# Base Rate Fallacy

Prof. Ravi Sandhu
Executive Director and Endowed Chair

Lecture 13

ravi.utsa@gmail.com
www.profsandhu.com

*World-Leading Research with Real-World Impact!*

S: Patient is **S**ick
(has the disease)

| | S | ¬S |
|---|---|---|
| **R** | R ∧ S<br><br>True positive | R ∧ ¬S<br><br>False positive |
| **¬R** | ¬R ∧ S<br><br>False negative | ¬R ∧ ¬S<br><br>True negative |

R: Test **R**esult
is positive

*World-Leading Research with Real-World Impact!*

# Base-Rate Fallacy

S: Patient is **S**ick
(has the disease)
System is under attack

|  | S | ¬S |
|---|---|---|
| **R** | R ∧ S<br><br>True positive | R ∧ ¬S<br><br>False positive |
| **¬R** | ¬R ∧ S<br><br>False negative | ¬R ∧ ¬S<br><br>True negative |

R: Test **R**esult
is positive
Alarm is raised

*World-Leading Research with Real-World Impact!*

# Malware Detection Techniques

I know what is bad and can detect it
False positives: none
False negatives: ever increasing

I will learn what is good and bad
False positives: incorrect learning
False negatives: incorrect learning

**Malware Detection**

**Signature-based**

**Anomaly-based**

static    dynamic    hybrid

static    dynamic    hybrid

**Specification-based**

static    dynamic    hybrid

I know what is good and can detect when you go beyond specification
False positives: incomplete specification
False negatives: incorrect specification

Nwokedi Idika and Aditya Mathur, A Survey of Malware Detection Techniques, Purdue University, Feb 2007.

S: Patient is **S**ick
(has the disease)

|  | S | ¬S |
|---|---|---|
| **R** | **R ∧ S**<br><br>True positive | **R ∧ ¬S**<br><br>False positive |
| **¬R** | **¬R ∧ S**<br><br>False negative | **¬R ∧ ¬S**<br><br>True negative |

R: Test **R**esult
is positive

*World-Leading Research with Real-World Impact!*

# Base-Rate Fallacy

S: Patient is **S**ick
(has the disease)

|  | S | ¬S |
|---|---|---|
| **R** | R ∧ S<br><br>True positive<br><br>P(R\|S) = 0.99 | R ∧ ¬S<br><br>False positive<br><br>P(R\|¬S) = 0.01 |
| **¬R** | ¬R ∧ S<br><br>False negative<br><br>P(¬R\|S) = 0.01 | ¬R ∧ ¬S<br><br>True negative<br><br>P(¬R\|¬S) = 0.99 |

R: Test **R**esult is positive

These probabilities can be empirically estimated

*World-Leading Research with Real-World Impact!*

2000 sick

1000 not sick

Test R
is positive

Test R
is negative

Test R
is positive

Test R
is negative

1980

20

10

990

estimate    $P(R|S) = 0.99$    $P(\neg R|S) = 0.01$    $P(R|\neg S) = 0.01$    $P(\neg R|\neg S) = 0.99$

Coincidentally equal

# Estimating P(R|S) etc

UTSA

2000 sick                                   1000 not sick

Test R              Test R              Test R              Test R
is positive         is negative         is positive         is negative

1980                20                  30                  970

estimate    P(R|S) = 0.99    P(¬R|S) = 0.01    P(R|¬S) = 0.03    P(¬R|¬S) = 0.97

In general will not be equal

# Base-Rate Fallacy

S: Patient is **S**ick
(has the disease)

|  | S | ¬S |
|---|---|---|
| **R** | R ∧ S<br><br>True positive<br><br>P(R\|S) = 0.99 | R ∧ ¬S<br><br>False positive<br><br>P(R\|¬S) = 0.03 |
| **¬R** | ¬R ∧ S<br><br>False negative<br><br>P(¬R\|S) = 0.01 | ¬R ∧ ¬S<br><br>True negative<br><br>P(¬R\|¬S) = 0.97 |

R: Test **R**esult
is positive

Rows must
total between
0 and 2

These probabilities
can be empirically
estimated

Columns must total 1

*World-Leading Research with Real-World Impact!*

# Base-Rate Fallacy

S: Patient is **S**ick
(has the disease)

We will continue
with these numbers

|  | S | ¬S |
|---|---|---|
| **R** | R ∧ S<br><br>True positive<br><br>P(R\|S) = 0.99 | R ∧ ¬S<br><br>False positive<br><br>P(R\|¬S) = 0.01 |
| **¬R** | ¬R ∧ S<br><br>False negative<br><br>P(¬R\|S) = 0.01 | ¬R ∧ ¬S<br><br>True negative<br><br>P(¬R\|¬S) = 0.99 |

R: Test **R**esult
is positive

These probabilities
can be empirically
estimated

# Real Interest

S: Patient is **S**ick
(has the disease)

|  | S | ¬S |
|---|---|---|
| **R** | R ∧ S<br><br>True positive<br><br>P(S\|R) = ?? | R ∧ ¬S<br><br>False positive<br><br>P(¬S\|R) = ?? |
| **¬R** | ¬R ∧ S<br><br>False negative<br><br>P(S\|¬R) = ?? | ¬R ∧ ¬S<br><br>True negative<br><br>P(¬S\|¬R) = ?? |

R: Test **R**esult
is positive

Rows must total 1

These probabilities can be computed by Bayes' theorem if we know P(S)

Columns must total between 0 and 2

*World-Leading Research with Real-World Impact!*

➢ P(S|R) =
(P(S) × P(R|S))/
(P(S) × P(R|S)+P(¬S) ) × P(R|¬S))

➢ P(¬S|R) = 1 - P(S|R)

➢ P(S|¬R) =
(P(S) × P(¬R|S))/
(P(S) × P(¬R|S)+P(¬S) ) × P(¬R|¬S))

➢ P(¬S|¬R) = 1 - P(S|¬R)

# Base-Rate Fallacy

S: Patient is **S**ick
(has the disease)

We will continue
with these numbers

|  | S | ¬S |
|---|---|---|
| **R** | R ∧ S<br><br>True positive<br><br>P(R\|S) = 0.99 | R ∧ ¬S<br><br>False positive<br><br>P(R\|¬S) = 0.01 |
| **¬R** | ¬R ∧ S<br><br>False negative<br><br>P(¬R\|S) = 0.01 | ¬R ∧ ¬S<br><br>True negative<br><br>P(¬R\|¬S) = 0.99 |

R: Test **R**esult
is positive

These probabilities
can be empirically
estimated

# Real Interest

**Assume**
**P(S)=0.0001**
**1 in 10,000 has**
**disease**

S: Patient is **S**ick
(has the disease)

|  | S | ¬S |
|---|---|---|
| **R** | **R ∧ S**<br><br>True positive<br><br>$P(S\|R) = 0.009804$ | **R ∧ ¬S**<br><br>False positive<br><br>$P(\neg S\|R) = 0.990196$ |
| **¬R** | **¬R ∧ S**<br><br>False negative<br><br>$P(S\|\neg R) = 0.000001$ | **¬R ∧ ¬S**<br><br>True negative<br><br>$P(\neg S\|\neg R) = 0.999999$ |

R: Test **R**esult
is positive

**Rows must**
**total 1**

**These probabilities**
**can be computed by**
**Bayes' theorem if we**
**know P(S)**

**Columns must total between 0 and 2**

# False Alarms Predominate!

Assume
P(S)=0.0001
1 in 10,000 has
disease

| P(S\|R) | requires | P(R\|¬S) |
|---------|----------|----------|
| 0.01    |          | 0.01     |
| 0.09    |          | 0.001    |
| 0.5     |          | 0.0001   |
| 0.9     |          | 0.00001  |
| 0.99    |          | 0.000001 |

# Base-Rate Fallacy

S: Patient is **S**ick
(has the disease)

Total population = 1,000,000
1 in 10,000 has disease

|  | S | ¬S |
|---|---|---|
|  | 100 | 999,900 |
| **R** | R ∧ S<br><br>True positive | R ∧ ¬S<br><br>False positive |
| **¬R** | ¬R ∧ S<br><br>False negative | ¬R ∧ ¬S<br><br>True negative |

R: Test **R**esult
is positive

R is 99% accurate
for sick and non-sick
populations

# Base-Rate Fallacy

S: Patient is **S**ick
(has the disease)

Total population = 1,000,000
1 in 10,000 has disease

|  | S | ¬S |
|---|---|---|
|  | 100 | 999,900 |
| **R** | R ∧ S<br><br>True positive<br><br>99 | R ∧ ¬S<br><br>False positive<br><br>9,999 |
| **¬R** | ¬R ∧ S<br><br>False negative<br><br>1 | ¬R ∧ ¬S<br><br>True negative<br><br>989,901 |

R: Test **R**esult
is positive

R is 99% accurate
for sick and non-sick
populations