

Knowledge Discovery from Structured Mammography Reports Using Inductive Logic Programming

Inês Dutra

ines@dcc.fc.up.pt

CRACS-INESC-Porto LA

DCC/FCUP

Joint work with

Elizabeth Burnside¹

Jesse Davis³

Chuck Kahn²

David Page¹

Vítor Santos Costa⁴

¹University of Wisconsin – Madison, WI, USA

²Medical College of Wisconsin – Milwaukee, WI, USA

³University of Leuven, Belgium

⁴University of Porto, Portugal

Objectives

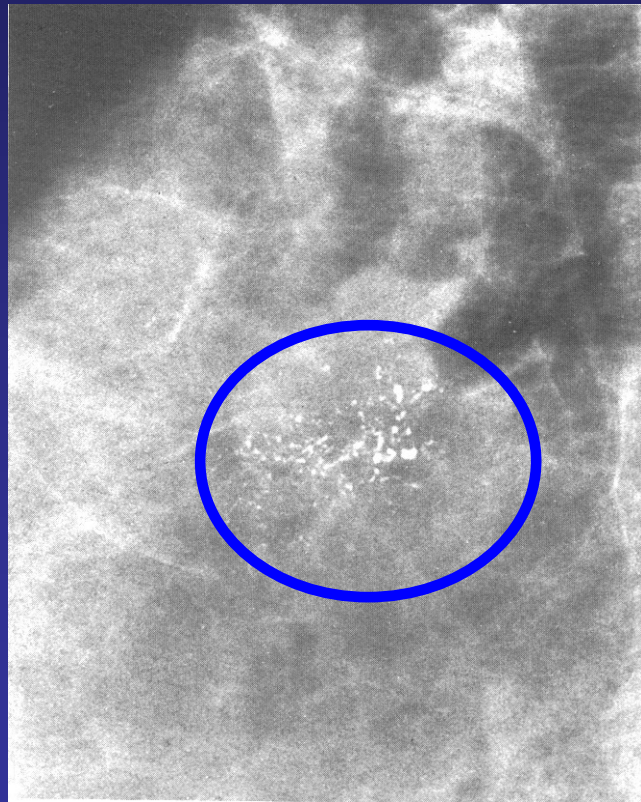
- Knowledge discovery using
 - Demographic information
 - Mammography findings
 - Inductive logic programming (ILP)
- Discuss domain of breast cancer imaging
- Describe Data
- ILP
- Unique features of data discovery in this domain
- Results
- Conclusions

Background

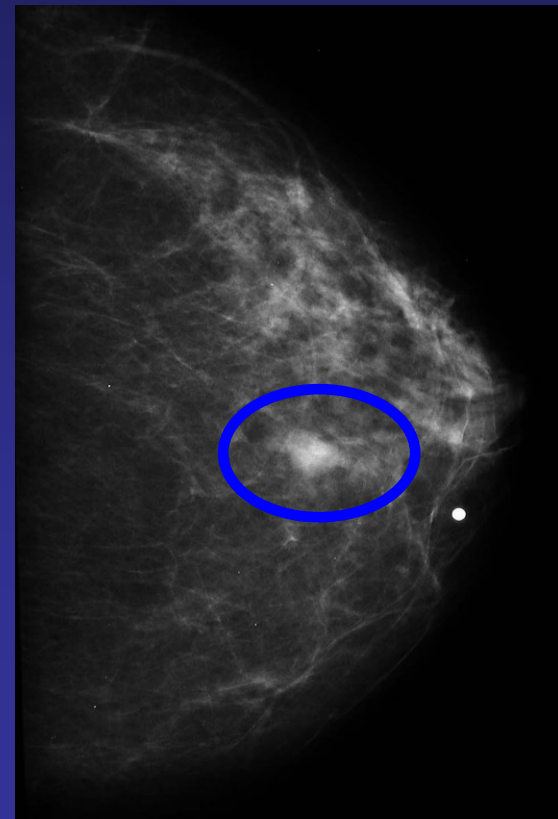
- Breast cancer is the most common cancer
- Mammography is the only proven screening test demonstrated in RCT to improve survival from breast cancer (also the cheapest)
- 20 million mammograms every 2 years

Common Mammogram Findings

Calcifications



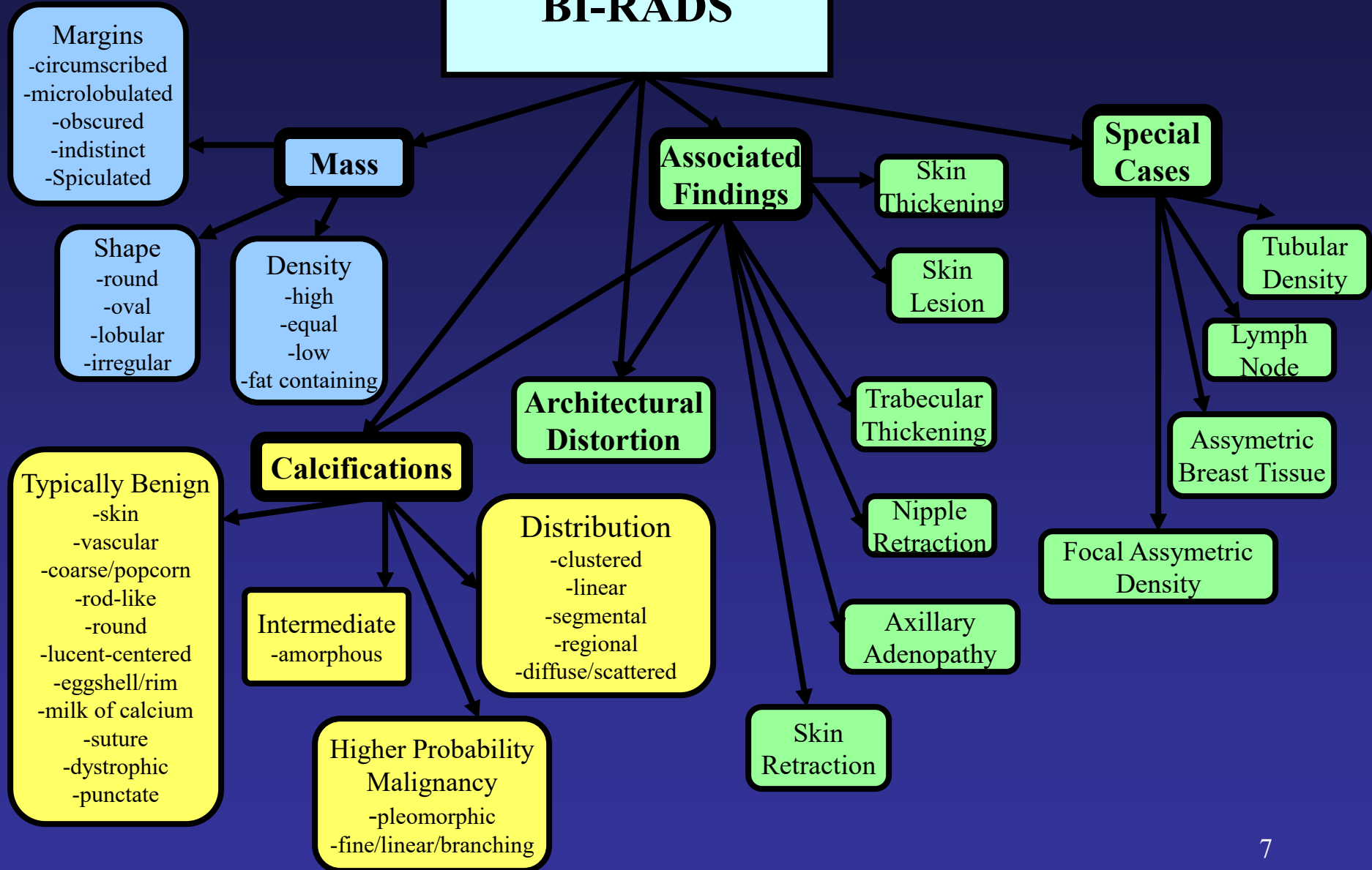
Masses



Mammography Lexicon

- BI-RADS
- A controlled terminology for mammography reporting
- A hierarchy of terms describing features of findings seen on mammograms
- Terms were selected among those most predictive of benign and malignant diseases

BI-RADS



BI-RADS Category

- Decisions are based on category:
 - BI-RADS 0: “Needs Additional Imaging”
 - BI-RADS 1: “Negative”
 - BI-RADS 2: “Benign”
 - BI-RADS 3: “Probably Benign”
 - BI-RADS 4: “Suspicious for malignancy”
 - BI-RADS 5: “Highly suggestive of malignancy”

Data

- National Mammography Database
- Defines a standard for reporting
 - Observed abnormalities on mammograms
 - Structured data that facilitates the use of computer technologies
- Our dataset contains
 - All abnormalities from 1999-2004 at MCW
 - 435 malignancies
 - 65,365 benign abnormalities

Data

Patient information	Abnormality location	Mass descriptors	Calcification descriptors
Age	Side	Shape	Shape
Hormone therapy	Depth	Density	Distribution
Family medical history	Clock location	Margins	Stability
Personal medical history	Quadrant location	Stability	

Mammography Database

Patient	Finding	Date	Calcification Fine/Linear	...	Mass Size	Loc	Benign/ Malignant
P1	1	5/02	No		0.03	RU4	B
P1	2	5/04	Yes		0.05	RU4	M
P1	3	5/04	No		0.04	LL3	B
P2	4	6/00	No		0.02	RL2	B
...

Important Change Over Time

Patient	Finding	Date	Calcification Fine/Linear	...	Mass Size	Loc	Benign/ Malignant
P1	1	5/02	No		0.03	RU4	B
P1	2	5/04	Yes		0.05	RU4	M
P1	3	5/04	No		0.04	LL3	B
P2	4	6/00	No		0.02	RL2	B
...

Inductive Logic Programming (ILP)

- Learn set of rules in 1st order logic
- Rules distinguish between positive and negative examples
- ILP algorithms: Aleph (Srinivasan), Progol (Muggleton), FOIL (Quinlan), etc.

Inductive Logic Programming (ILP)

- Assumption 1
 - Background knowledge B
 - Form of a Prolog program

Inductive Logic Programming (ILP)

- Assumption 2
 - Language specification
 - Clause representation
- ```
is_malignant(A) :-
 'Age'(A,age6570),
 'MassesShape'(A,spiculated),
 'BIRADS_category'(A,b5)
```

# Inductive Logic Programming (ILP)

---

- Assumption 3
  - Constraints on acceptable clauses
  - Example:  
clause cannot have more than six literals



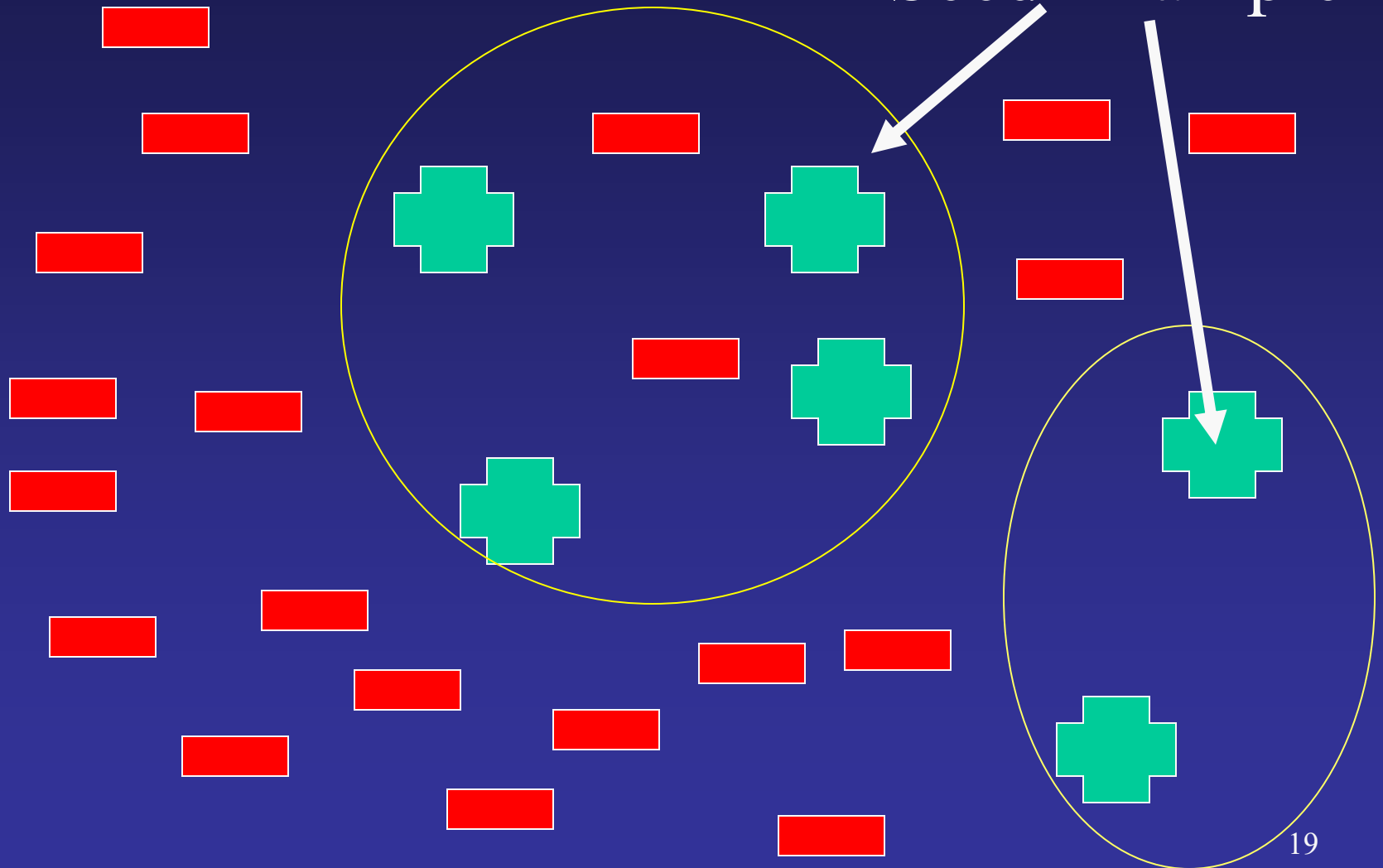
# Inductive Logic Programming (ILP)

---

- Assumption 4
  - Finite set of examples
    - $E^+$  are positive examples (malignant)
    - $E^-$  are negative examples (benign)

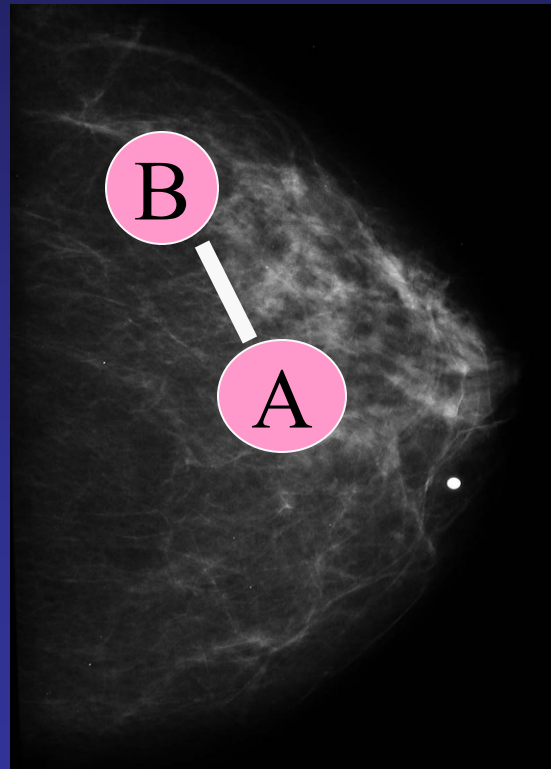
# ILP as in Aleph

Seed Example



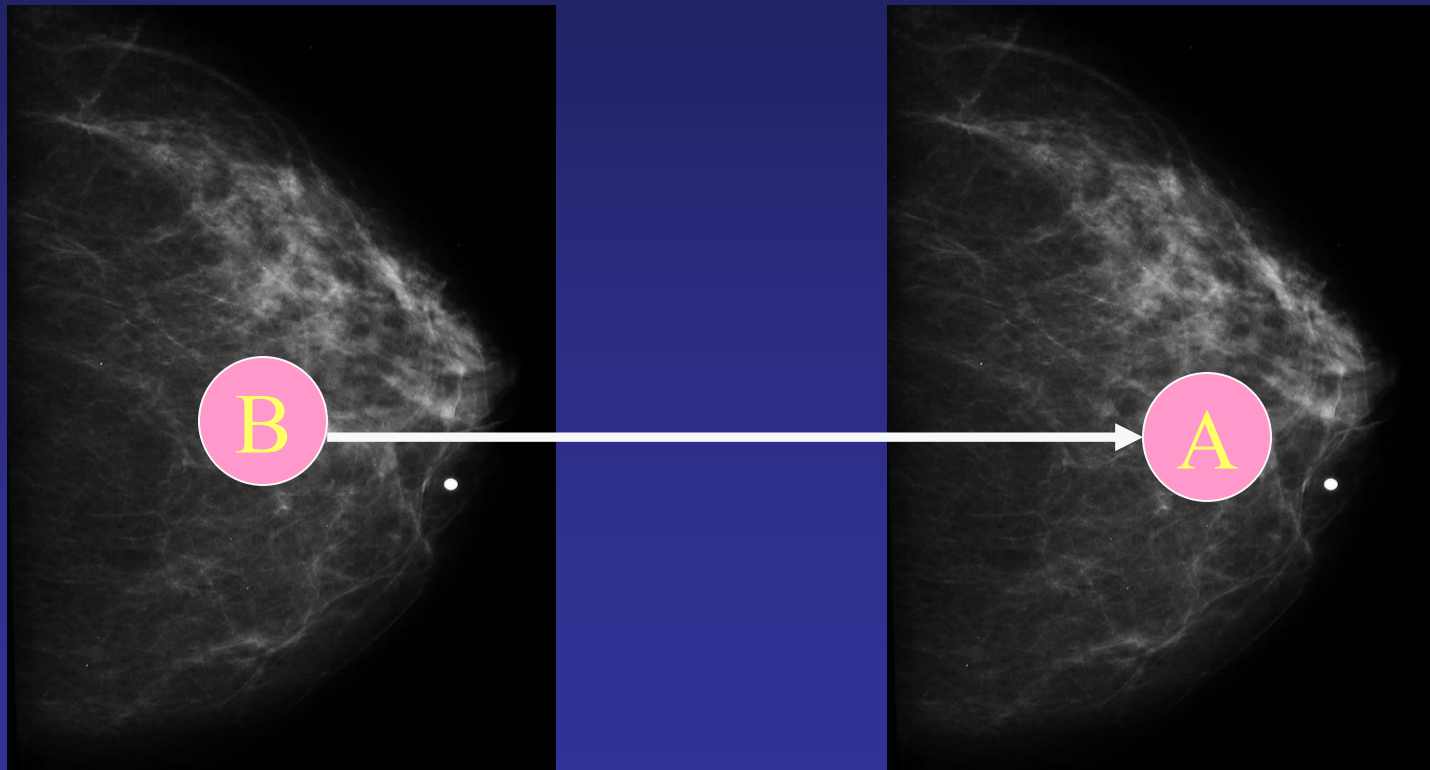
# Incorporating related information

`in_same_mammogram(A,B)`



# Incorporating related information

prior\_mammogram(A,B)



# Results

---

- Aleph discovered millions of rules
- Evaluated each rule by how well it covered E+ and not E-
- We quantify performance using the m-estimate, a smoothed ratio between the number of positive E+ covered and the total E covered
- Selected the top 130 rules

# Review of Results

---

- Radiologist reviewed these rules
- Found 2 to be interesting
- Significance of these rules was validated using the data

# Rule #1

---

is\_malignant(A) :-

'BIRADS\_category'(A,b5),

'MassPAO'(A,present),

'MassesDensity'(A,high),

'HO\_BreastCA'(A,hxDCorLC),

in\_same\_mammogram(A,B),

'Calc\_Pleomorphic'(B,notPresent),

'Calc\_Punctate'(B,notPresent).

# Rule #1

---

is\_malignant(A) :-

'BIRADS\_category'(A,b5),

'MassPAO'(A,present),

'MassesDensity'(A,high),

'HO\_BreastCA'(A,hxDCorLC),

in\_same\_mammogram(A,B),

'Calc\_Pleomorphic'(B,notPresent),

'Calc\_Punctate'(B,notPresent).

42 malignant and 11 benign findings



# Significance of Rule #1

---

Mass density has not previously been significantly associated with breast cancer

*Jackson VP, Dines KA, Bassett LW, Gold RH, Reynolds HE, Diagnostic importance of the radiographic density of noncalcified breast masses: analysis of 91 lesions. Am J Roentgenol. 1991 Jul;157(1):25-8.*

# Validity of Result #1

---

| Mass Density | Malignant (%) |        | Total |
|--------------|---------------|--------|-------|
| Fat-density  | 0             | (0)    | 493   |
| Low          | 2             | (.1)   | 3408  |
| Equal        | 17            | (3.3)  | 513   |
| High         | 103           | (31.7) | 324   |
| Total        | 122           | (2.6)  | 4738  |

# Validity of Result #1

---

| Mass Density | Malignant (%) |        | Total |
|--------------|---------------|--------|-------|
| Fat-density  | 0             | (0)    | 493   |
| Low          | 2             | (.1)   | 3408  |
| Equal        | 17            | (3.3)  | 513   |
| High         | 103           | (31.7) | 324   |
| Total        | 122           | (2.6)  | 4738  |



# Validity of Result #1

---

| Mass Density | Malignant (%) |        | Total |
|--------------|---------------|--------|-------|
| Fat-density  | 0             | (0)    | 493   |
| Low          | 2             | (.1)   | 3408  |
| Equal        | 17            | (3.3)  | 513   |
| High         | 103           | (31.7) | 324   |
| Total        | 122           | (2.6)  | 4738  |



P < .001

# Validity of Result #1

---

| Mass Density | Malignant (%) |        | Total |
|--------------|---------------|--------|-------|
| Fat-density  | 0             | (0)    | 493   |
| Low          | 2             | (.1)   | 3408  |
| Equal        | 17            | (3.3)  | 513   |
| High         | 103           | (31.7) | 324   |
| Total        | 122           | (2.6)  | 4738  |



P < .001

# Validity of Result #1

| Mass Density | Malignant (%) |        | Total |
|--------------|---------------|--------|-------|
| Fat-density  | 0             | (0)    | 493   |
| Low          | 2             | (.1)   | 3408  |
| Equal        | 17            | (3.3)  | 513   |
| High         | 103           | (31.7) | 324   |
| Total        | 122           | (2.6)  | 4738  |



P < .001

# Rule #2

---

```
is_malignant(A) :-
 'BIRADS_category'(A,b5),
 'Mass'(A,present),
 'Age'(A,age6570),
 previous_finding(A,B),
 'Calc_Punctate'(B,notPresent),
 'BIRADS_category'(B,b3).
```

# Rule #2

---

is\_malignant(A) :-

'BIRADS\_category'(A,b5),

'Mass'(A,present),

'Age'(A,age6570),

previous\_finding(A,B),

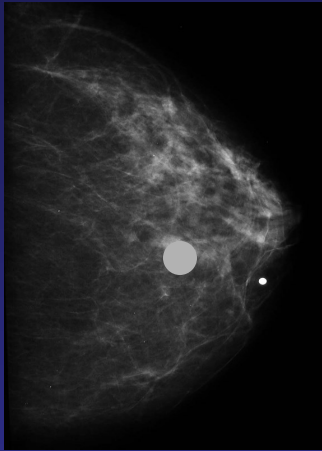
'Calc\_Punctate'(B,notPresent),

'BIRADS\_category'(B,b3).

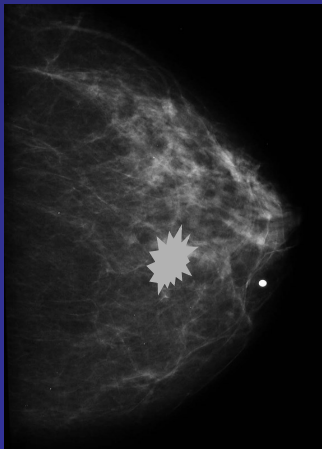


# Significance of Rule #2

---



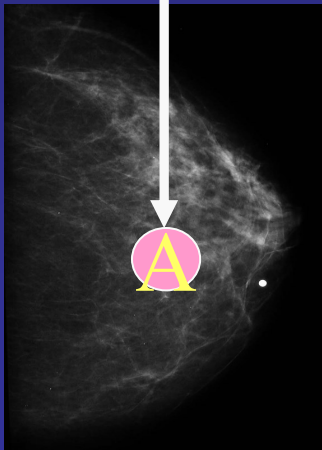
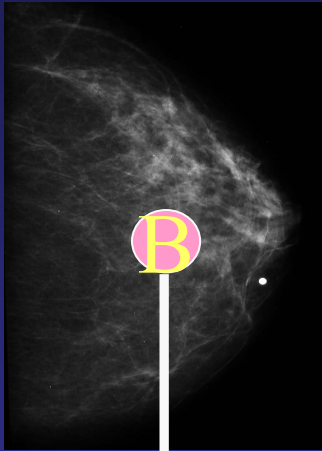
BI-RADS 3:  
Probably Benign



BI-RADS 5: Highly  
Suggestive of  
Malignancy

# Significance of Rule #2

---



B

A

BI-RADS 3:  
Probably Benign



Delay in  
Diagnosis!!

BI-RADS 5: Highly  
Suggestive of  
Malignancy

# Validity of Result #2

Analysing all cases labeled as BI-RADS 3 and later diagnosed with cancer

| BI-RADS 3 abnormality    |      |       |       |      | BI-RADS 5 abnormality        |      |       |       |      |            |
|--------------------------|------|-------|-------|------|------------------------------|------|-------|-------|------|------------|
| abnormality              | side | clock | depth | quad | abnormality                  | side | clock | depth | quad | match      |
| Clustered calcifications | L    | 12    | M     | UO   | High density spiculated mass | L    | C     | M     | *    | possible   |
| Ill-defined oval mass    | R    | 11    | M     | UO   | High density spiculated mass | R    | 11    | M     | UO   | <b>yes</b> |
| Oval circumscribed mass  | R    | 12    | A     | UI   | Oval spiculated mass         | R    | 5     | P     | UI   | no         |
| *                        | R    | 4     | M     | *    | Round spiculated mass        | R    | 4     | M     | LI   | <b>yes</b> |
| Oval mass                | R    | 12    | P     | UO   | Irregular spiculated mass    | R    | 12    | P     | UO   | <b>yes</b> |
| Ill-defined oval mass    | R    | 2     | P     | LI   | Irregular high density mass  | R    | 2     | P     | LI   | <b>yes</b> |
| *                        | L    | 12    | M     | UO   | Irregular spiculated mass    | L    | 1     | M     | UO   | possible   |

**Important: location!**

# Conclusion

---

- With large amounts of data, ILP holds significant promise in the domain of mammography to discover novel hypothesis and provide quality assurance.

Thank you!

# Connecting Abnormalities

May  
2002

Patient 1

May  
2004

