

Wakacje 2020

Jakub Wiśniewski

Plan prezentacji

- Fairness
- fairmodels
- Blog
- Moduł fairness w dalex

Czym jest fairness - podsumowanie

- Część Responsible ML
- Zajmuje się oceną sprawiedliwości decyzji podejmowanych przez modele
- W perspektywie fairness model to 3 wektory:
 - y - target (binary)
 - \hat{y} - response
 - A - protected

Czy można jednoznacznie ocenić fairness? Na co uważać?

- Nie można, ale trzeba próbować!
- Zasada 4/5 (80%) w przypadku *statistical parity* [UNIFORM GUIDELINES ON EMPLOYEE SELECTION PROCEDURES (1978)]
- Możliwe sprawdzanie nie tylko binarnych *sensitive attributes* ale także tych niebinarnych oraz ich przecięć

Gender Shades: Intersectional Accuracy Disparities in
Commercial Gender Classification
Joy Buolamwini, Timnit Gebru 2018

4.1. Key Findings on Evaluated Classifiers

- All classifiers perform better on male faces than female faces (8.1% – 20.6% difference in error rate)
- All classifiers perform better on lighter faces than darker faces (11.8% – 19.2% difference in error rate)
- All classifiers perform worst on darker female faces (20.8% – 34.7% error rate)
- Microsoft and IBM classifiers perform best on lighter male faces (error rates of 0.0% and 0.3% respectively)
- Face++ classifiers perform best on darker male faces (0.7% error rate)
- The maximum difference in error rate between the best and worst classified groups is 34.4%

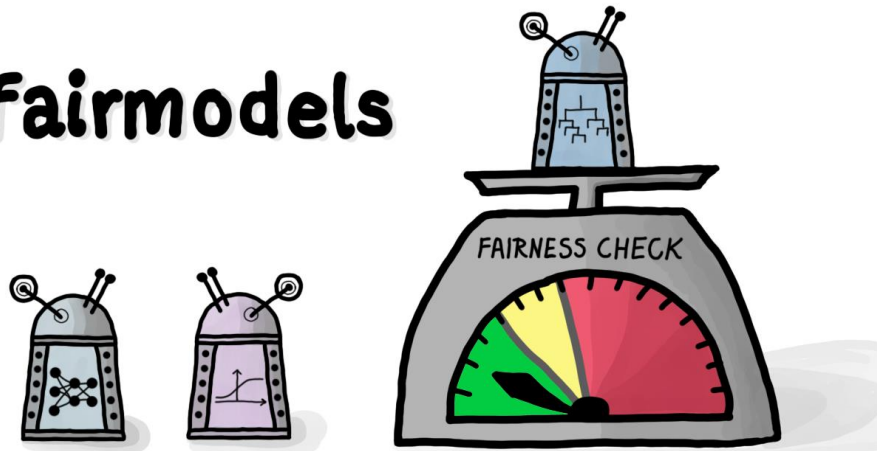
fairmodels

- Jedyne* takie narzędzie w R
- Na CRAN od 20 sierpnia
- Pojawia się w *top 40 CRAN new packages (August)*
- Wylistowane w “*awesome-machine-learning-interpretability*”

☆ Star 34

CRAN 0.2.2 downloads 727/month

fairmodels

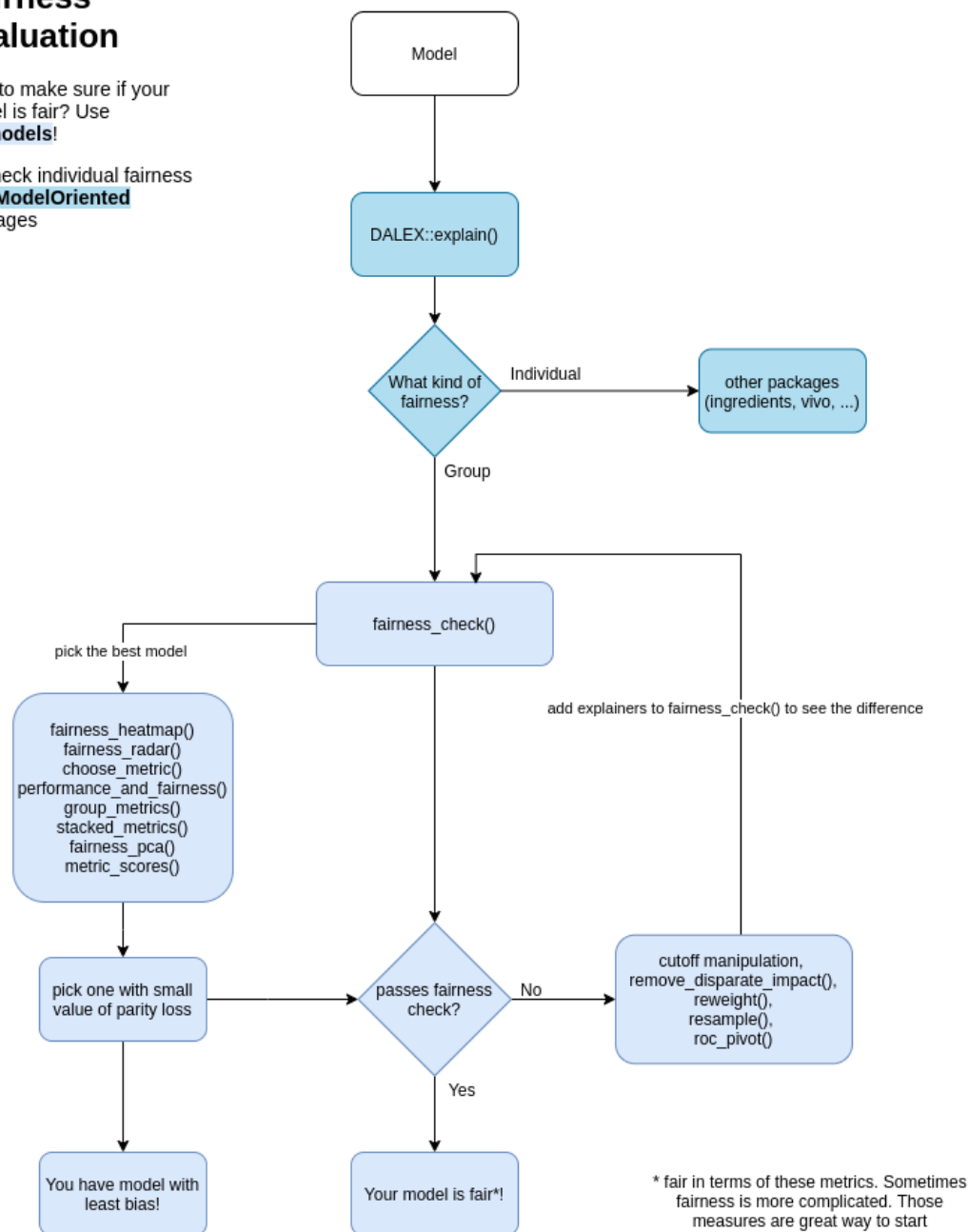


* które obsługuje modele oraz jest używalne

Fairness evaluation

How to make sure if your model is fair? Use **fairmodels!**

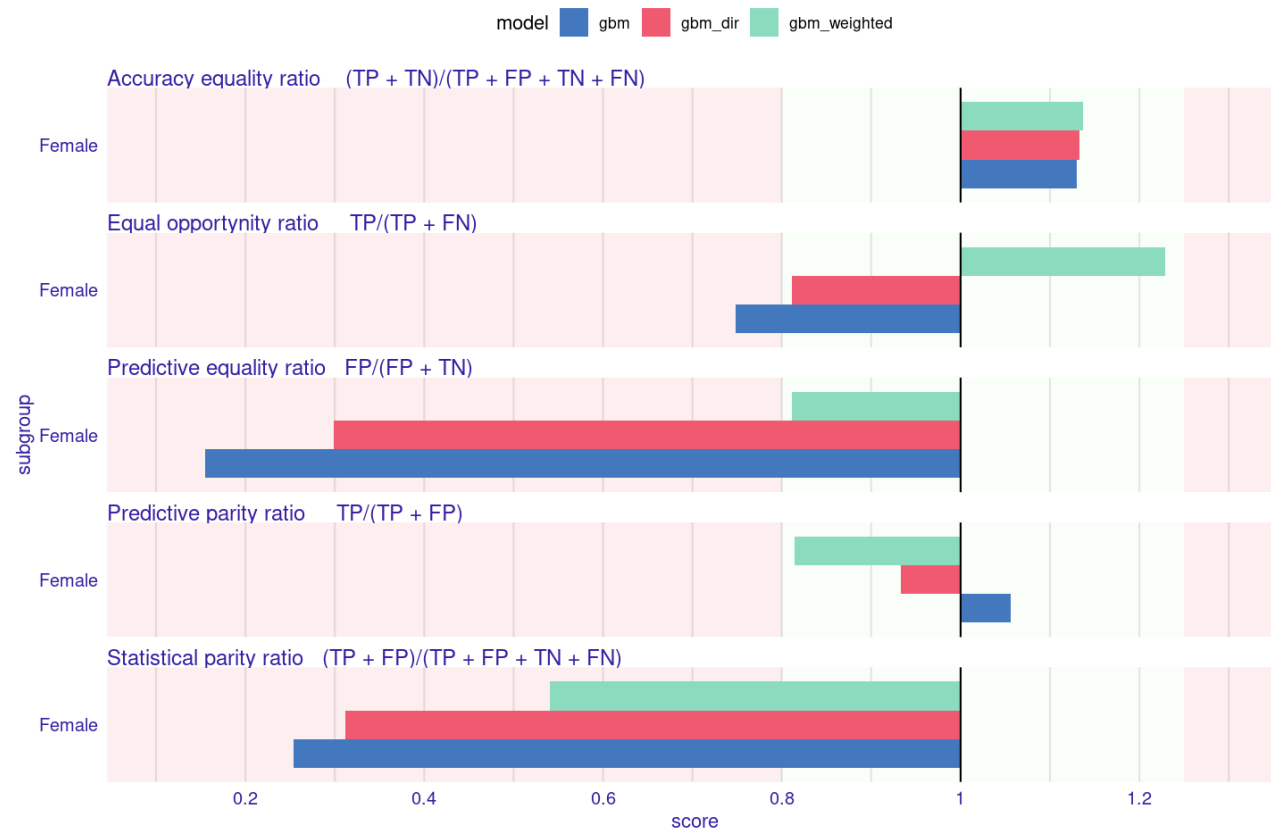
Or check individual fairness with **ModelOriented** packages



Jak używać fairmodels?

Fairness check

Created with gbm_weighted, gbm, gbm_dir



* fair in terms of these metrics. Sometimes fairness is more complicated. Those measures are great way to start

Metryki

- Metryki bazujące na macierzy pomyłek
- Pochodzące z nich metryki fairness
 - Equal opportunity (TPR)
 - Predictive parity (PPV)
 - Predictive equality (FPR)
 - Accuracy (ACC)
 - Statistical parity ($STP = (TP+FP)/(TP+FP+TN+FN)$)

Parity loss – dlaczego zmiana?

- Kiedyś:

$$TPR_{parity-loss} = \sum_i |TPR_{A=i} - TPR_{A=privileged}|$$

- Teraz:

$$TPR_{parity-loss} = \sum_i \left| \ln \left(\frac{TPR_{A=i}}{TPR_{A=privileged}} \right) \right|$$

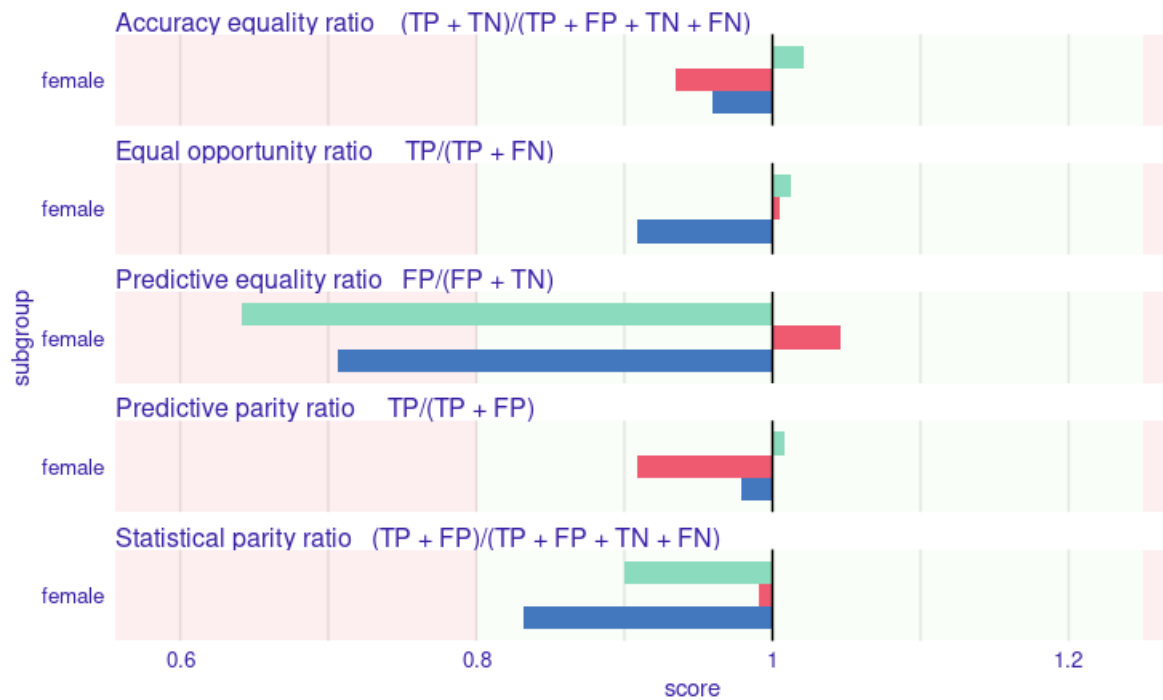
- Po co?
 - Aby dopasować się do Amerykańskiego prawa (four-fifths rule)
 - Aby dla dużych oraz małych wartości metryk interpretacja była ta sama.

Wizualizacja fairness

Fairness check

Created with lm_changed, lm, ranger

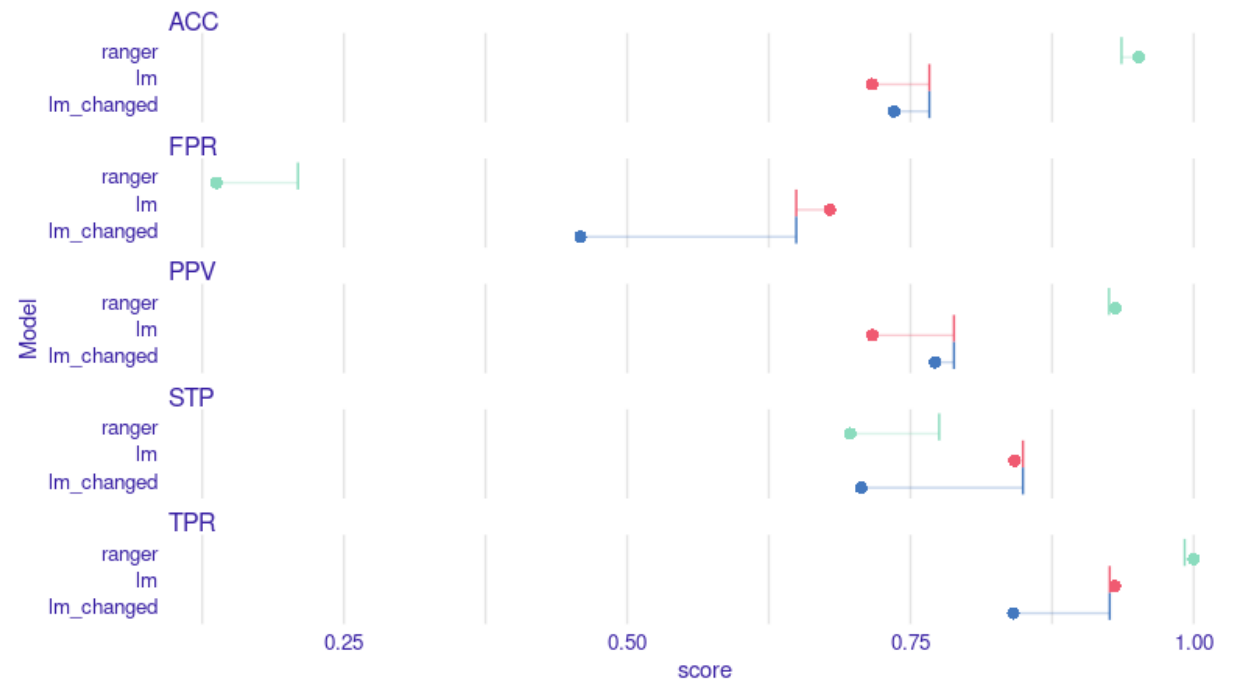
model lm lm_changed ranger



Metric scores plot

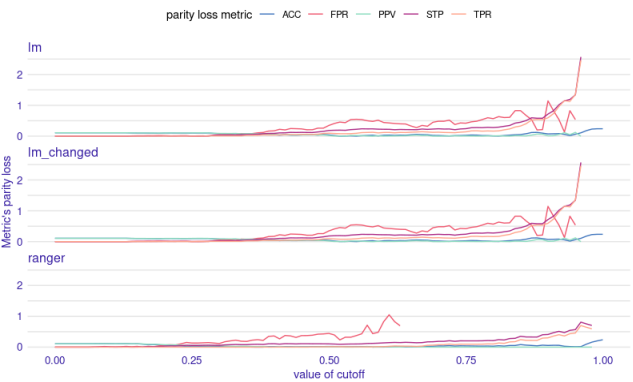
Created with lm_changed, lm, ranger

subgroup ● female model ● lm ● lm_changed ● ranger

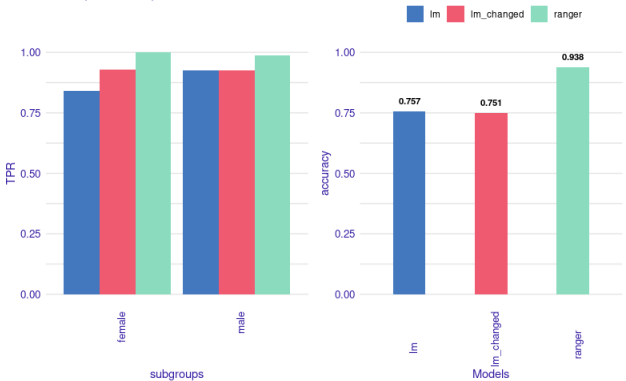


All cutoffs plot

created with lm_changed, lm, ranger

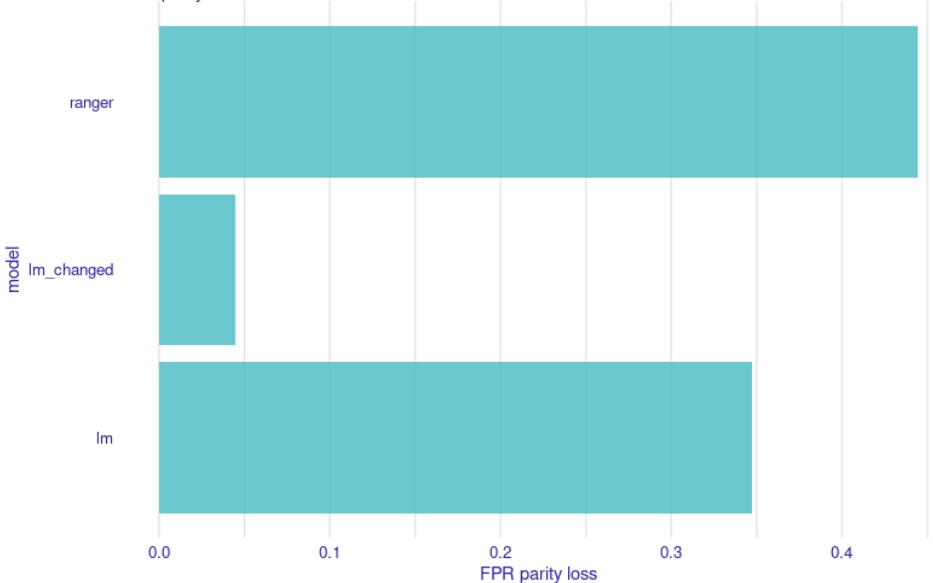


Group metric plot



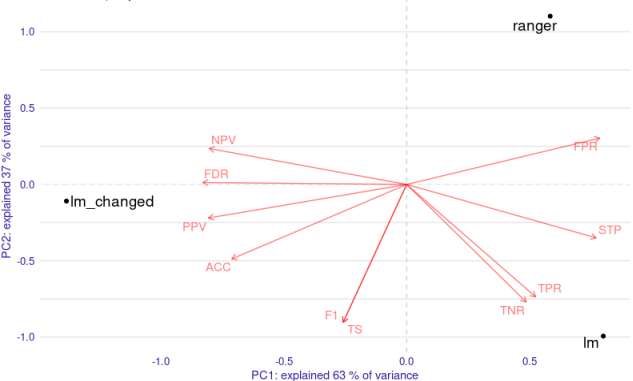
Chosen metric plot

FPR parity loss in all models

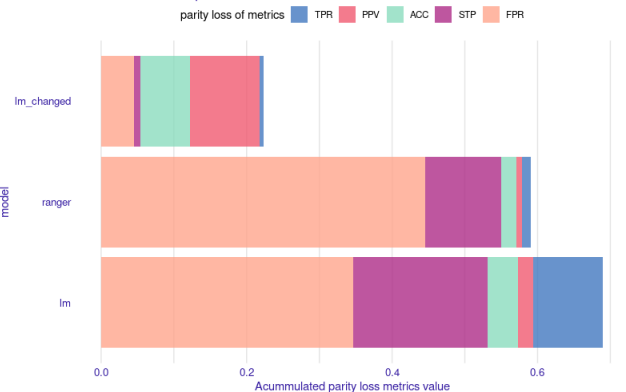


Fairness PCA plot

created with parity loss metrics

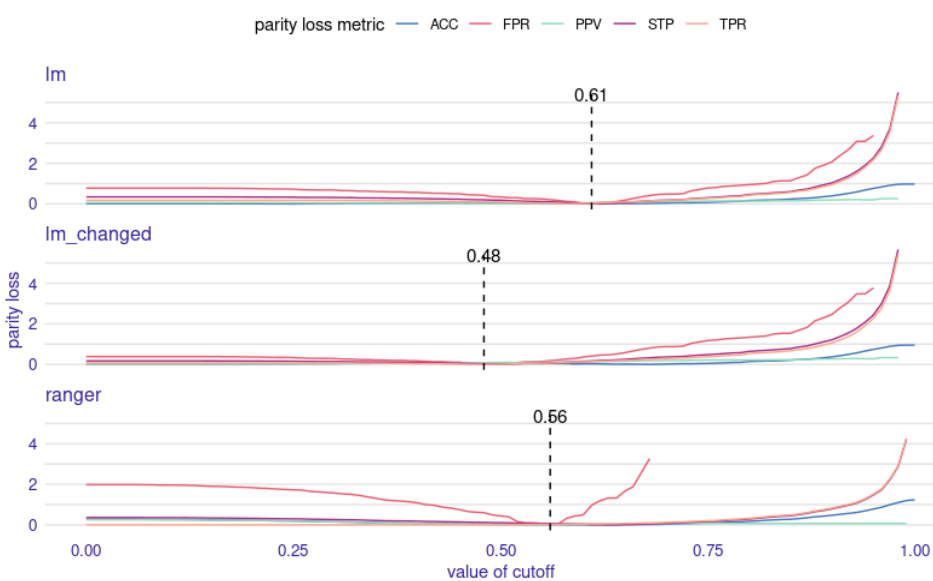


Stacked Metric plot

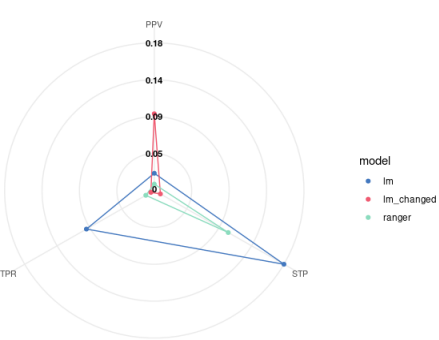


Ceteris paribus cutoff plot

Based on male

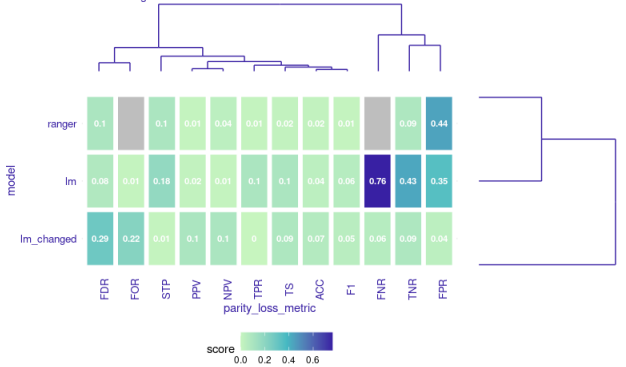


Parity loss metric radar plot



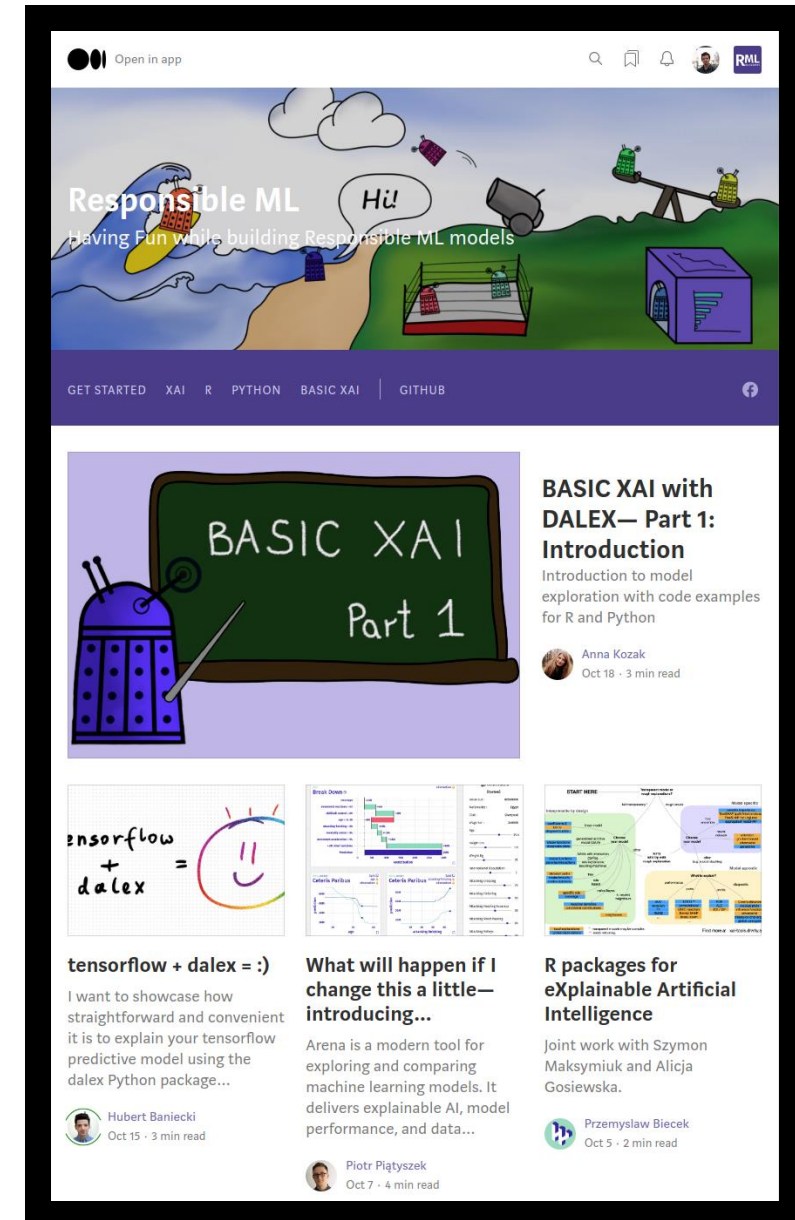
Heatmap

With dendrograms



Blog - miesięcznie

- 3027 wyświetleń
- Ponad 30 godzin spędzonych na czytaniu
- Średnio 63 unikalnych użytkowników dziennie
- Maksymalnie 294 użytkowników dziennie



Blog – statystyki

- Najwięcej zasięgów miały blogi Kasi oraz S&A&P
- Duży wpływ ma LinkedIn Przemka: +230 wyświetleń

OCTOBER 2020

| | | | | |
|---|----------|-----|-----|----|
| BASIC XAI with DALEX— Part 1: Introduction <small>3 min read · In ResponsibleML · View story · Details</small> | 176 +23K | 115 | 65% | 4 |
| tensorflow + dalex = :) <small>3 min read · In ResponsibleML · View story · Details</small> | 162 | 95 | 59% | 6 |
| What will happen if I change this a little— i... <small>4 min read · In ResponsibleML · View story · Details</small> | 245 +34K | 108 | 44% | 5 |
| R packages for eXplainable Artificial Intellig... <small>2 min read · In ResponsibleML · View story · Details</small> | 800 +28K | 468 | 59% | 10 |
| Imputing missing data in mlr3 with EMMA <small>4 min read · In ResponsibleML · View story · Details</small> | 261 +24K | 137 | 52% | 6 |

SEPTEMBER 2020

| | | | | |
|--|----------|-----|-----|---|
| What is new in fairmodels? <small>4 min read · In ResponsibleML · View story · Details</small> | 231 +35K | 115 | 50% | 6 |
| What's new in DALEX and DALEXtra <small>7 min read · In ResponsibleML · View story · Details</small> | 435 +36K | 155 | 36% | 8 |
| Explaining models with Triplot, part 1 <small>4 min read · In ResponsibleML · View story · Details</small> | 658 +41K | 325 | 49% | 8 |
| How to use the Arena for exploration of ML ... <small>3 min read · In ResponsibleML · View story · Details</small> | 119 | 82 | 69% | 7 |
| Introduction to ResponsibleML publication <small>3 min read · In ResponsibleML · View story · Details</small> | 151 | 87 | 58% | 1 |

← RXAI →

← Triplot

VIEWS BY TRAFFIC SOURCE 800

| | |
|---------------------------------|------|
| Internal ⓘ | 0% |
| External referrals | 100% |
| RSS readers (full text) | 28K |
| linkedin.com | 266 |
| email, IM, and direct | 182 |
| r-bloggers.com | 152 |
| Android device (not Medium app) | 65 |
| google.com | 16 |
| feedly.com | 16 |
| news.google.com | 6 |
| twitter.com | 5 |
| mail.google.com | 4 |
| All other external referrals | 12 |

Moduł do dalex (in progress)



LEPIEJ ZAIMPLEMENTOWANY
FAIRMODELS



MNIEJ FUNKCJONALNOŚCI,
SKUPIENIE SIĘ NA GŁÓWNYCH
FEATURES

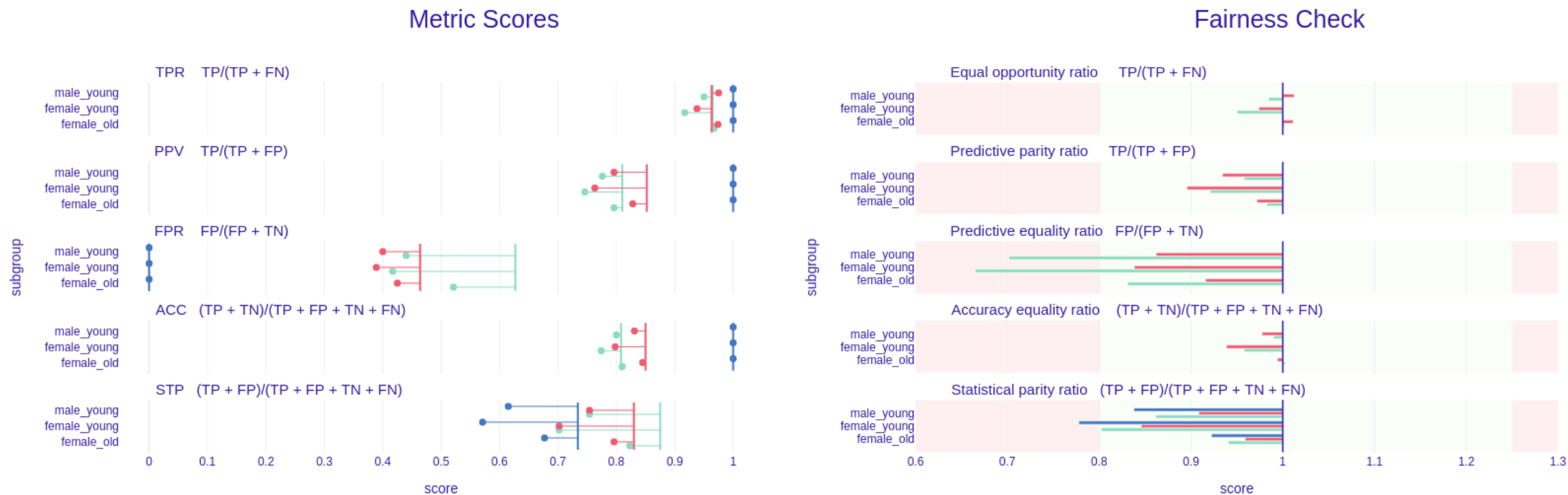


DOBRE PODSTAWY POD
INDIVIDUAL FAIRNESS & FAIRNESS
IN REGRESSION



ZAIMPLEMENTOWANY
BAKCEND ORAZ 2 (MOIM ZDANIEM)
NAJWAŻNIEJSZE WYKRESY

Interaktywne wykresy plotly



Koniec :)