



A Platform for Generating and Validating Breast Risk Models from Clinical Data: Towards Patient-Centered Risk Stratified Screening

ucla mi²

Nova F. Smedley¹, Ngan Chau², Antonia Petrus², Alex A. T. Bui PhD¹, Arash Naeim MD PhD², William Hsu, PhD¹

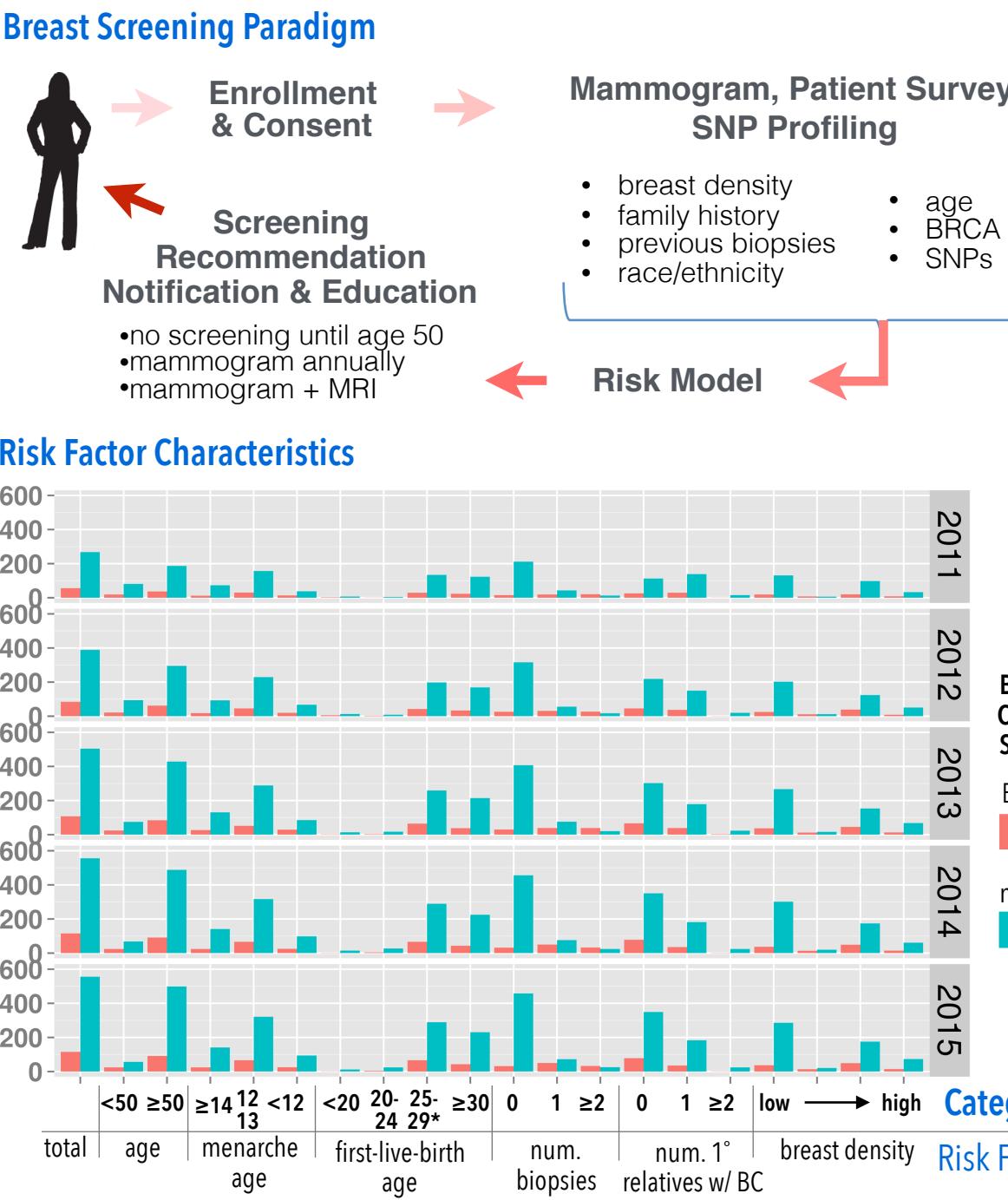
¹Medical Imaging Informatics, Departments of Bioengineering and Radiological Sciences, UCLA, Los Angeles, CA; ²Athena Breast Health Network, Jonsson Comprehensive Cancer Center, UCLA, Los Angeles, CA

Introduction

The purpose of this platform is to explore observational clinical data to stratify breast cancer screening population into risk groups by:

1. Characterizing risk factors
2. Evaluating classification performance of current risk models on UCLA population
3. Assessing the value of adding other risk factors (e.g., breast density)

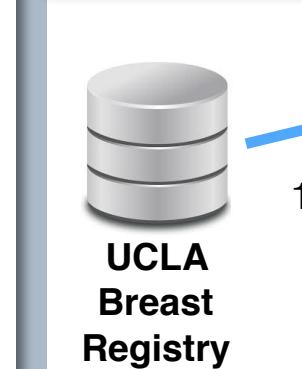
Data Collection



Data quality considerations

- Data cleaning was performed by cross-referencing sources.
- Changes in distribution over time can provide insight into potential systematic biases.
- Patients diagnosed with breast cancer are consistently underrepresented.

Methods



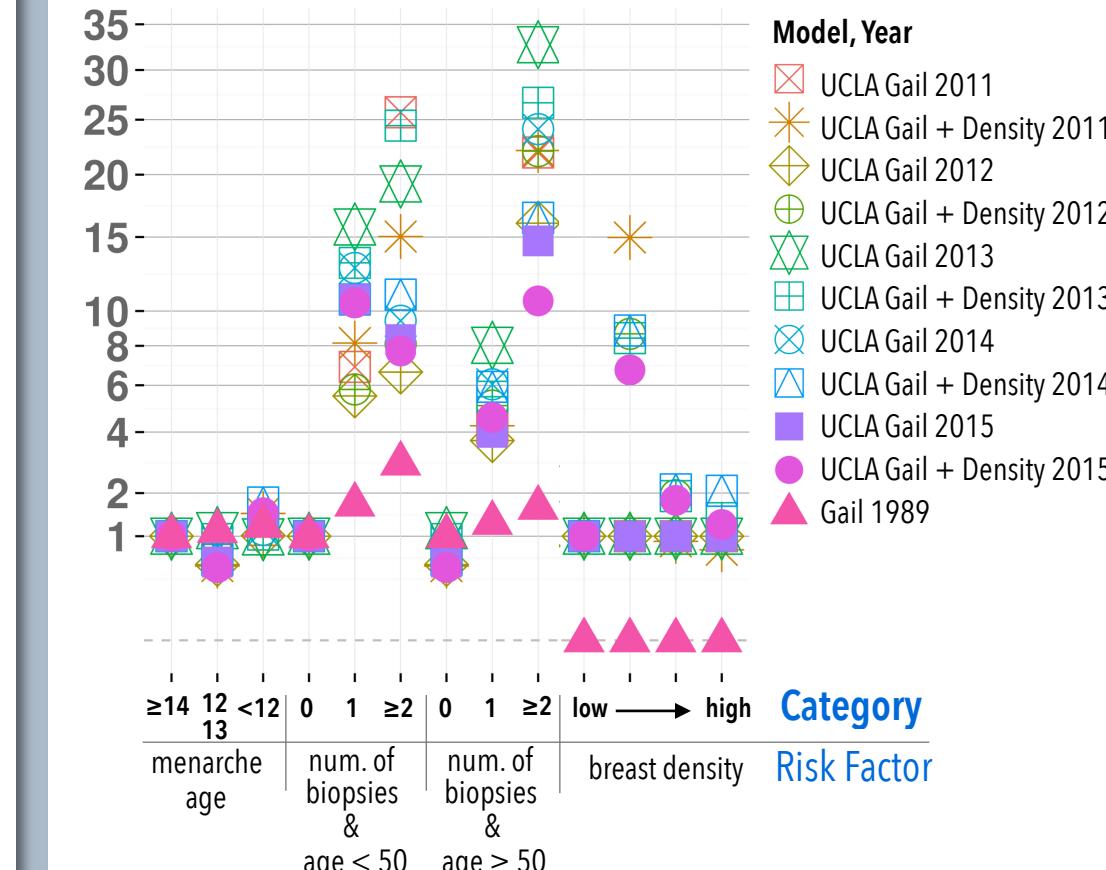
- 0.74948(intercept) → 80% training
+ 0.09401(menarche age) → 20% testing
+ 0.52916(num. biopsies)
+ 0.21863(first-live-birth age)
+ 0.95830(num. 1° relatives with breast cancer)
+ 0.01081(age category)
- 0.28804(num. biopsies)(age category)
- 0.19081(first-live-birth age)(num. 1° relatives with breast cancer)

estimate age-specific baseline hazard rate
where relative risk equals $e^{\text{coefficient}}$ from the logit model
+ age-specific breast cancer hazard rates estimated by (1)
+ age-specific competing hazard rates from death statistics estimated by (2)

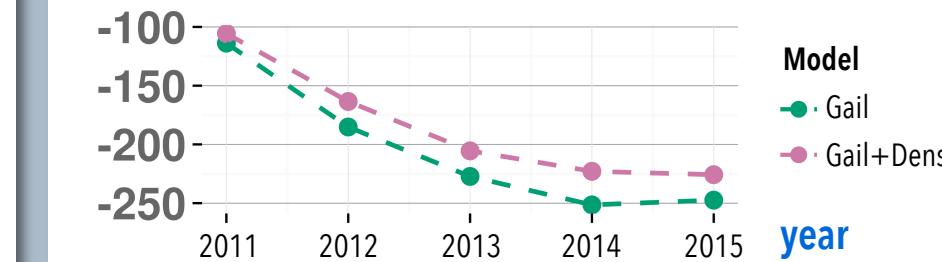
project relative risk to absolute risk of developing breast cancer is $P(a, \tau, r)$, where
 a = age of a patient
 τ = years to project
 r = relative risk associated with patient at a age

Results

Relative Risk



Log Likelihood



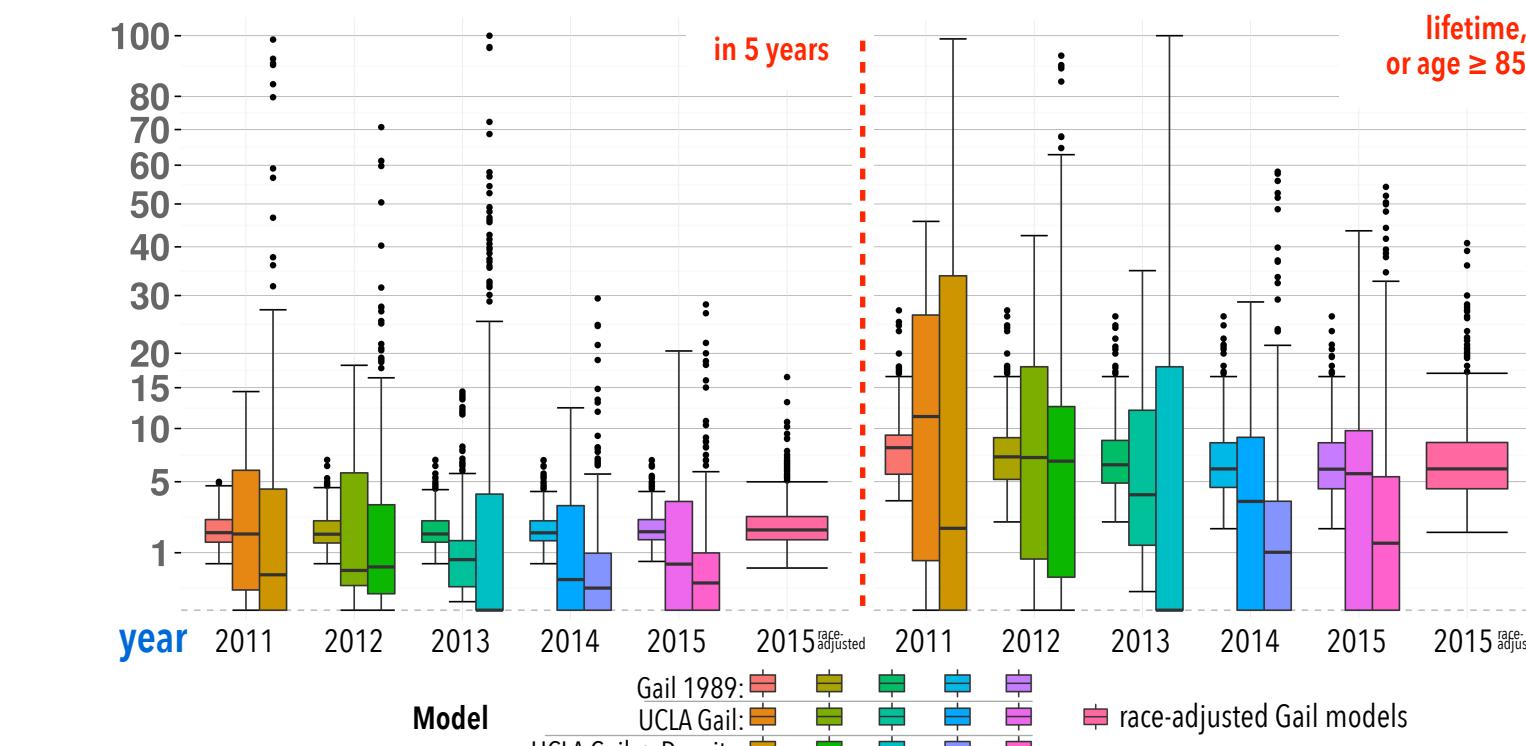
Likelihood of the data given the model:

- Likelihood decreases over time, possibly reflecting the poor fit of the model to UCLA's data.
- Including breast density improves the model fit.

References

1. Gail MH et al. Projecting individualized probabilities of developing breast cancer for white females who are being examined annually. J Natl Cancer Inst 1989; 81(24):1879-1886
2. Breast Cancer Risk Assessment Tool [Internet]. Bethesda: National Cancer Institute; [updated 2011 May 16]. Available from: www.cancer.gov/bcrisktool/Default.aspx

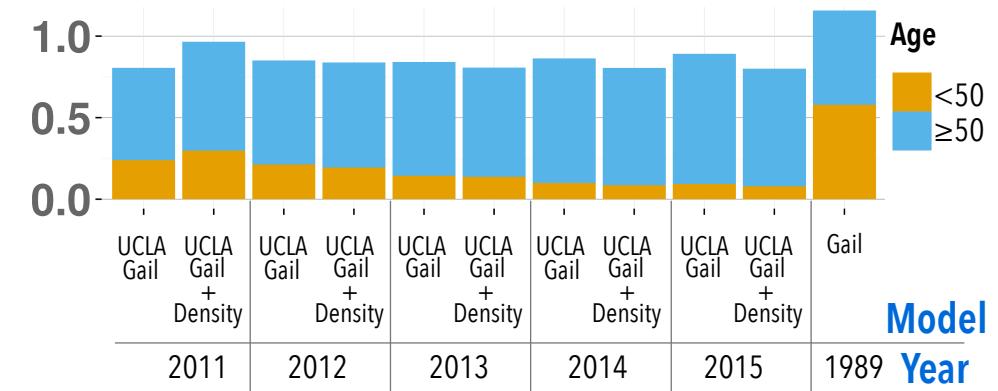
Probability of Cancer (%)



Interpreting relative risks:

- In the UCLA Gail 2013 model, a woman aged ≥ 50 with ≥ 2 breast biopsies has 33x more risk than a woman aged ≥ 50 with no breast biopsies.
- Slightly dense breast tissue in mammograms have the highest risk relative to other density categories.
- By 2015, the models are more similar to the Gail 1989 model.
- The addition of new risk factors requires larger sample sizes.
- An initial 80 SNPs for breast cancer risk was evaluated through principal component analysis, but showed uniform variance among components.

Population Attributable Factor



Population attribution factor (PAF) describes age category by its contribution to the relative risks associated with breast cancer.

- The Gail 1989 model has a balanced PAF
- At UCLA, women age ≥ 50 are adding the most risk information into the modeling cohort.

Summary of Conclusions

By building this platform, many of the challenges in noisy clinical data can be evaluated and problems can be tackled.

Information can be cataloged in real time — starting with smaller datasets as pilot studies, then iteratively improving data pipelines and workflows to accommodate decision making and research.

A learning health infrastructure enables an institution to generate and validate factors for risk stratification for their own study population.