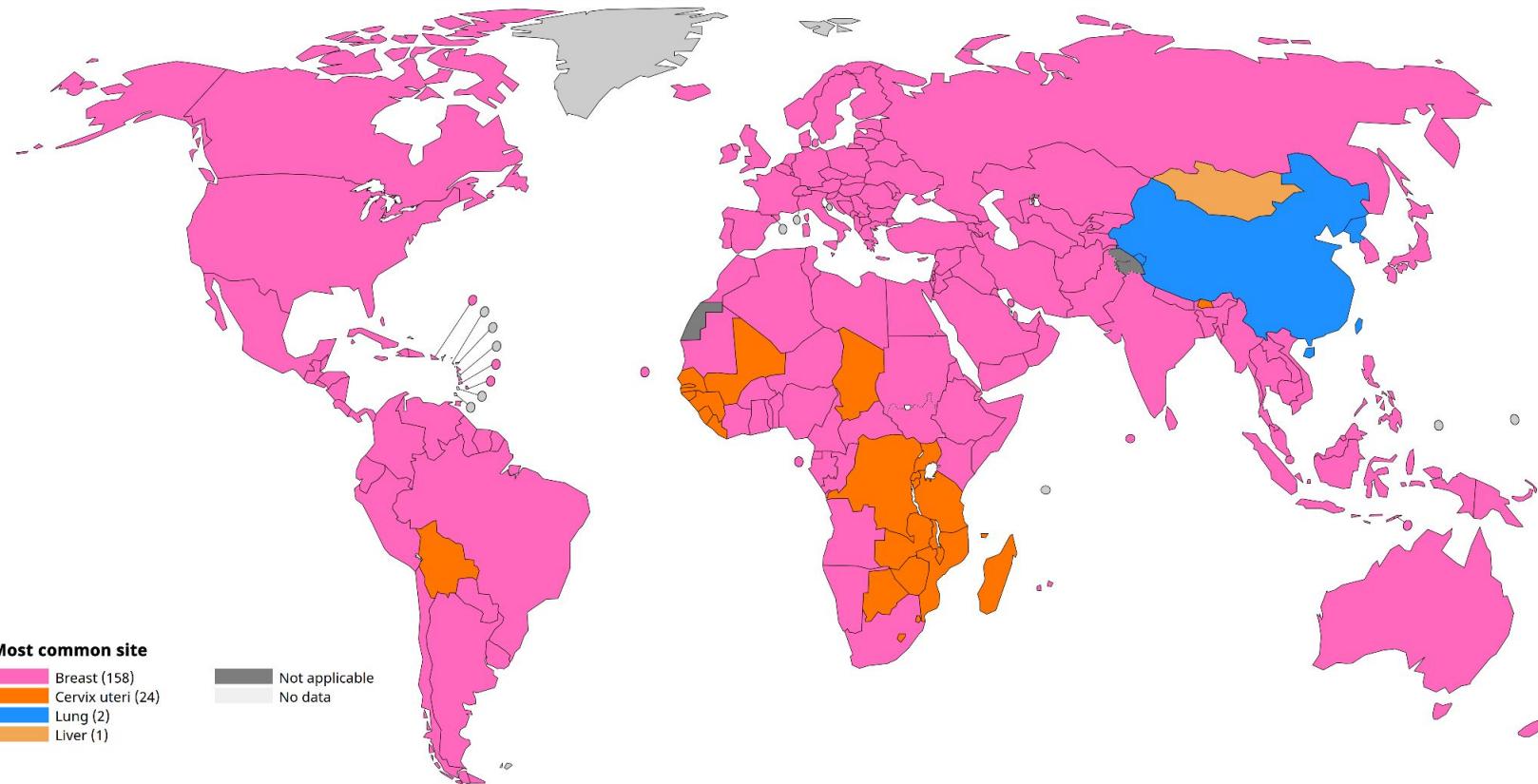


Precision Medicine for All: Using Tidymodels to Validate Breast Cancer PRS in Brazil

Flávia E. Rius, PhD

Data Scientist - Mendelics

Most common site per country, Absolute numbers, Incidence, Females, in 2022 (excl. NMSC)



Most common site

Pink	Breast (158)
Orange	Cervix uteri (24)
Blue	Lung (2)
Tan	Liver (1)

Not applicable
No data

All rights reserved. The designations employed and the presentation of the material in this publication do not imply the expression of any opinion whatsoever on the part of the World Health Organization / International Agency for Research on Cancer concerning the legal status of any country, territory, city or area or of its authorities, or concerning the delimitation of its frontiers or boundaries. Dotted and dashed lines on maps represent approximate borderlines for which there may not yet be full agreement.

Cancer TODAY | IARC

<https://gco.iarc.who.int/today>

Data version: GLOBOCAN 2022 (version 1.1) - 08.02.2024

© All Rights Reserved 2025

International Agency
for Research on Cancer



Main Risk Factors

- Sex ♀

Main Risk Factors

- Sex ♀
- Age ↑

Main Risk Factors

- Sex ♀
- Age ↑
- Lifestyle 🍷

Main Risk Factors

- Sex ♀
- Age ↑
- Lifestyle 🍷
- Hormonal

Main Risk Factors

- Sex ♀
- Age ↑
- Lifestyle 🍷
- Hormonal
- Family history

Main Risk Factors

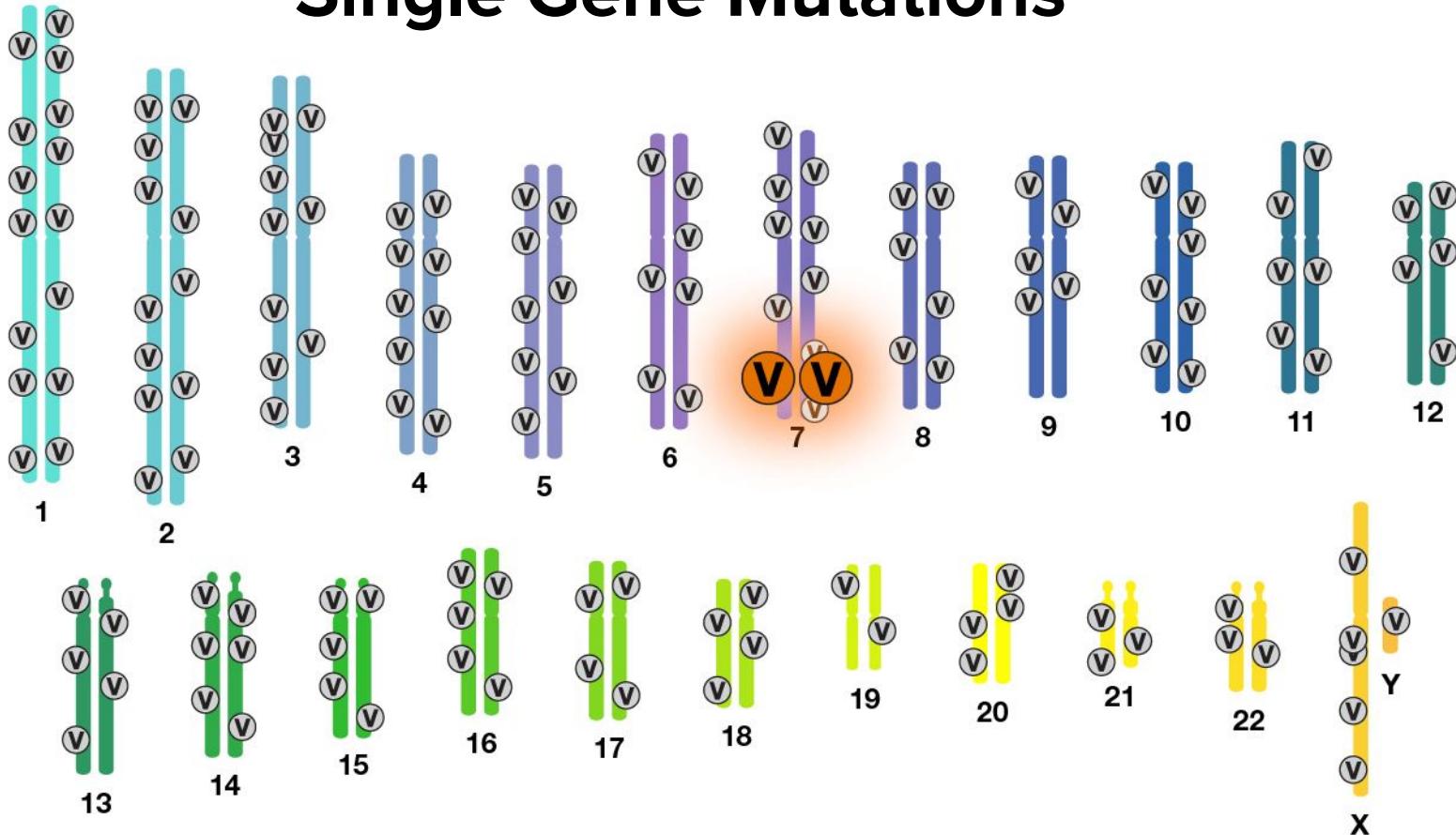
- Sex 
- Age 
- Lifestyle 
- Hormonal
- Family history
- Genetics 

Genetics of Breast Cancer

- *BRCA1* mutation

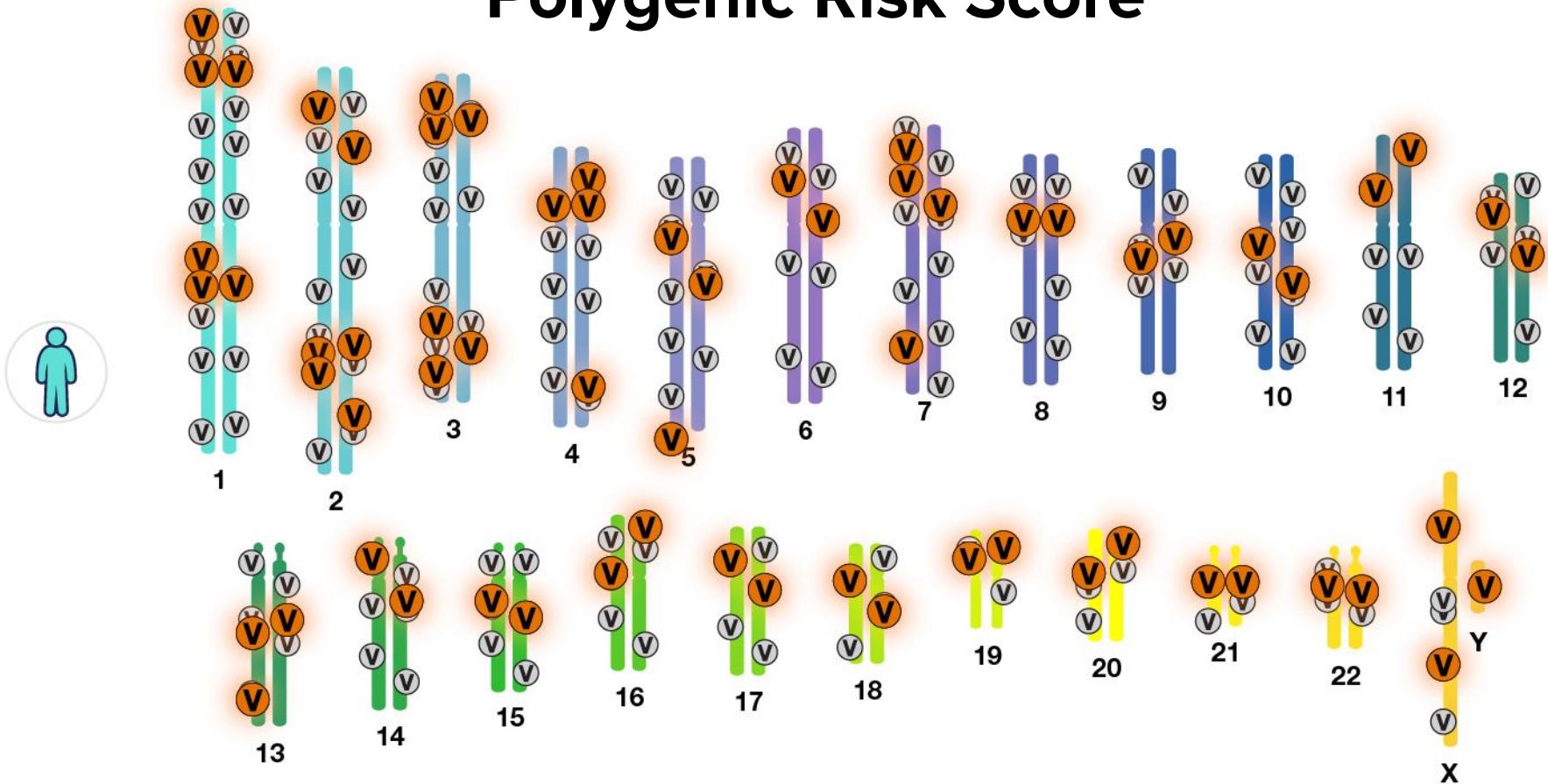


Single Gene Mutations

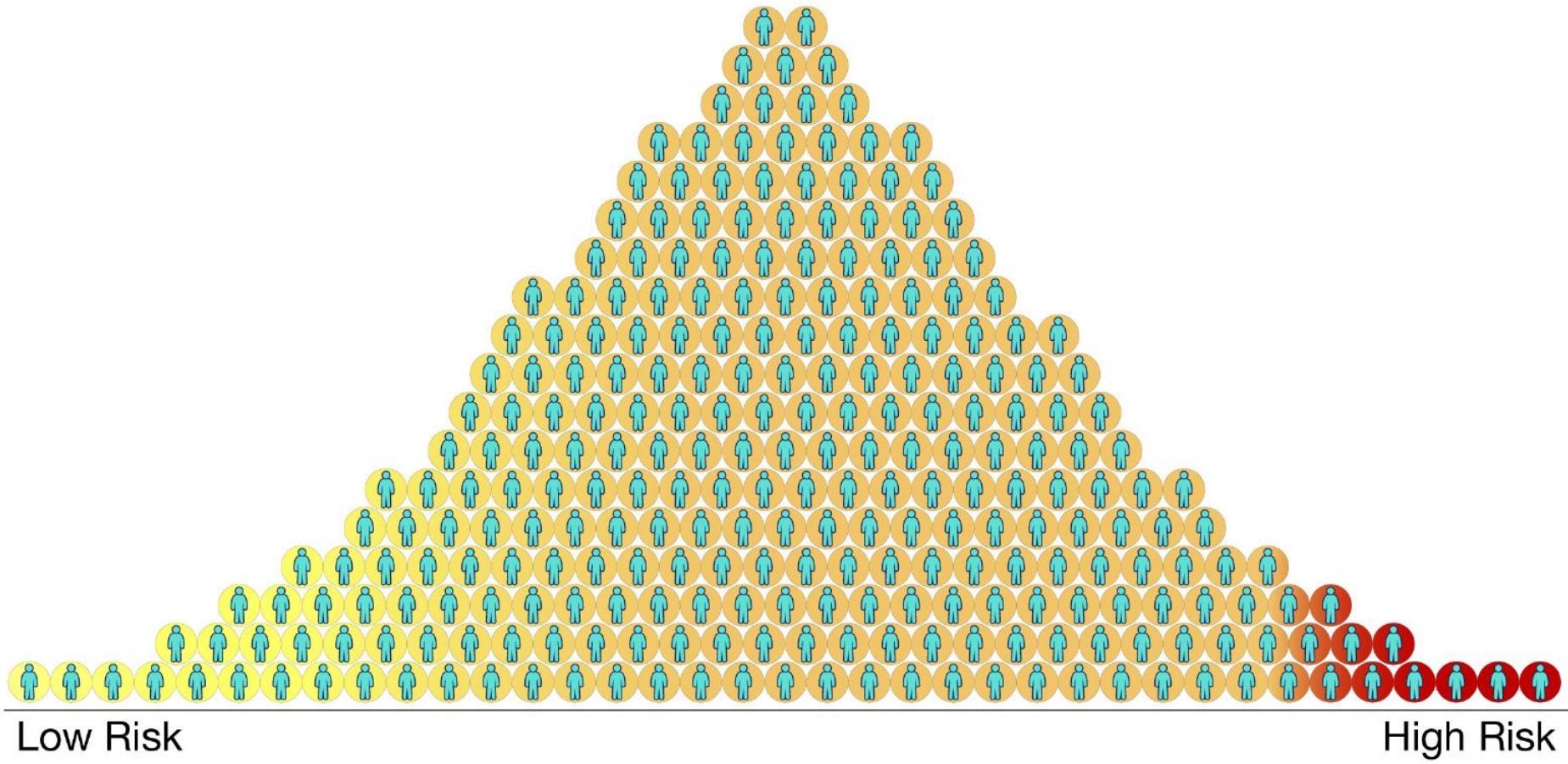


Courtesy: [National Human Genome Research Institute](#)

Polygenic Risk Score



Courtesy: [National Human Genome Research Institute](#)



Courtesy: [National Human Genome Research Institute](#)

Genome Wide Association Studies (GWAS)



Total GWAS participants diversity

Version 1.0.0. Last check for data: 2025-09-05 00:34:26 .



African American or Afro-Caribbean Hispanic or Latin American Other/Mixed



Breast Cancer PRS Validation in Brazilians



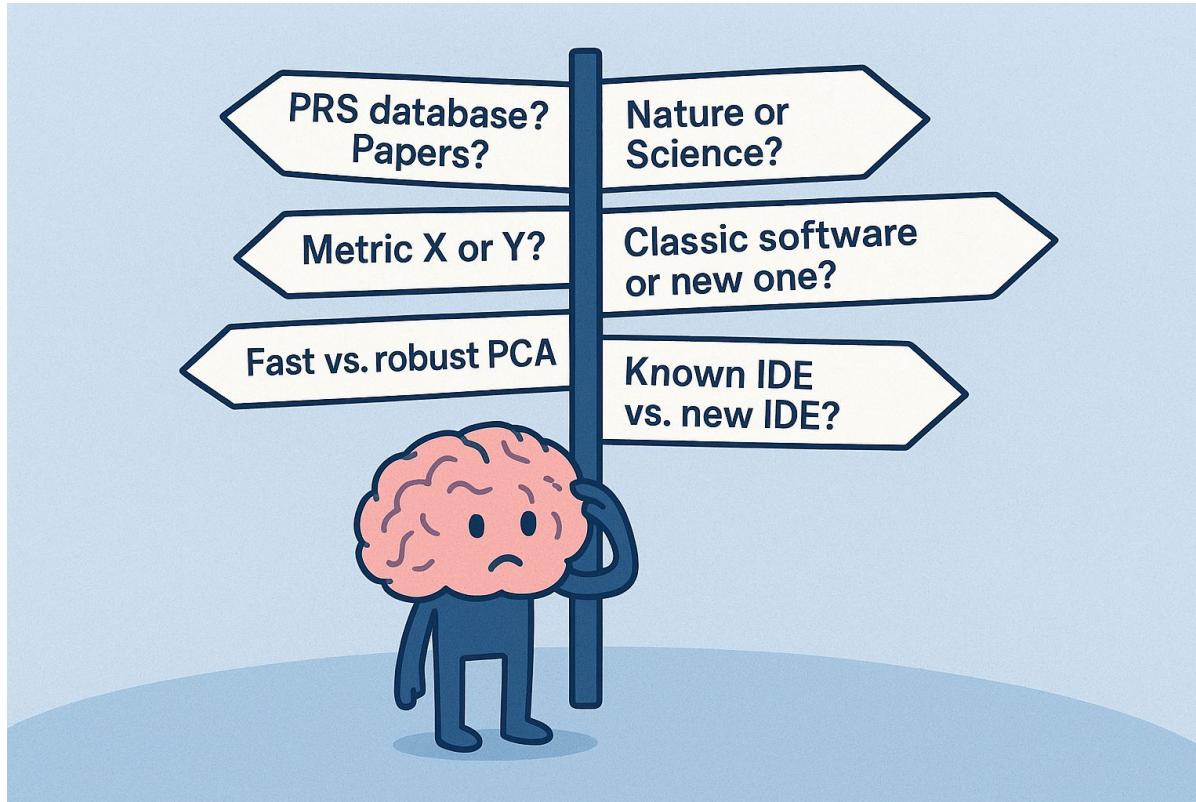


Greek Yogurt





Questions, questions, questions



**Can European-based PRS be used in Brazil
to identify a high breast cancer risk?**

Tidymodels



Tidymodels

- Integratable tools



Tidymodels

- Integratable tools
- Compatible with {tidyverse}



Tidymodels

- Integratable tools
- Compatible with {tidyverse}
- Concept of reusable building blocks



Read the data

```
1 breast <- readr::read_csv(here::here(  
2   "data",  
3   "output",  
4   "v2",  
5   "breast_preprocessed_no_mut_v2.csv"  
6 ))  
7  
8 head(breast, 3)  
  
# A tibble: 3 × 17  
#> id    status sex    prs_broad prs_313 prs_3820 prs_ukbb pc1     pc2     pc3  
#> <chr> <fct>  <fct>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>    <dbl>  
1 1      1       F        0.501    0.488    0.505  0.000230  343. -2120. -487.  
2 2      0       F        0.504    0.488    0.506  0.000234  555. -2568.  297.  
3 6      1       F        0.508    0.488    0.502  0.00409  5502. -1035. -602.  
#> # i 7 more variables: pc4 <dbl>, pc5 <dbl>, pc6 <dbl>, pc7 <dbl>, pc8 <dbl>,  
#> #   pc9 <dbl>, pc10 <dbl>
```





Split the sample

```
1 set.seed(28)
2
3 split <- initial_split(breast, strata = status, prop = 0.75)
4
5 train <- training(split)
6 test <- testing(split)
```

Normalize PRSs

```
1 control_stats_train <- get_control_stats(train)
2
3 train_ctrl_stats <- train |>
4   # create norm prss
5   dplyr::mutate(
6     prs_3820_norm = (prs_3820 - control_stats_train$prs_3820_mean) /
7       control_stats_train$prs_3820_sd,
8     prs_broad_norm = (prs_broad - control_stats_train$prs_broad_mean) /
9       control_stats_train$prs_broad_sd,
10    prs_313_norm = (prs_313 - control_stats_train$prs_313_mean) /
11      control_stats_train$prs_313_sd,
12    prs_ukbb_norm = (prs_ukbb - control_stats_train$prs_ukbb_mean) /
13      control_stats_train$prs_ukbb_sd
14  ) |>
15  select(
16    !all_of(c("prs_broad", "prs_313", "prs_3820", "prs_ukbb")))
17  )
```





Adjust the data prior to modeling

```
1 rec_3820 <- recipe(  
2   status ~ .,  
3   data = train_ctrl_stats  
4 ) |>  
5   step_rm(  
6     id,  
7     sex,  
8     prs_313_norm,  
9     prs_broad_norm,  
10    prs_ukbb_norm  
11  ) |>  
12  step_normalize(starts_with("pc"))
```



Adjust the data prior to modeling

```
1 rec_3820 <- recipe(  
2   status ~ .,  
3   data = train_ctrl_stats  
4 ) |>  
5   step_rm(  
6     id,  
7     sex,  
8     prs_313_norm,  
9     prs_broad_norm,  
10    prs_ukbb_norm  
11  ) |>  
12  step_normalize(starts_with("pc"))
```



Adjust the data prior to modeling

```
1 rec_3820 <- recipe(  
2   status ~ .,  
3   data = train_ctrl_stats  
4 ) |>  
5   step_rm(  
6     id,  
7     sex,  
8     prs_313_norm,  
9     prs_broad_norm,  
10    prs_ukbb_norm  
11  ) |>  
12  step_normalize(starts_with("pc"))
```



Set engine of the model

```
1 log_reg <- logistic_reg() |>  
2   set_engine("glm") |>  
3   set_mode("classification")
```



Put everything together

```
1 wflow_3820 <- workflow() |>
2   add_model(log_reg) |>
3   add_recipe(rec_3820)
```

And fit

```
1 fit_3820 <- wflow_3820 |>
2   fit(data = train_ctrl_stats)
```



Get statistics

```
1 tidy(fit_3820, exponentiate = T) |>  
2   filter(grepl("prs", term)) |>  
3   select(!all_of("std.error", "statistic"))
```

```
# A tibble: 1 × 3  
  term      estimate  p.value  
  <chr>        <dbl>    <dbl>  
1 prs_3820_norm     1.44 3.15e-72
```



Update recipe to remove PCs

```
1 new_rec_3820 <- recipe(  
2   status ~ .,  
3   data = train_ctrl_stats  
4 ) |>  
5 step_rm(  
6   id,  
7   sex,  
8   prs_313_norm,  
9   prs_broad_norm,  
10  prs_ukbb_norm,  
11  starts_with("pc"))
```

```
1 new_wfflow_rec_3820 <- wfflow_3820 |>  
2 update_recipe(new_rec_3820)  
3  
4 fit_3820 <- new_wfflow_rec_3820 |>  
5 fit(data = train_ctrl_stats)
```



Get metrics

```
1 # Predict probabilities
2 pred_prob_3820 <- predict(new_fit_3820, new_data = test_ctrl_stats,
3   type = "prob")
4
5 # Add to testing data
6 results_3820 <- test_ctrl_stats |>
7   bind_cols(pred_prob_3820)
8
9 # Obtain ROC AUC
10 roc_auc(results_3820, truth = status, .pred_0)

# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>   <chr>        <dbl>
1 roc_auc binary      0.593
```



Repeat with other PRSs

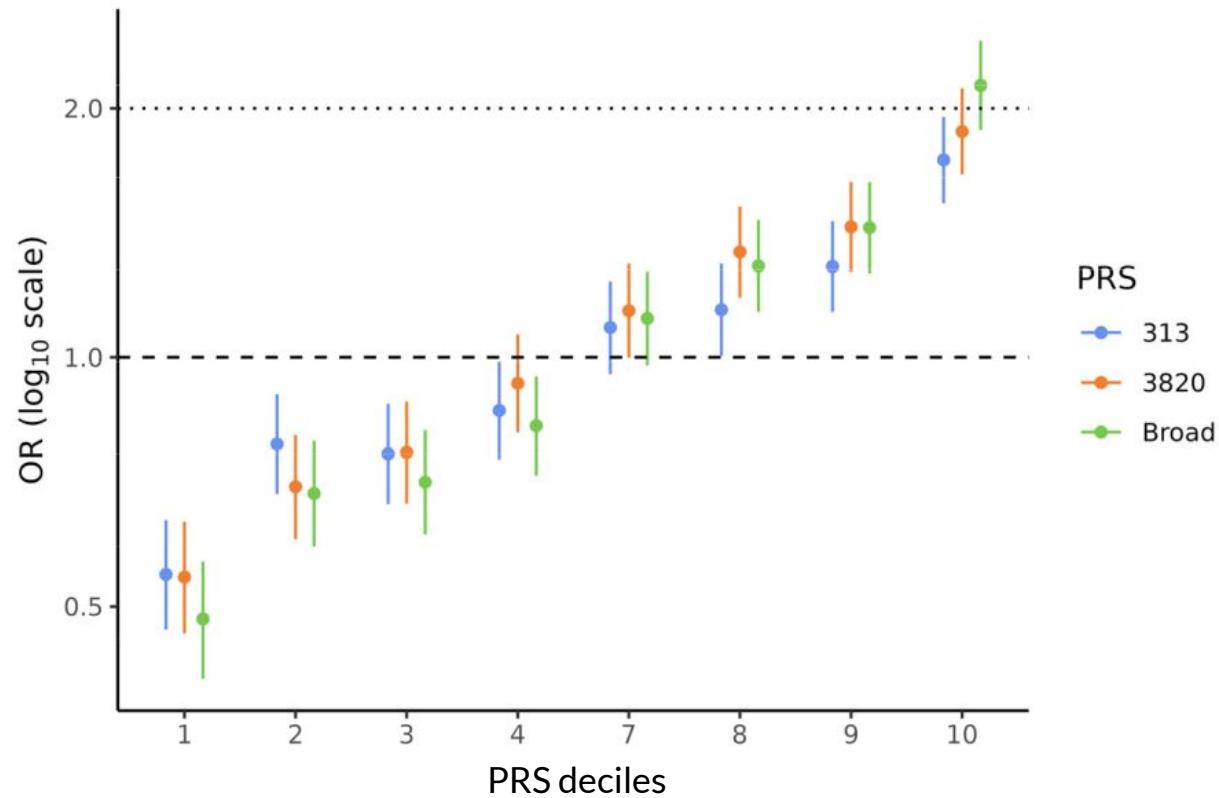
Repeat with other PRSs

Repeat by decile

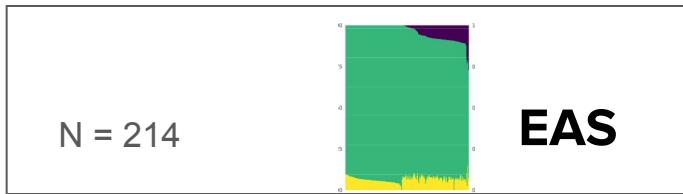
Repeat with other PRSs

Repeat by decile

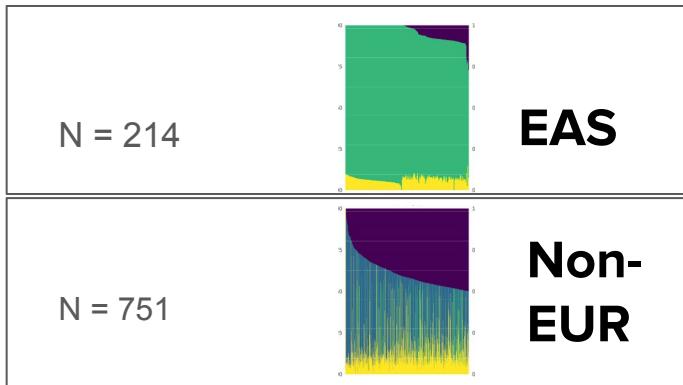
Repeat by ancestry composition



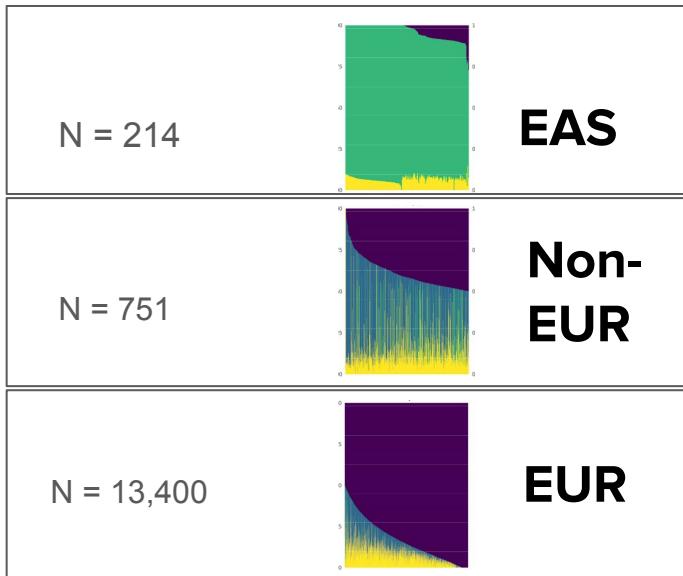
Per-Ancestry Results



Per-Ancestry Results

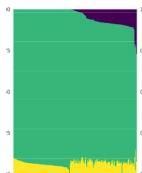


Per-Ancestry Results



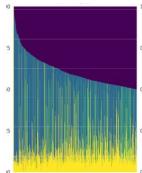
Per-Ancestry Results

N = 214



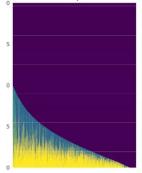
EAS

N = 751

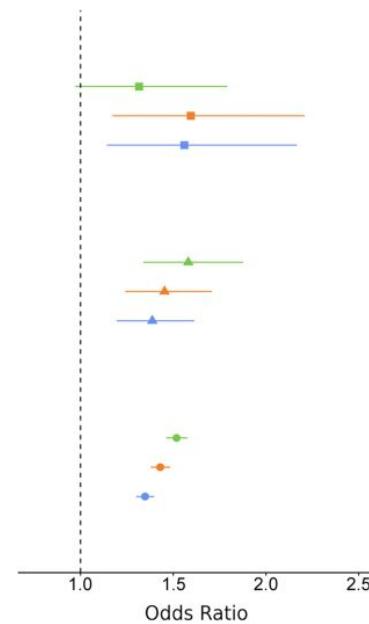


Non-
EUR

N = 13,400



EUR



PRS

● Broad ● 3820 ● 313

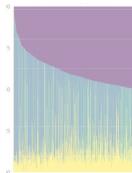
Per-Ancestry Results

N = 214



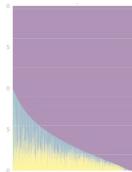
EAS

N = 751

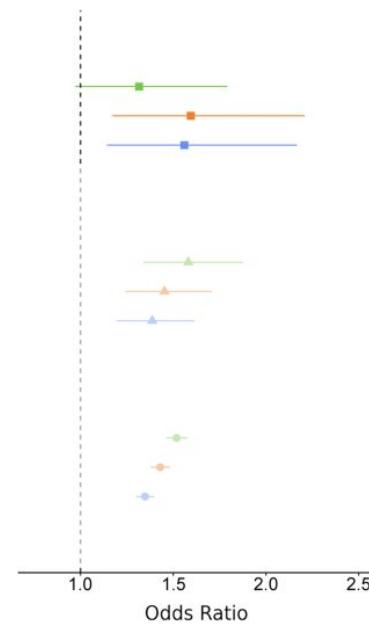


Non-
EUR

N = 13,400



EUR



PRS

● Broad ● 3820 ● 313

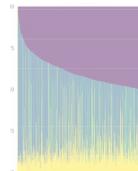
Per-Ancestry Results

N = 214



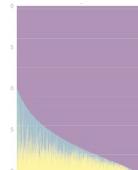
EAS

N = 751

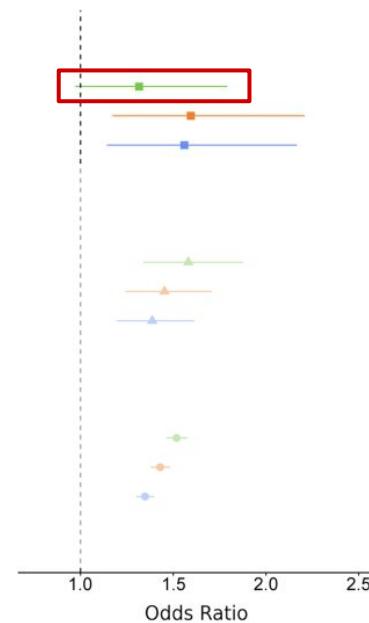


Non-
EUR

N = 13,400



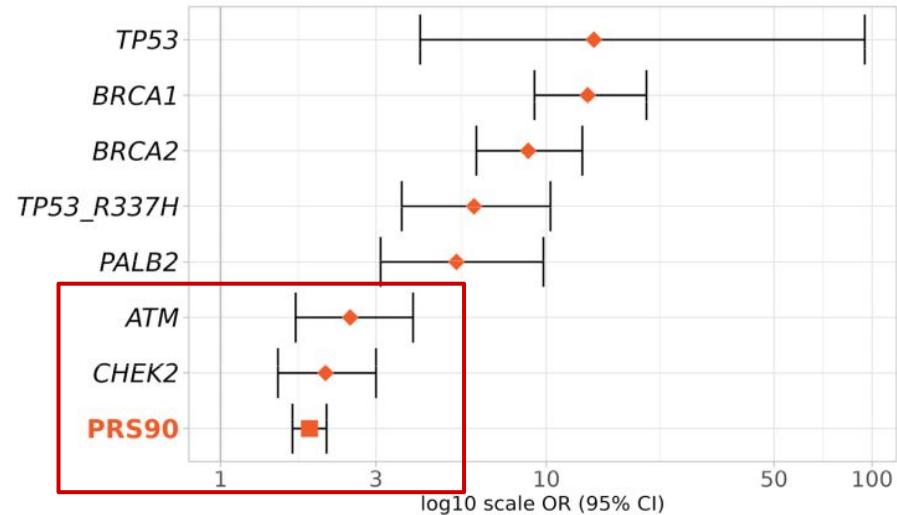
EUR



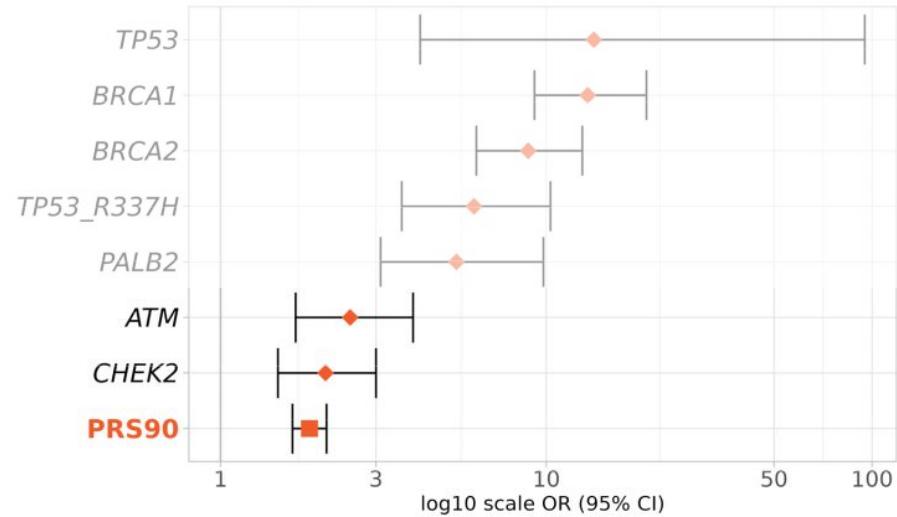
PRS

● B 1d ✗ 3820 313

PRS and Variants in Breast Cancer Genes



PRS and Variants in Breast Cancer Genes



Article

A Breast Cancer Polygenic Risk Score Validation in 15,490 Brazilians Using Exome Sequencing

Flávia Eichemberger Rius ^{1,2}, Rodrigo Santa Cruz Guindalini ³, Danilo Viana ¹, Júlia Salomão ¹, Laila Gallo ¹, Renata Freitas ¹, Cláudia Bertolacini ¹, Lucas Taniguti ¹ , Danilo Imparato ¹ , Flávia Antunes ¹, Gabriel Sousa ¹, Renan Achjian ¹, Eric Fukuyama ¹, Cleandra Gregório ¹ , Iuri Ventura ¹, Juliana Gomes ¹, Nathália Taniguti ¹ , Simone Maistro ² , José Eduardo Krieger ⁴ , Yonglan Zheng ⁵, Dezheng Huo ⁶, Olufunmilayo I. Olopade ⁵, Maria Aparecida Azevedo Koike Folgueira ² and David Schlesinger ^{1,*}

Gold Standard Hereditary Breast and Ovarian Cancer Panel



Mendelics

Key Takeaways

- Projects in data science are overwhelming 

Key Takeaways

- Projects in data science are overwhelming 
- {Tidymodels} orchestrates ML section fluidly 

Key Takeaways

- Projects in data science are overwhelming 
- {Tidymodels} orchestrates ML section fluidly 
- Enabled PRS validation in Brazil 

Thank you!



 [mendelics/prs_validation_posit_conf](https://github.com/mendelics/prs_validation_posit_conf)

Acknowledgements: Mendelics and ICESP-FMUSP team, paper collaborators, and all patients and controls who made this study possible.

Credits: Angelina Jolie Stock photos by Vecteezy.



@flaviaerius

flaviaerius.com