

# Mammograms, Machine Learning and Cumulus

Daniel F. Schmidt, Enes Makalic, Jennifer Stone and John  
L. Hopper

Centre for MEGA Epidemiology  
The University of Melbourne

“Why Study Mammographic Density?”, 16–18th August,  
2010

# Outline

- 1 Introduction
- 2 Method
- 3 Preliminary Results

# Outline

- 1 Introduction
- 2 Method
- 3 Preliminary Results

# Background

- We are from computer science background
- Want to approach the problem from a machine learning perspective

# Background

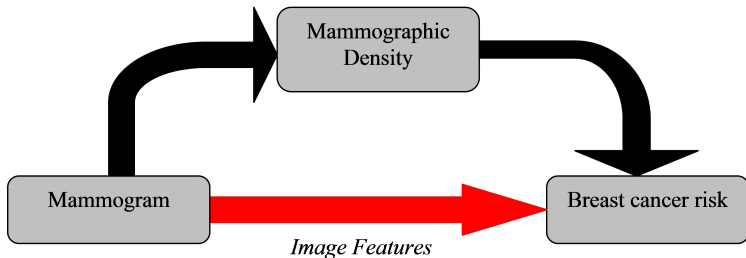
- We are from computer science background
- Want to approach the problem from a machine learning perspective

# Background

- We are from computer science background
- Want to approach the problem from a machine learning perspective

# Background

- We are from computer science background
- Want to approach the problem from a machine learning perspective



# CUMULUS

- One the most important tools for measuring MD is **CUMULUS**
- Transforms an image (mammogram) into a single summary statistic
- If  $\mathbf{Z}$  is a mammogram, then the **CUMULUS** measure is

$$\text{CUMULUS}(\mathbf{Z}) = \sum_{i,j} \mathbf{I}(Z_{i,j} > t)$$

where the threshold  $t$  is chosen manually

- Though simple, **CUMULUS** appears to be the “gold standard” for breast cancer risk  
⇒ Would be of great use if **CUMULUS** could be fully automated



# Outline

- 1 Introduction
- 2 Method**
- 3 Preliminary Results

# Training Dataset

- Aged matched film mammograms from Cambridge study
- Total people,  $n = 343$ ; multiple mammograms per person
  - 161 cases
  - 182 controls
- Restricted to MLOs only
  - ⇒ MLOs showed stronger association on this dataset
- Multiple mammograms per person handled by averaging
- For each mammogram, we had
  - Age (AGE)
  - Percent Dense Area (PDA)
  - Square-root Dense Area (DA)
  - Square-root Non Dense Area (NDA)

## Our Procedure (CIRRUSv1)

- Extract simple “features” from mammograms  
⇒ Total number of extracted features exceeded 2,500
- Use features in a linear logistic regression model

$$p(y_i = \text{case} | \mathbf{x}_i) = \frac{1}{1 + \exp(-\mathbf{x}_i \beta)}$$

where  $\mathbf{x}_i$  is the vector of features and  $y_i$  is the case-control status for individual  $i$  respectively

- Number of features is large; must use something more sophisticated than maximum likelihood/stepwise selection  
⇒ Bayesian regularised logistic regression employed

## Our Procedure (CIRRUSv1)

- Extract simple “features” from mammograms  
⇒ Total number of extracted features exceeded 2,500
- Use features in a linear logistic regression model

$$p(y_i = \text{case} | \mathbf{x}_i) = \frac{1}{1 + \exp(-\mathbf{x}_i \boldsymbol{\beta})}$$

where  $\mathbf{x}_i$  is the vector of features and  $y_i$  is the case-control status for individual  $i$  respectively

- Number of features is large; must use something more sophisticated than maximum likelihood/stepwise selection  
⇒ Bayesian regularised logistic regression employed

## Our Procedure (CIRRUSv1)

- Extract simple “features” from mammograms  
⇒ Total number of extracted features exceeded 2,500
- Use features in a linear logistic regression model

$$p(y_i = \text{case} | \mathbf{x}_i) = \frac{1}{1 + \exp(-\mathbf{x}_i \boldsymbol{\beta})}$$

where  $\mathbf{x}_i$  is the vector of features and  $y_i$  is the case-control status for individual  $i$  respectively

- Number of features is large; must use something more sophisticated than maximum likelihood/stepwise selection  
⇒ Bayesian regularised logistic regression employed

## “Smoothing” CUMULUS Measurements

- A CUMULUS measurement of MD is of the form

$$\text{CUMULUS} = \text{MD} + \text{bias} + \text{noise}$$

- MD is the true, underlying amount of MD (person dependent)
  - bias represents the limitations of the CUMULUS measurement scheme (person dependent)
  - noise is measurement noise due to various factors (measurement specific)
- Use PDA as “outcome”, our image features as “exposures”  
⇒ Regularised linear regression

## “Smoothing” CUMULUS Measurements

- A CUMULUS measurement of MD is of the form

$$\text{CUMULUS} = \text{MD} + \text{bias} + \text{noise}$$

- MD is the true, underlying amount of MD (person dependent)
- bias represents the limitations of the CUMULUS measurement scheme (person dependent)
- noise is measurement noise due to various factors (measurement specific)
- Use PDA as “outcome”, our image features as “exposures”  
⇒ Regularised linear regression

# Outline

- 1 Introduction
- 2 Method
- 3 Preliminary Results**



# Procedure

- Based on the Cambridge dataset
- Pectoral muscle was manually removed
- Compare **CIRRUSv1** against
  - **CUMULUS**: PDA adjusted for NDA and AGE
  - **SmoothedCUMULUS**: Smoothed PDA adjusted for AGE
- These are implemented in logistic regression models estimated using maximum likelihood
- Comparison metrics
  - Classification accuracy
  - Odds ratios
  - Correlation
- All metrics are computed *within-sample*

# Classification

- Classifying with logistic regression  
⇒ if  $p(y_i = \text{case} | \mathbf{x}_i) > 1/2$ , classify individual  $i$  as CASE
- Compare methods on number of correct classifications (classification accuracy)
- Can produce a “confusion matrix”; example:

	<i>CON</i>	<i>CAS</i>
<i>CON</i>	<i>TN</i>	<i>FP</i>
<i>CAS</i>	<i>FN</i>	<i>TP</i>

# Classification Results

- Baseline = 0.53
- **CUMULUS** = 0.6122,

	<i>CON</i>	<i>CAS</i>
<i>CON</i>	129	53
<i>CAS</i>	80	81

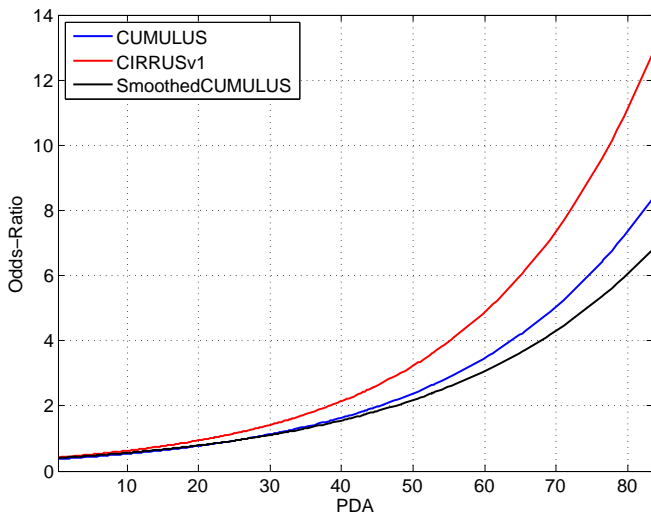
- **CIRRUSv1** = 0.6735,

	<i>CON</i>	<i>CAS</i>
<i>CON</i>	161	21
<i>CAS</i>	91	70

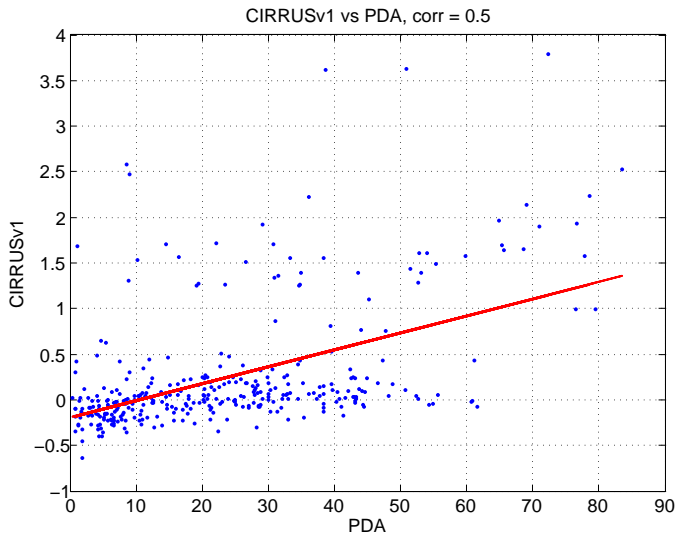
- **SmoothedCUMULUS** = 0.6501,

	<i>CON</i>	<i>CAS</i>
<i>CON</i>	136	46
<i>CAS</i>	74	87

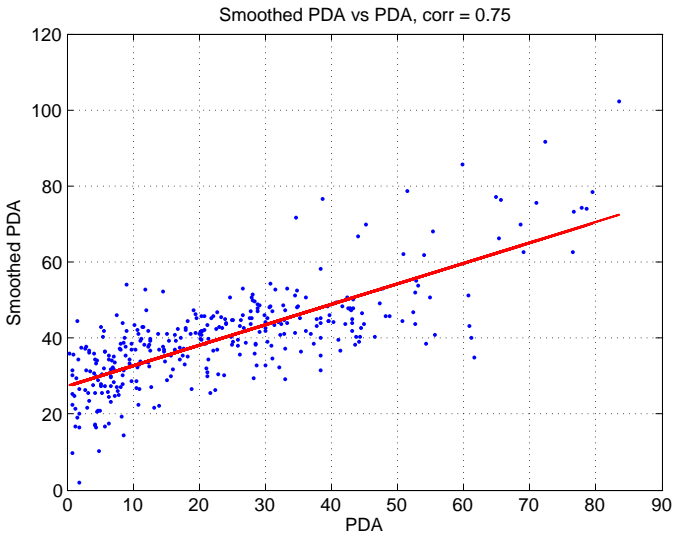
# Odds-Ratios



# CUMULUS vs CIRRUSv1



# CUMULUS vs SmoothedCUMULUS



## Future Work and Conclusion

- More ideas to try; more advanced features
- Try on larger datasets/validation datasets
- Questions?