda \sum_{j=1}^p \left(\alpha|\beta_j| + \frac{1-\alpha}{2}\beta_j^2\right)$$

 $\lambda \in \mathbb{R}$ controls overall sparsity

The elastic net objective function is given as:

$$J(\beta_0,\beta) = \frac{1}{2}||\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta||_2^2 + \lambda \sum_{j=1}^p \left(\alpha|\beta_j| + \frac{1-\alpha}{2}\beta_j^2\right)$$

- $\lambda \in \mathbb{R}$ controls overall sparsity

The elastic net objective function is given as:

$$J(\beta_0,\beta) = \frac{1}{2}||\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta||_2^2 + \lambda \sum_{j=1}^p \left(\alpha|\beta_j| + \frac{1-\alpha}{2}\beta_j^2\right)$$

- $\lambda \in \mathbb{R}$ controls overall sparsity
- ► Higher \Rightarrow larger penalty on all β_j , equally

 $\alpha \in [0,1]$ is the mixing parameter

The elastic net objective function is given as:

$$J(\beta_0,\beta) = \frac{1}{2}||\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta||_2^2 + \lambda \sum_{j=1}^p \left(\alpha|\beta_j| + \frac{1-\alpha}{2}\beta_j^2\right)$$

- $\lambda \in \mathbb{R}$ controls overall sparsity
- ► Higher \Rightarrow larger penalty on all β_i , equally

- $\alpha \in [0,1]$ is the mixing parameter
- $1 \Rightarrow \text{only } \ell^1 \text{ penalty}$ (LASSO)

The elastic net objective function is given as:

$$J(\beta_0,\beta) = \frac{1}{2}||\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta||_2^2 + \lambda \sum_{j=1}^p \left(\alpha|\beta_j| + \frac{1-\alpha}{2}\beta_j^2\right)$$

- $\lambda \in \mathbb{R}$ controls overall sparsity
- ▶ Higher \Rightarrow larger penalty on all β_j , equally

 $\begin{array}{c}
1 \Rightarrow \text{only } \ell^1 \text{ penalty} \\
\text{(LASSO)}
\end{array}$

parameter

 $\alpha \in [0,1]$ is the mixing

 $ightharpoonup 0 \Rightarrow \text{only } \ell_2^2 \text{ penalty}$

Introduced by Tay et al. (2020)

- Introduced by Tay et al. (2020)
- ▶ Multiple applications: External information & feature-grouping

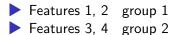
- Introduced by Tay et al. (2020)
- ▶ Multiple applications: External information & feature-grouping
- Sometimes it's desirable to adjust penalization weights on individual or groups of coefficients

- Introduced by Tay et al. (2020)
- ▶ Multiple applications: External information & feature-grouping
- Sometimes it's desirable to adjust penalization weights on individual or groups of coefficients
- Assign groups via matrix $\mathbf{Z} \in \mathbb{R}^{p \times K}$

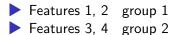
$$\mathbf{Z} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$

Features 1, 2 group 1

$$\mathbf{Z} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$



$$\mathbf{Z} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$



$$\mathbf{Z} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}$$

Features 3, 4 group 2

Now imagine this, but with $p>>1000\ \mathrm{and}$ e.g. genetic data.

Feature-Weighting

$$J(\beta_0,\beta) = \frac{1}{2}||\mathbf{y} - \beta_0\mathbf{1} - \mathbf{X}\beta||_2^2 + \lambda \sum_{j=1}^p w_j(\theta) \left(\alpha|\beta_j| + \frac{1-\alpha}{2}\beta_j^2\right)$$

where

Feature-Weighting

$$J(\beta_0,\beta) = \frac{1}{2}||\mathbf{y} - \beta_0\mathbf{1} - \mathbf{X}\beta||_2^2 + \lambda \sum_{j=1}^p w_j(\theta) \left(\alpha|\beta_j| + \frac{1-\alpha}{2}\beta_j^2\right)$$

where

$$w_j(\theta) = \frac{\sum_{l=1}^{p} \exp(\mathbf{z}_l^T \theta)}{p \exp(\mathbf{z}_j^T \theta)}$$

Defines the **penalization weight** of coefficient j based on its corresponding value in ${\bf Z}$ and additionally hyper-parameter θ

 $igwedge w_j(heta)$ is chosen heuristically, suggested by the authors for desirable properties

- $igwedge w_j(heta)$ is chosen heuristically, suggested by the authors for desirable properties
- $ightharpoonup \mathbf{z}_{j}^{T} \theta$ functions as a *score*

- $igwedge w_j(heta)$ is chosen heuristically, suggested by the authors for desirable properties
- $\triangleright \mathbf{z}_{j}^{T} \theta$ functions as a *score*
 - ightharpoonup = 0, reduces to original elastic net

- $lackbox{ } w_j(heta)$ is chosen heuristically, suggested by the authors for desirable properties
- $\triangleright \mathbf{z}_{i}^{T} \theta$ functions as a score
 - ightharpoonup = 0, reduces to original elastic net
 - ► Higher score → lower penalization weight, feature is "more important"

- $lackbox{ } w_j(heta)$ is chosen heuristically, suggested by the authors for desirable properties
- $\triangleright \mathbf{z}_{i}^{T} \theta$ functions as a score
 - ightharpoonup = 0, reduces to original elastic net
 - ► Higher score → lower penalization weight, feature is "more important"

- $igwedge w_j(heta)$ is chosen heuristically, suggested by the authors for desirable properties
- $\triangleright \mathbf{z}_{i}^{T} \theta$ functions as a score
 - ightharpoonup = 0, reduces to original elastic net
 - ► Higher score → lower penalization weight, feature is "more important"
- In the "feature grouping" setting, this just allows group-specific penalization weights.

- $lackbox{ } w_j(heta)$ is chosen heuristically, suggested by the authors for desirable properties
- $\triangleright \mathbf{z}_{i}^{T} \theta$ functions as a score
 - ightharpoonup = 0, reduces to original elastic net
 - ► Higher score → lower penalization weight, feature is "more important"
- In the "feature grouping" setting, this just allows group-specific penalization weights.
- Related to the "group lasso" (Jacob, Obozinski, and Vert 2009)

Z can also be $p \times 1$: No groups, just weights

- $ightharpoonup {f Z}$ can also be p imes 1: No groups, just weights
- ightharpoonup Example from paper: ${f Z}$ set to noisy version of true $|\beta|$ in simulation study

- $ightharpoonup {f Z}$ can also be p imes 1: No groups, just weights
- ightharpoonup Example from paper: ${f Z}$ set to noisy version of true $|\beta|$ in simulation study

- **Z** can also be $p \times 1$: No groups, just weights
- \blacktriangleright Example from paper: ${\bf Z}$ set to noisy version of true $|\beta|$ in simulation study
- lackbrack Higher $|eta_j| \Rightarrow$ lower penalization for \hat{eta}_j

- **Z** can also be $p \times 1$: No groups, just weights
- \blacktriangleright Example from paper: ${\bf Z}$ set to noisy version of true $|\beta|$ in simulation study
- $lackbr{\blacktriangleright}$ Higher $|eta_j| \Rightarrow$ lower penalization for \hat{eta}_j
- $lackbrack |eta_j|pprox 0 \Rightarrow$ higher penalization for \hat{eta}_j

Authors suggest multi-task learning algorithm, outline:

- Authors suggest multi-task learning algorithm, outline:
- 1. $\beta_1^{(0)},\beta_2^{(0)} \colon \texttt{glmnet}$ solution for $(\mathbf{X},\mathbf{y}_1),(\mathbf{X},\mathbf{y}_2)$ respectively

- Authors suggest multi-task learning algorithm, outline:
- 1. $\beta_1^{(0)}, \beta_2^{(0)}$: glmnet solution for $(\mathbf{X}, \mathbf{y}_1), (\mathbf{X}, \mathbf{y}_2)$ respectively
- 2. For k = 0, 1, ...:

- Authors suggest multi-task learning algorithm, outline:
- 1. $\beta_1^{(0)}, \beta_2^{(0)}$: glmnet solution for $(\mathbf{X}, \mathbf{y}_1), (\mathbf{X}, \mathbf{y}_2)$ respectively
- 2. For k = 0, 1, ...:
 - a) $\mathbf{Z}_2 = \left| eta_1^{(k)} \right|$. Fit fwelnet with $(\mathbf{X}, \mathbf{y}_2, \mathbf{Z}_2)$

- Authors suggest multi-task learning algorithm, outline:
- 1. $\beta_1^{(0)}, \beta_2^{(0)}$: glmnet solution for $(\mathbf{X}, \mathbf{y}_1), (\mathbf{X}, \mathbf{y}_2)$ respectively
- 2. For k = 0, 1, ...:
 - a) $\mathbf{Z}_2 = \left| eta_1^{(k)} \right|$. Fit fwelnet with $(\mathbf{X}, \mathbf{y}_2, \mathbf{Z}_2)$
 - $lackbox{ Set } \left|eta_2^{(k+1)}
 ight|$ to solution with optimal lambda

- Authors suggest multi-task learning algorithm, outline:
- 1. $\beta_1^{(0)}, \beta_2^{(0)}$: glmnet solution for $(\mathbf{X}, \mathbf{y}_1), (\mathbf{X}, \mathbf{y}_2)$ respectively
- 2. For k = 0, 1, ...:
 - a) $\mathbf{Z}_2 = \left| eta_1^{(k)} \right|$. Fit fwelnet with $(\mathbf{X}, \mathbf{y}_2, \mathbf{Z}_2)$
 - $lackbox{Set} \left| eta_2^{(k+1)}
 ight|$ to solution with optimal lambda
 - b) $\mathbf{Z}_1 = \left| \beta_2^{(k+1)} \right|$. Fit fwelnet with $(\mathbf{X}, \mathbf{y}_1, \mathbf{Z}_1)$

- Authors suggest multi-task learning algorithm, outline:
- 1. $\beta_1^{(0)}, \beta_2^{(0)}$: glmnet solution for $(\mathbf{X}, \mathbf{y}_1), (\mathbf{X}, \mathbf{y}_2)$ respectively
- 2. For k = 0, 1, ...:
 - a) $\mathbf{Z}_2 = \left| eta_1^{(k)} \right|$. Fit fwelnet with $(\mathbf{X}, \mathbf{y}_2, \mathbf{Z}_2)$
 - lacksquare Set $\left|eta_2^{(k+1)}
 ight|$ to solution with optimal lambda
 - b) $\mathbf{Z}_1 = \left| eta_2^{(k+1)} \right|$. Fit fwelnet with $(\mathbf{X}, \mathbf{y}_1, \mathbf{Z}_1)$
 - ightharpoonup Set $\left|eta_1^{(k+1)}\right|$ to solution with optimal lambda

Transfer to Competing Risks

 \blacktriangleright Setting with two causes: $\mathbf{y}_1,\mathbf{y}_2\in\{0,1\}$

- \blacktriangleright Setting with two causes: $\mathbf{y}_1,\mathbf{y}_2\in\{0,1\}$
- Assumption: Shared information for both causes:

- ▶ Setting with two causes: $y_1, y_2 \in \{0, 1\}$
- Assumption: Shared information for both causes:
 - If \mathbf{x}_j is important for \mathbf{y}_1 , also for \mathbf{y}_2 ?

- ▶ Setting with two causes: $\mathbf{y}_1, \mathbf{y}_2 \in \{0, 1\}$
- Assumption: Shared information for both causes:
 - If \mathbf{x}_j is important for \mathbf{y}_1 , also for \mathbf{y}_2 ?
 - ightharpoonup \Rightarrow Avoid shrinking it to 0 in cause-specific models

- ▶ Setting with two causes: $y_1, y_2 \in \{0, 1\}$
- Assumption: Shared information for both causes:
 - If \mathbf{x}_i is important for \mathbf{y}_1 , also for \mathbf{y}_2 ?
 - ightharpoonup \Rightarrow Avoid shrinking it to 0 in cause-specific models
- ▶ Basic idea: Adapt previous algorithm to Cox regression

- ▶ Setting with two causes: $y_1, y_2 \in \{0, 1\}$
- Assumption: Shared information for both causes:
 - If \mathbf{x}_j is important for \mathbf{y}_1 , also for \mathbf{y}_2 ?
 - ightharpoonup \Rightarrow Avoid shrinking it to 0 in cause-specific models
- ▶ Basic idea: Adapt previous algorithm to Cox regression
- Multi-task ⇒ "Multi-cause"

▶ 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet

- ▶ 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet
- ▶ 2: Implement algorithm for Cox/CR setting

- 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet
- ▶ 2: Implement algorithm for Cox/CR setting
- ▶ 3: Apply algorithm to some "easy" simulated data settings:

- 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet
- ▶ 2: Implement algorithm for Cox/CR setting
- ▶ 3: Apply algorithm to some "easy" simulated data settings:
 - Large effects ($\beta_j = 1$ or 0.25 for "small effect")

- 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet
- ▶ 2: Implement algorithm for Cox/CR setting
- ▶ 3: Apply algorithm to some "easy" simulated data settings:
 - ▶ Large effects ($\beta_i = 1$ or 0.25 for "small effect")
 - N = 1000

- 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet
- ▶ 2: Implement algorithm for Cox/CR setting
- ▶ 3: Apply algorithm to some "easy" simulated data settings:
 - Large effects ($\beta_i = 1$ or 0.25 for "small effect")
 - N = 1000
 - ▶ 11-14 noise variables

- 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet
- ▶ 2: Implement algorithm for Cox/CR setting
- ▶ 3: Apply algorithm to some "easy" simulated data settings:
 - Large effects ($\beta_i = 1$ or 0.25 for "small effect")
 - N = 1000
 - ▶ 11-14 noise variables

- ▶ 1: Adapt fwelnet for Coxnet/Surv endpoint via glmnet
- ▶ 2: Implement algorithm for Cox/CR setting
- ▶ 3: Apply algorithm to some "easy" simulated data settings:
 - Large effects ($\beta_i = 1$ or 0.25 for "small effect")
 - N = 1000
 - ▶ 11-14 noise variables

First goal: See if we find some improvement over cause-specific glmnet

What Comes Next

▶ Wider set of simulation cases

What Comes Next

- Wider set of simulation cases
- Evaluate predictive performance (non-trivial in CR/censored setting)

References

- Jacob, Laurent, Guillaume Obozinski, and Jean-Philippe Vert. 2009. "Group Lasso with Overlap and Graph Lasso." In *Proceedings of the 26th Annual International Conference on Machine Learning ICML '09*, 1–8. ACM Press. https://doi.org/10.1145/1553374.1553431.
- Tay, J. Kenneth, Nima Aghaeepour, Trevor Hastie, and Robert Tibshirani. 2020. "Feature-Weighted Elastic Net: Using 'Features of Features' for Better Prediction." arXiv:2006.01395 [Cs, Stat], June. http://arxiv.org/abs/2006.01395.
- Zou, Hui, and Trevor Hastie. 2005. "Regularization and Variable Selection via the Elastic Net." *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 67 (2): 301–20. https://www.jstor.org/stable/3647580.