



# ACHIEVING FAIRNESS WITH A SIMPLE RIDGE PENALTY ESTIMATES AND UNCERTAINTY

Marco Scutari  
[scutari@bnlearn.com](mailto:scutari@bnlearn.com)

Dalle Molle Institute for  
Artificial Intelligence (IDSIA)

December 12, 2025

## → INTRODUCTION

FAIR LINEAR MODELS

FAIR RIDGE REGRESSION

BAYESIAN FAIR REGRESSION

UNCERTAINTY QUANTIFICATION

CONCLUSIONS & ACKNOWLEDGEMENTS

- Machine learning (statistical?) models are being used in applications where it is crucial to ensure the **accountability and fairness** of the decisions made based on their outputs.
- Models are trained on historical data that contain various forms of bias, **capture those biases and carry them over** into current applications resulting in unfair discrimination of certain groups of people.
- The **concept of fairness** itself is difficult to define because it depends on the type of distortion we wish to limit and how we characterise it mathematically.
- How can we specify **fair models** that capture the non-discriminating information in the data and disregard the discriminating information?

Say that  $y$  is our response,  $\hat{y}$  are fitted values from the model,  $S$  are the sensitive attributes containing the discriminating information and  $X$  are the other predictors.

- **Group fairness:** predictions should be similar across the groups identified by the sensitive attributes.
  - Statistical or demographic parity ( $\hat{y} \perp\!\!\!\perp S$ ).
  - Equality of opportunity ( $\hat{y} \perp\!\!\!\perp S \mid y$ ).
- **Individual fairness:** individuals that are similar receive similar predictions

$$f(y, S) = \sum_{i,j} d_1(y_i, y_j) d_2(s_i, s_j).$$

Many, many mathematical characterisations in the literature [10, 4, 11].

We can enforce fairness at different stages of the model selection, estimation and validation process [2]:

- **Pre-processing:** transform the data to remove the underlying discrimination so that models are guaranteed to be fair.
- **In-processing:** modify model estimation to remove discrimination, either by changing its objective function (typically the log-likelihood) or by imposing constraints on its parameters.
- **Post-processing:** assess a previously-estimated model, treated as a black box, and alter its predictions to make them fair.

✓ INTRODUCTION

→ FAIR LINEAR MODELS

FAIR RIDGE REGRESSION

BAYESIAN FAIR REGRESSION

UNCERTAINTY QUANTIFICATION

CONCLUSIONS & ACKNOWLEDGEMENTS

Komiyama *et al.* [7] did:

1. remove the association between  $\mathbf{X}$  and  $\mathbf{S}$  with  $\mathbf{X} = \mathbf{B}^T \mathbf{S} + \mathbf{U}$ , estimating  $\widehat{\mathbf{B}}_{\text{OLS}} = (\mathbf{S}^T \mathbf{S})^{-1} \mathbf{S}^T \mathbf{X}$ ;
2. take the **decorrelated predictors**  $\widehat{\mathbf{U}} = \mathbf{X} - \widehat{\mathbf{B}}_{\text{OLS}}^T \mathbf{S}$  which contain the component of  $\mathbf{X}$  that cannot be explained by  $\mathbf{S}$  ( $\widehat{\mathbf{U}} \perp \mathbf{S}$ );
3. formulate the regression model  $\mathbf{y} = \mathbf{S}\boldsymbol{\alpha} + \widehat{\mathbf{U}}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ;
4. formulate the fairness constraint

$$R_{\mathbf{S}}^2(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\text{VAR}(\mathbf{S}\boldsymbol{\alpha})}{\text{VAR}(\widehat{\mathbf{y}})} = \frac{\boldsymbol{\alpha}^T \text{VAR}(\mathbf{S})\boldsymbol{\alpha}}{\boldsymbol{\alpha}^T \text{VAR}(\mathbf{S})\boldsymbol{\alpha} + \boldsymbol{\beta}^T \text{VAR}(\widehat{\mathbf{U}})\boldsymbol{\beta}};$$

5. solve the optimisation problem

$$\min_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \mathbb{E}[(\mathbf{y} - \widehat{\mathbf{y}})^2] \quad \text{such that} \quad R_{\mathbf{S}}^2(\boldsymbol{\alpha}, \boldsymbol{\beta}) \leq r, r \in [0, 1].$$

### PROS:

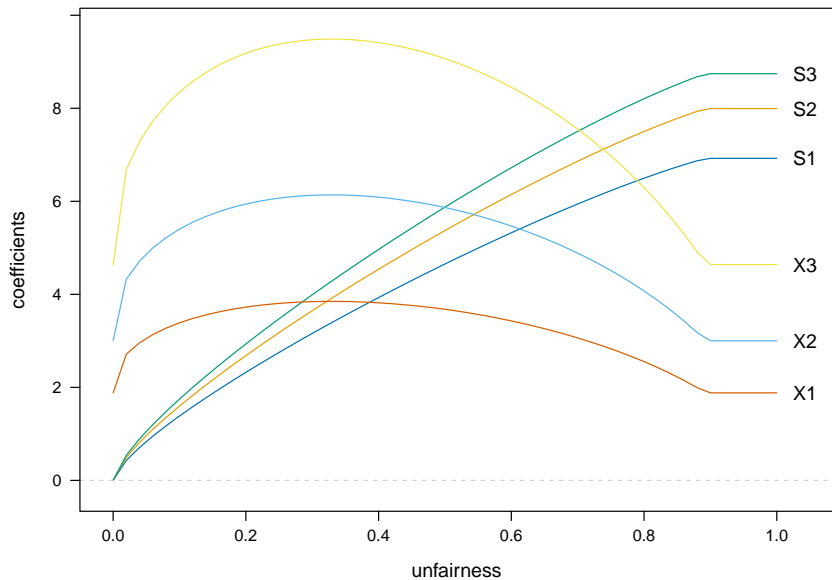
- The formulation is simple.
- Discriminating and non-discriminating information are separated.
- The optimisation problem is QCQP, for which there are solvers.
- The fairness constraint is defined in terms of explained variance, the natural measure of information in a linear model.
- The bound is interpretable: 0 = complete fairness, 1 = no constraint.

### CONS:

- No distributional assumptions.
- Cannot be extended without losing the ability to use QCQP solvers.
- The behaviour of the estimated coefficients is weird.



## COEFFICIENT PROFILE PLOTS IN KOMIYAMA ET AL.



✓ INTRODUCTION

✓ FAIR LINEAR MODELS

→ FAIR RIDGE REGRESSION

BAYESIAN FAIR REGRESSION

UNCERTAINTY QUANTIFICATION

CONCLUSIONS & ACKNOWLEDGEMENTS

Take two vintage pieces of statistics from the 1970s-1980s:

1. ridge regression (RR) [6];
2. generalised linear models (GLMs) [9].

With them, we build a **Fair (Generalised) Ridge Regression Model** (F(G)RRM) that fixes the CONS above and keep all the PROS:

- **Modular:** swappable characterisation of fairness.
- **Versatile:** supports all generalised linear models.
- **Interpretable:** both the model and the fairness constraints are interpretable, and all the best practices from the literature apply.
- **Statistical:** model selection, model validation, hypothesis testing, confidence intervals, etc. are already available in the literature.

Let's start again from  $\mathbf{y} = \mathbf{S}\boldsymbol{\alpha} + \widehat{\mathbf{U}}\boldsymbol{\beta} + \varepsilon$ . We want to re-create the shrinkage effects on the coefficients  $\boldsymbol{\alpha}$  associated with  $\mathbf{S}$  that we see in Komiyama *et al.*: we can do that with a ridge penalty,

$$(\widehat{\boldsymbol{\alpha}}_{\text{FRRM}}, \widehat{\boldsymbol{\beta}}_{\text{FRRM}}) = \underset{\boldsymbol{\alpha}, \boldsymbol{\beta}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{S}\boldsymbol{\alpha} - \widehat{\mathbf{U}}\boldsymbol{\beta}\|_2^2 + \lambda(r)\|\boldsymbol{\alpha}\|_2^2,$$

which we only apply to  $\boldsymbol{\alpha}$  because by construction there is no discriminating information in  $\widehat{\mathbf{U}}$ . The parameter estimates are in closed form:

$$\begin{bmatrix} \widehat{\boldsymbol{\alpha}}_{\text{FRRM}} \\ \widehat{\boldsymbol{\beta}}_{\text{FRRM}} \end{bmatrix} = \begin{bmatrix} (\mathbf{S}^T\mathbf{S} + \lambda(r)\mathbf{I})^{-1} \mathbf{S}^T\mathbf{y} \\ (\widehat{\mathbf{U}}^T\widehat{\mathbf{U}})^{-1} \widehat{\mathbf{U}}^T\mathbf{y} \end{bmatrix}.$$

But how do we control the fairness of the model?

For a given level of fairness  $r \in [0, 1]$ :

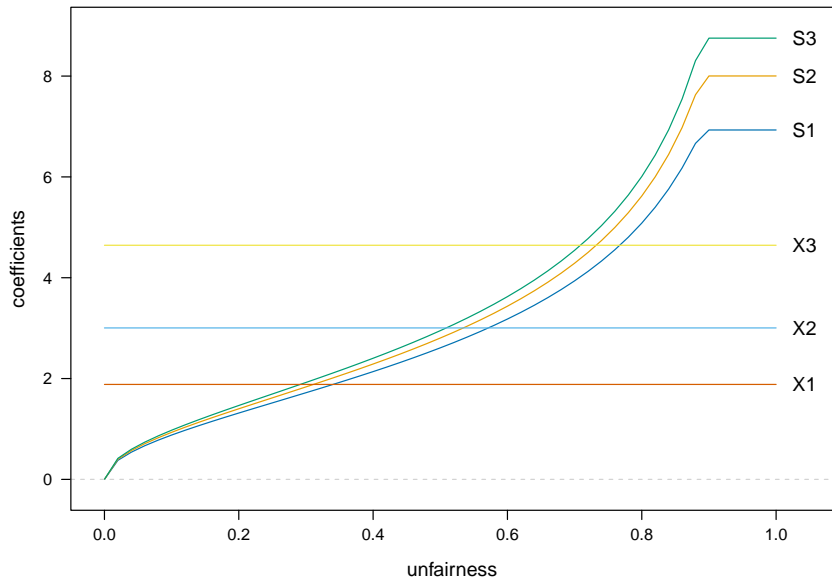
1. Compute  $\widehat{\mathbf{U}}$  from  $\mathbf{X}, \mathbf{S}$ .
2. Estimate  $\widehat{\boldsymbol{\beta}}_{\text{FRRM}} = (\widehat{\mathbf{U}}^T \widehat{\mathbf{U}})^{-1} \widehat{\mathbf{U}}^T \mathbf{y}$ .
3. Estimate  $\widehat{\boldsymbol{\alpha}}_{\text{OLS}} = (\mathbf{S}^T \mathbf{S})^{-1} \mathbf{S}^T \mathbf{y}$ . Then:
  - 3.1 If  $R_{\mathbf{S}}^2(\widehat{\boldsymbol{\alpha}}_{\text{OLS}}, \widehat{\boldsymbol{\beta}}_{\text{OLS}}) \leq r$ , set  $\widehat{\boldsymbol{\alpha}}_{\text{FRRM}} = \widehat{\boldsymbol{\alpha}}_{\text{OLS}}$ .
  - 3.2 Otherwise, find the value of  $\lambda(r)$  that satisfies

$$\boldsymbol{\alpha}^T \text{VAR}(\mathbf{S}) \boldsymbol{\alpha} = \frac{r}{1-r} \widehat{\boldsymbol{\beta}}_{\text{FRRM}}^T \text{VAR}(\widehat{\mathbf{U}}) \widehat{\boldsymbol{\beta}}_{\text{FRRM}}$$

and estimate the associated  $\widehat{\boldsymbol{\alpha}}_{\text{FRRM}}$  in the process.

**A single solution**, requiring a simple **univariate** root finding algorithm regardless of the number of variables involved.

## COEFFICIENTS IN FRRM



We can easily **replace  $R_S^2(\alpha, \beta)$  with other constraints.**

1.  $R_{EO}^2(\phi, \psi) = \frac{\text{VAR}(\mathbf{S}\phi)}{\text{VAR}(\mathbf{y}\psi + \mathbf{S}\phi)}$ , from  $\hat{\mathbf{y}} = \mathbf{y}\psi + \mathbf{S}\phi + \epsilon^*$ .
  2.  $f(\alpha, \mathbf{y}, \mathbf{S}) = \sum_{i,j} d(y_i, y_j)(\mathbf{s}_i\alpha - \mathbf{s}_j\alpha)^2$  and  
 $D_{IF} = f(\hat{\alpha}_{\text{FRRM}}, \mathbf{y}, \mathbf{S}) / f(\hat{\alpha}_{\text{OLS}}, \mathbf{y}, \mathbf{S})$ ,
  3. Any convex combination of  $R_S^2(\cdot)$ ,  $R_{EO}^2(\cdot)$ ,  $D_{IF}(\cdot)$  and others.
- 

We can draw on [5, 12, 13] to estimate the **FGRRM**

$$(\hat{\alpha}_{\text{FRRM}}, \hat{\beta}_{\text{FRRM}}) = \underset{\alpha, \beta}{\operatorname{argmin}} D(\alpha, \beta) + \lambda(r) \|\alpha\|_2^2.$$

where  $D(\cdot)$  is the **deviance** of a GLM + RR, choosing  $\lambda(r)$  to give

$$\frac{D(\alpha, \beta) - D(\mathbf{0}, \beta)}{D(\alpha, \beta) - D(\mathbf{0}, \mathbf{0})} \leq r.$$

✓ INTRODUCTION

✓ FAIR LINEAR MODELS

✓ FAIR RIDGE REGRESSION

→ **BAYESIAN FAIR REGRESSION**

UNCERTAINTY QUANTIFICATION

CONCLUSIONS & ACKNOWLEDGEMENTS



## A BAYESIAN CONSTRUCTION SIMILAR TO FGRRM?

Take a logistic FGRRM with linear component  $\boldsymbol{\eta} = \mathbf{S}\boldsymbol{\alpha} + \widehat{\mathbf{U}}\boldsymbol{\beta}$ :

$$\mathbf{y} \mid \widehat{\mathbf{U}}, \mathbf{S}, \boldsymbol{\alpha}, \boldsymbol{\beta} \sim \text{Ber}(\boldsymbol{\pi}), \quad \boldsymbol{\pi} = \text{logit}^{-1}(\boldsymbol{\eta}) = \frac{\exp(\boldsymbol{\eta})}{1 + \exp(\boldsymbol{\eta})}.$$

We can implement it with **Bayesian logistic regression** where

$$\boldsymbol{\alpha} \sim N\left(0, (\sigma_{\boldsymbol{\alpha}}^2/\lambda)\mathbf{I}_q\right), \quad \boldsymbol{\beta} \sim N\left(0, \sigma_{\boldsymbol{\beta}}^2\mathbf{I}_p\right),$$

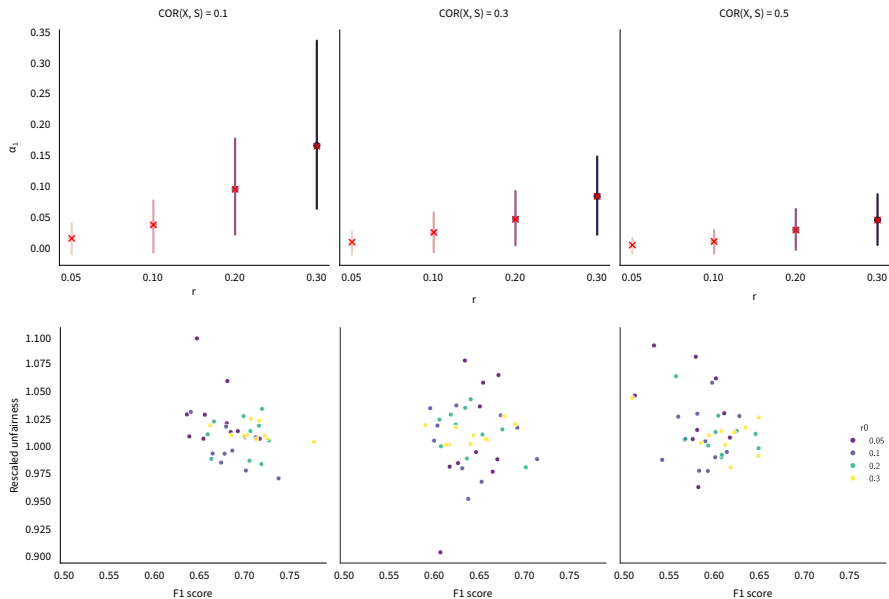
and tie everything together with a prior on  $\lambda(r)$

$$P\left(\boldsymbol{\alpha}, \boldsymbol{\beta}, \lambda(r) \mid \mathbf{y}, \widehat{\mathbf{U}}, \mathbf{S}\right) \propto P\left(\mathbf{y} \mid \widehat{\mathbf{U}}, \mathbf{S}, \boldsymbol{\alpha}, \boldsymbol{\beta}\right) P\left(\boldsymbol{\beta}\right) P\left(\boldsymbol{\alpha} \mid \lambda(r)\right) P\left(\lambda(r)\right).$$

$\boldsymbol{\alpha}, \boldsymbol{\beta}$  and  $\lambda(r)$  are stochastic, and induce a **posterior distribution on  $r$**  through the fairness constraint expression.

$P(r \mid \widehat{\mathbf{U}}, \mathbf{S} < r_0)$  then controls the fairness of the model.

# MARGINAL PLOTS FROM THE POSTERIOR DISTRIBUTION



### PROs:

- Putting Zellner's g-priors on  $\alpha$  and  $\beta$  gives the MLE for  $g \rightarrow \infty$ , which is useful for comparisons.
- Posterior inference (jointly) on  $\alpha$ ,  $\beta$  and  $r$  via MCMC.
- Posterior sensitivity analysis on the choice of  $r$ .

### CONs:

- It appears to be less numerically stable when  $\mathbf{X}$  and  $\mathbf{S}$  are collinear.
- MCMC convergence to its stationary distribution can be problematic in this and other less-than-smooth settings.

How can we approximate the MCMC inference in the original FGRRM?

✓ INTRODUCTION

✓ FAIR LINEAR MODELS

✓ FAIR RIDGE REGRESSION

✓ BAYESIAN FAIR REGRESSION

→ UNCERTAINTY QUANTIFICATION

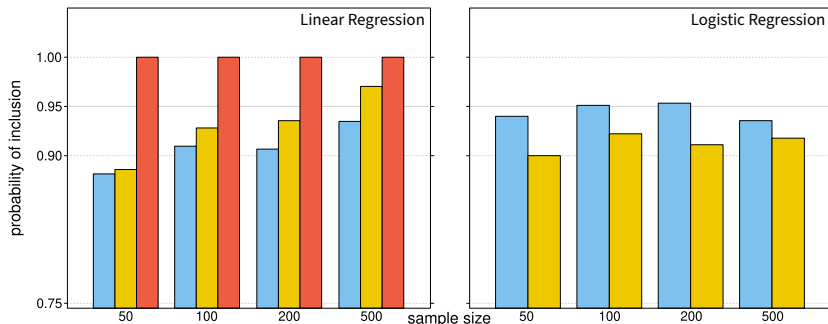
CONCLUSIONS & ACKNOWLEDGEMENTS

# WHAT CAN WE DO ON THE FREQUENTIST SIDE?

Much has been written on **confidence intervals for RR coefficients**.

For FGRRM:

- Nonparametric bootstrap [3] works best.
- Double bootstrap [14] is costly, and works only for FRRM.
- Residual bootstrap [8] does not work.

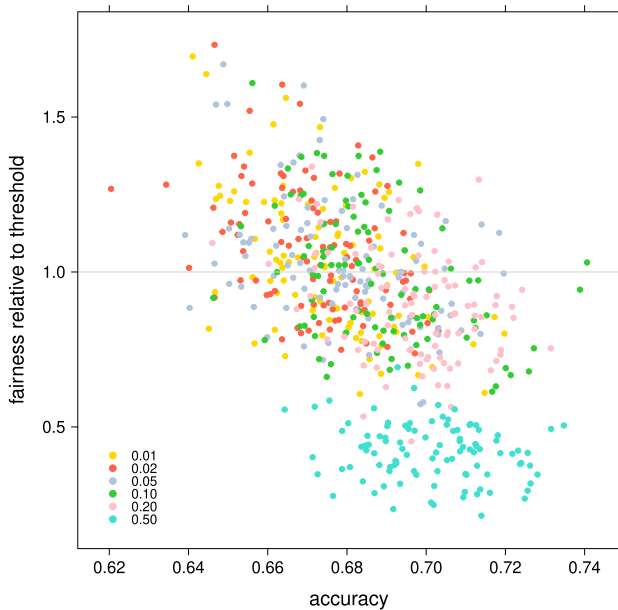


At some computational cost, nested bootstrap allows us to **marginalise**  $\alpha, \beta, \lambda(r)$  and training-validation split variability.

---

1. For  $b = 1, \dots, B_1$ :
    - 1.1 Draw a bootstrap sample  $\mathbf{y}^{(b)}, \mathbf{X}^{(b)}$  and  $\mathbf{S}^{(b)}$  from  $\mathbf{y}, \mathbf{X}, \mathbf{S}$ .
    - 1.2 Split training  $(\mathbf{y}_{\text{TR}}^{(b)}, \mathbf{X}_{\text{TR}}^{(b)}, \mathbf{S}_{\text{TR}}^{(b)})$  and validation  $(\mathbf{y}_{\text{VA}}^{(b)}, \mathbf{X}_{\text{VA}}^{(b)}, \mathbf{S}_{\text{VA}}^{(b)})$ .
    - 1.3 For  $b' = 1, \dots, B_2$ :
      - 1.3.1 Draw a bootstrap sample  $\mathbf{y}_{\text{TR}}^{(b')}, \mathbf{X}_{\text{TR}}^{(b')}, \mathbf{S}_{\text{TR}}^{(b')}$  from  $\mathbf{y}_{\text{TR}}^{(b)}, \mathbf{X}_{\text{TR}}^{(b)}, \mathbf{S}_{\text{TR}}^{(b)}$ .
      - 1.3.2 Estimate a model  $\mathcal{M}^{(b')}$  with fairness  $r$  from  $\mathbf{y}_{\text{TR}}^{(b')}, \mathbf{X}_{\text{TR}}^{(b')}, \mathbf{S}_{\text{TR}}^{(b')}$ .
      - 1.3.3 Estimate the fairness loss for  $\mathcal{M}^{(b')}$  on the validation set  $(\mathbf{y}_{\text{VA}}^{(b)}, \mathbf{X}_{\text{VA}}^{(b)}, \mathbf{S}_{\text{VA}}^{(b)})$ .
      - 1.3.4 Predict  $\mathbf{y}_{\text{VA}}^{(b)}$  from  $\mathbf{X}_{\text{VA}}^{(b)}, \mathbf{S}_{\text{VA}}^{(b)}$  using  $\mathcal{M}^{(b')}$  to estimate predictive accuracy.
  2. Assess the joint distribution of empirical fairness and predictive accuracy either by visual inspection or using a pinball loss. [1]
-

# FAIRNESS-ACCURACY TRADE-OFF PLOT



✓ INTRODUCTION

✓ FAIR LINEAR MODELS

✓ FAIR RIDGE REGRESSION

✓ BAYESIAN FAIR REGRESSION

✓ UNCERTAINTY QUANTIFICATION

→ CONCLUSIONS & ACKNOWLEDGEMENTS



- Fairness is increasingly a concern as machine learning models become an integral part of automated decision support systems.
- Explainable AI investigates black-box models such as neural networks, but simpler models are also in common use and should be made fair.
- The literature studies fairness as an optimisation problem, producing models whose statistical properties and best practices are unknown.
- Classical statistics provides all the tools to formulate versatile fair models that we know how to interpret and use.
- Uncertainty quantification of fair models remains key, and is challenging because of the nature of the fairness constraint.



Università  
della  
Svizzera  
italiana

Ernst Wit

Faculty of Informatics

*Università della Svizzera Italiana*



SAPIENZA  
UNIVERSITÀ DI ROMA

Francesca Panero

Department of Methods and Models for  
Economics, Territory and Finance

*Sapienza Università di Roma*



Manuel Proissl

Quantum Industry Applications Lead

*IBM*

Software: <https://cran.r-project.org/web/packages/fairml/>

THAT'S ALL!

HAPPY TO DISCUSS IN MORE DETAIL.

[scutari@bnlearn.com](mailto:scutari@bnlearn.com)

# REFERENCES I

- ◆ Y. Chung, W. Neiswanger, I. Char, and J. Schneider.  
[1] [Beyond Pinball Loss: Quantile Methods for Calibrated Uncertainty Quantification.](#)  
In *Advances in Neural Information Processing Systems*, pages 10971–10984, 2021.
- ◆ B. D'Alessandro, C. O'Neil, and T. LaGatta.  
[2] [Conscientious Classification: A Data Scientist's Guide to Discrimination-Aware Classification.](#)  
*Big Data*, 5(2):120–134, 2017.
- ◆ A. C. Davison and D. V. Hinkley.  
[3] [Bootstrap Methods and Their Application.](#)  
Cambridge University Press, 1997.
- ◆ E. Del Barrio, P. Gordaliza, and J. M. Loubes.  
[4] [Review of Mathematical Frameworks for Fairness in Machine Learning](#), 2020.
- ◆ J. Friedman, T. Hastie, and R. Tibshirani.  
[5] [Regularization Paths for Generalized Linear Models via Coordinate Descent.](#)  
*Journal of Statistical Software*, 33(1):1–22, 2010.
- ◆ A. E. Hoerl and R. W. Kennard.  
[6] [Ridge Regression: Biased Estimation for Nonorthogonal Problems.](#)  
*Technometrics*, 12(1):55–67, 1970.

- ◆ J. Komiyama, A. Takeda, J. Honda, and H. Shimao.  
[7] [Nonconvex Optimization for Regression with Fairness Constraints.](#)  
*Proceedings of Machine Learning Research*, 80:2737–2746, 2018.  
35th International Conference on Machine Learning.
- ◆ P. Legendre.  
[8] [Comparison of Permutation Methods for the Partial Correlation and Partial Mantel Tests.](#)  
*Journal of Statistical Computation and Simulation*, 67(1):37–73, 2000.
- ◆ P. McCullagh and J. A. Nelder.  
[9] [Generalized Linear Models.](#)  
CRC press, 2nd edition, 1989.
- ◆ N. Mehrabi, F. Morstatter, N. Saxena, et al.  
[10] [A Survey on Bias and Fairness in Machine Learning.](#)  
*ACM Computing Surveys*, 54(6):115, 2021.
- ◆ D. Pessach and E. Shmueli.  
[11] [A Review on Fairness in Machine Learning.](#)  
*ACM Computing Surveys*, 55(3):51, 2022.

- ◆ N. Simon, J. Friedman, T. Hastie, and R. Tibshirani.  
[12] [Regularization Paths for Cox's Proportional Hazards Model via Coordinate Descent.](#)  
*Journal of Statistical Software*, 39(5):1–13, 2011.
- ◆ J. K. Tay, B. Narasimhan, and T. Hastie.  
[13] [Elastic Net Regularization Paths for All Generalized Linear Models.](#)  
*Journal of Statistical Software*, 106(1):1–31, 2023.
- ◆ H. Vinod.  
[14] [Double Bootstrap for Shrinkage Estimators.](#)  
*Journal of Econometrics*, 68(2):287–302, 1995.