

# Module 5.1 Base Rate Fallacy

Ravi Sandhu

Spring 2021

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

S

$\neg$ S

	R	$R \wedge S$ True positive	$R \wedge \neg S$ False positive
	$\neg$ R	$\neg R \wedge S$ False negative	$\neg R \wedge \neg S$ True negative

R: Test **R**esult  
is positive

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

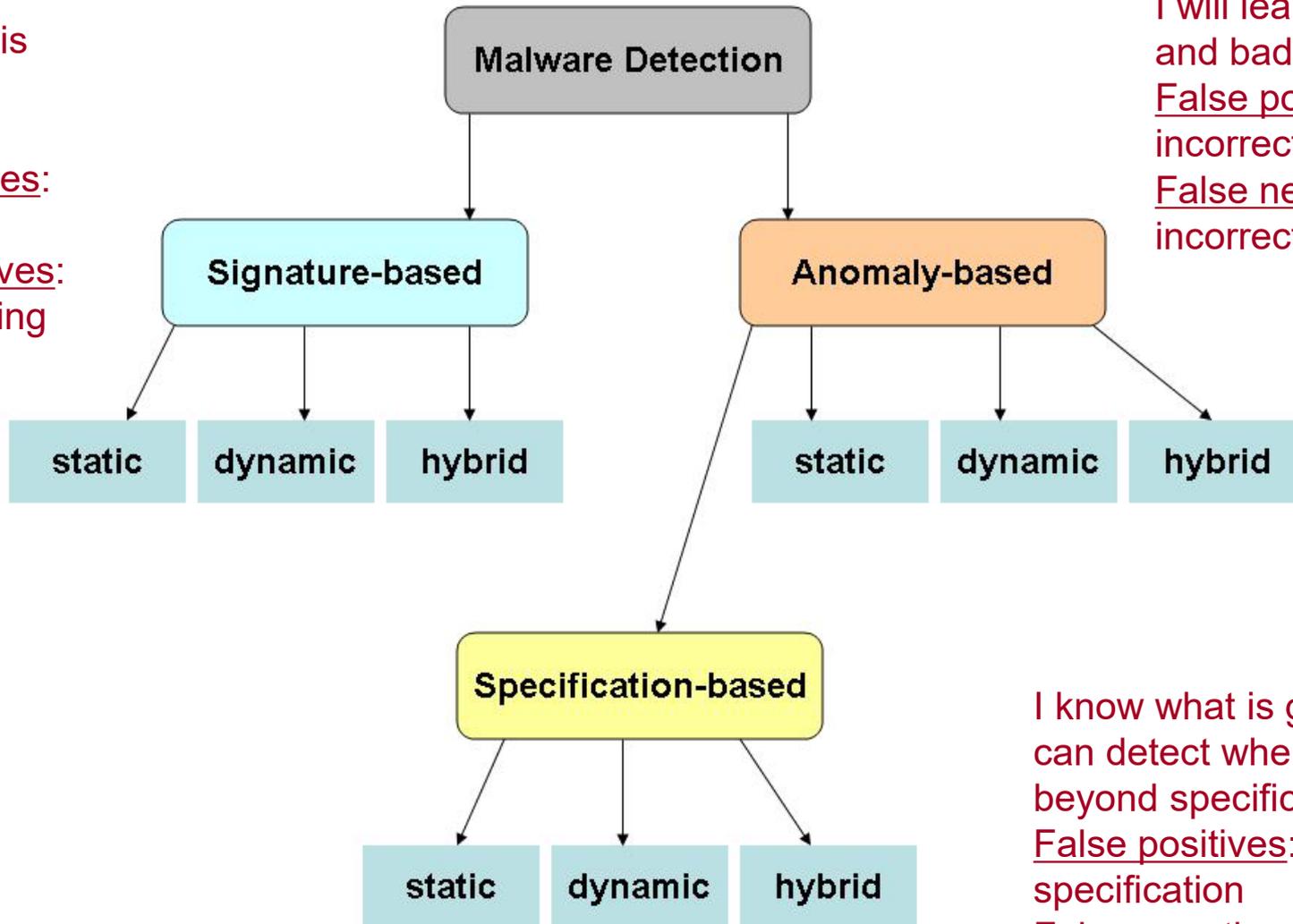
S

$\neg$ S

	$R \wedge S$	$R \wedge \neg S$
R	True positive	False positive
	$\neg R \wedge S$	$\neg R \wedge \neg S$
$\neg$ R	False negative	True negative

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

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

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

$\neg$ S

	R	$R \wedge S$ True positive	$R \wedge \neg S$ False positive
	$\neg$ R	$\neg R \wedge S$ False negative	$\neg R \wedge \neg S$ True negative

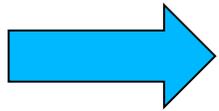
R: Test **R**esult  
is positive

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

S

$\neg$ S

	R $\wedge$ S	R $\wedge$ $\neg$ S
R	True positive $P(R S) = 0.99$	False positive $P(R \neg S) = 0.01$
	$\neg$ R $\wedge$ S	$\neg$ R $\wedge$ $\neg$ S
$\neg$ R	False negative $P(\neg R S) = 0.01$	True negative $P(\neg R \neg S) = 0.99$

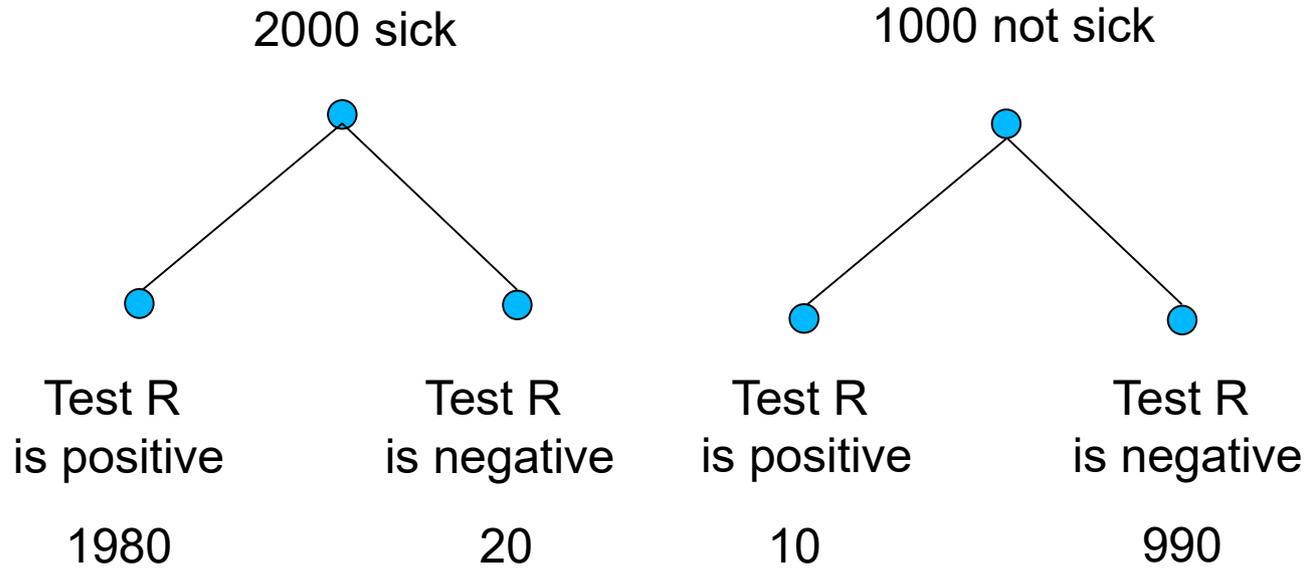


R: Test **R**esult  
is positive

Rows must  
total between  
0 and 2

These probabilities  
can be empirically  
estimated

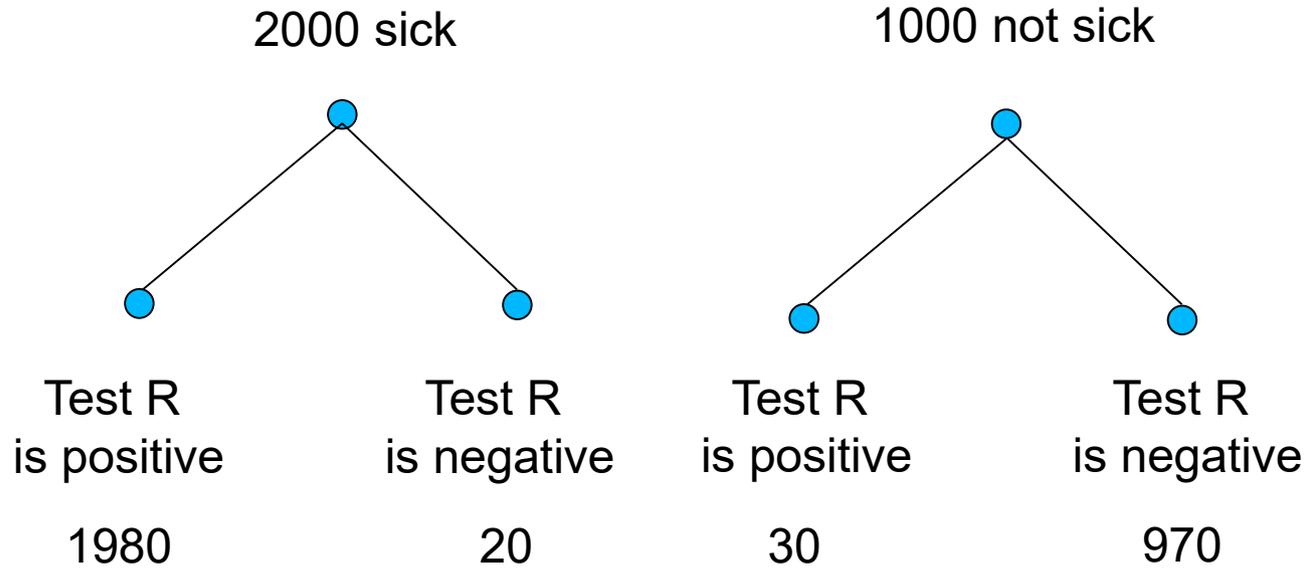
Columns must total 1



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



estimate  $P(R|S) = 0.99$   $P(\neg R|S) = 0.01$   $P(R|\neg S) = 0.03$   $P(\neg R|\neg S) = 0.97$



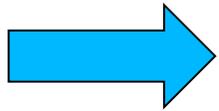
In general will not be equal

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

S

$\neg$ S

	R $\wedge$ S	R $\wedge$ $\neg$ S
R	True positive $P(R S) = 0.99$	False positive $P(R \neg S) = 0.03$
$\neg$ R	False negative $P(\neg R S) = 0.01$	True negative $P(\neg R \neg S) = 0.97$
	$\neg$ R $\wedge$ S	$\neg$ R $\wedge$ $\neg$ S



R: Test **R**esult  
is positive

Rows must  
total between  
0 and 2

These probabilities  
can be empirically  
estimated

Columns must total 1

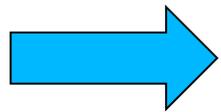
We will continue with these numbers

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

S

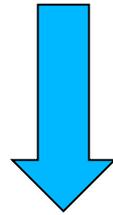
$\neg$ S

		S	$\neg$ S
R	$R \wedge S$	True positive $P(R S) = 0.99$	$R \wedge \neg S$ False positive $P(R \neg S) = 0.01$
$\neg R$	$\neg R \wedge S$	False negative $P(\neg R S) = 0.01$	$\neg R \wedge \neg S$ True negative $P(\neg R \neg S) = 0.99$



R: Test **R**esult is positive

These probabilities can be empirically estimated



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

S

$\neg$ S

	R $\wedge$ S	R $\wedge$ $\neg$ S
R	True positive $P(S R) = ??$	False positive $P(\neg S R) = ??$
	$\neg$ R $\wedge$ S	$\neg$ R $\wedge$ $\neg$ S
$\neg$ R	False negative $P(S \neg R) = ??$	True negative $P(\neg S \neg R) = ??$

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

- $P(S|R) = \frac{(P(S) \times P(R|S))}{(P(S) \times P(R|S) + P(\neg S) \times P(R|\neg S))}$
- $P(\neg S|R) = 1 - P(S|R)$
- $P(S|\neg R) = \frac{(P(S) \times P(\neg R|S))}{(P(S) \times P(\neg R|S) + P(\neg S) \times P(\neg R|\neg S))}$
- $P(\neg S|\neg R) = 1 - P(S|\neg R)$

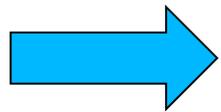
We will continue with these numbers

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

S

$\neg$ S

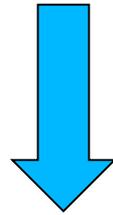
	R $\wedge$ S	R $\wedge$ $\neg$ S
R	True positive $P(R S) = 0.99$	False positive $P(R \neg S) = 0.01$
	$\neg$ R $\wedge$ S	$\neg$ R $\wedge$ $\neg$ S
$\neg$ R	False negative $P(\neg R S) = 0.01$	True negative $P(\neg R \neg S) = 0.99$



R: Test **R**esult is positive

These probabilities can be empirically estimated

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



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

S

$\neg S$

	R $\wedge$ S	R $\wedge$ $\neg S$
R	True positive $P(S R) = 0.009804$	False positive $P(\neg S R) = 0.990196$
	$\neg R \wedge S$	$\neg R \wedge \neg S$
$\neg R$	False negative $P(S \neg R) = 0.000001$	True negative $P(\neg S \neg R) = 0.999999$

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

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

$P(S R)$	requires	$P(R \neg S)$
0.01		0.01
0.09		0.001
0.5		0.0001
0.9		0.00001
0.99		0.000001



