# Learning Objectives

1. Learn how to calculate basic accuracy statistics such as sensitivity, specificity, likelihood ratios and AUC

2. Understand reasons for differences in diagnostic accuracy:  real differences, bias, random variation, cut-offs.

3. Understand the difference between tests conducted under ideal conditions vs real conditions

4. Understand the role of higher-level approaches to performance evaluation
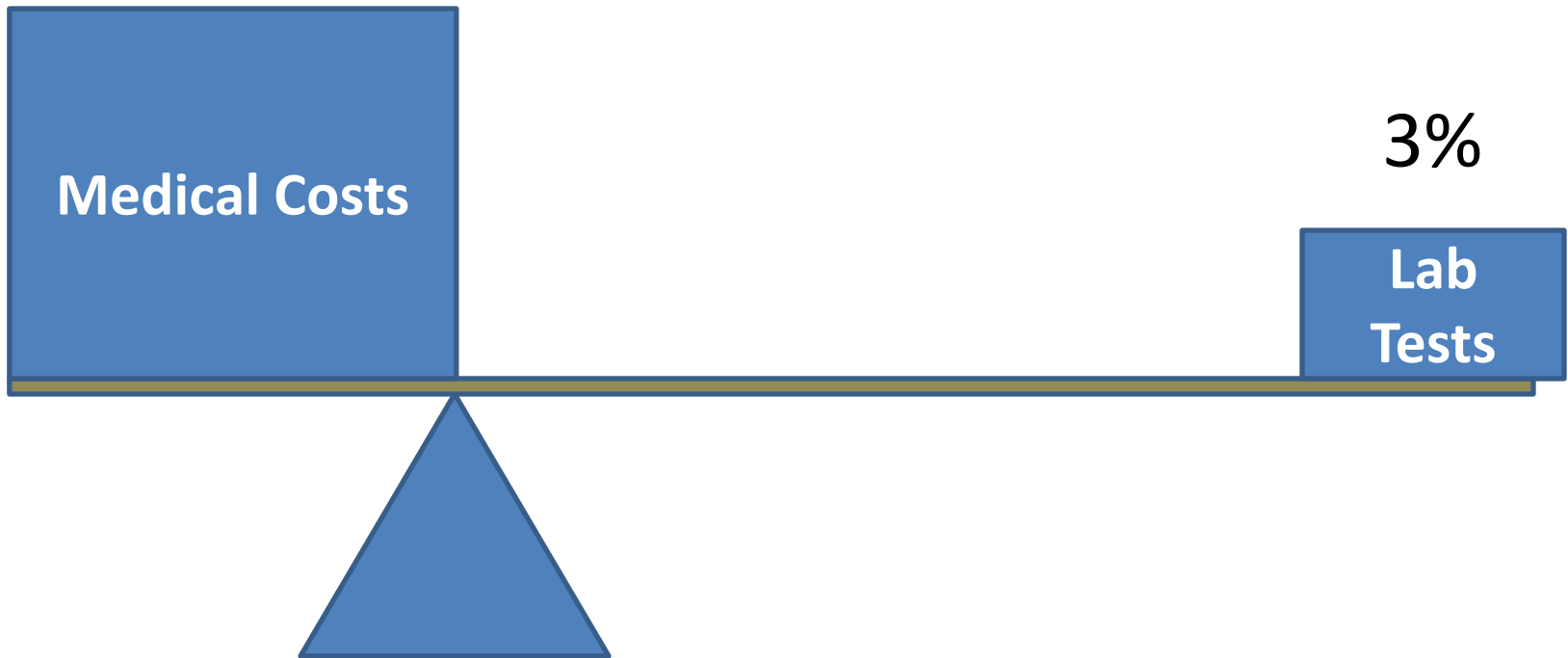
# No Disclosures

# Testing a Test :
# Beyond Sensitivity and Specificity:

**Robert Schmidt MBA MD PhD MMed**

# Tests are Central to Medicine

- Diagnosis

- Prognosis

- Monitoring

- Management

# Tests Exert Great Leverage

**Medical Costs**

3%

**Lab Tests**

# Hierarchy of Effectiveness

Societal Impact

Cost effectiveness

Clinical effectiveness

Clinical performance

Analytical performance

# What this talk is about

**Evaluating Tests:**
- Accuracy
- Usefulness
- Test Comparisons
- Limitations
- Future Directions

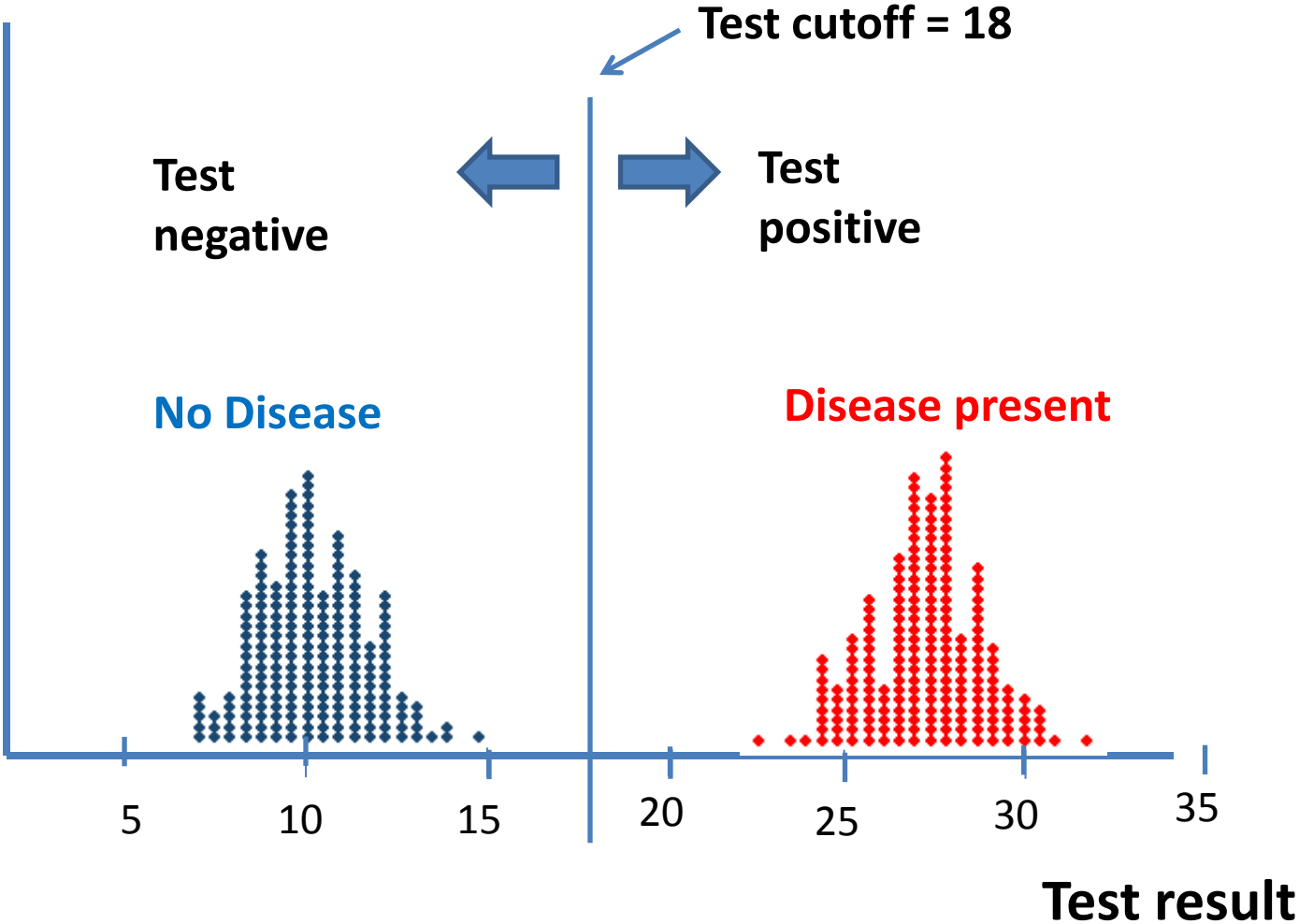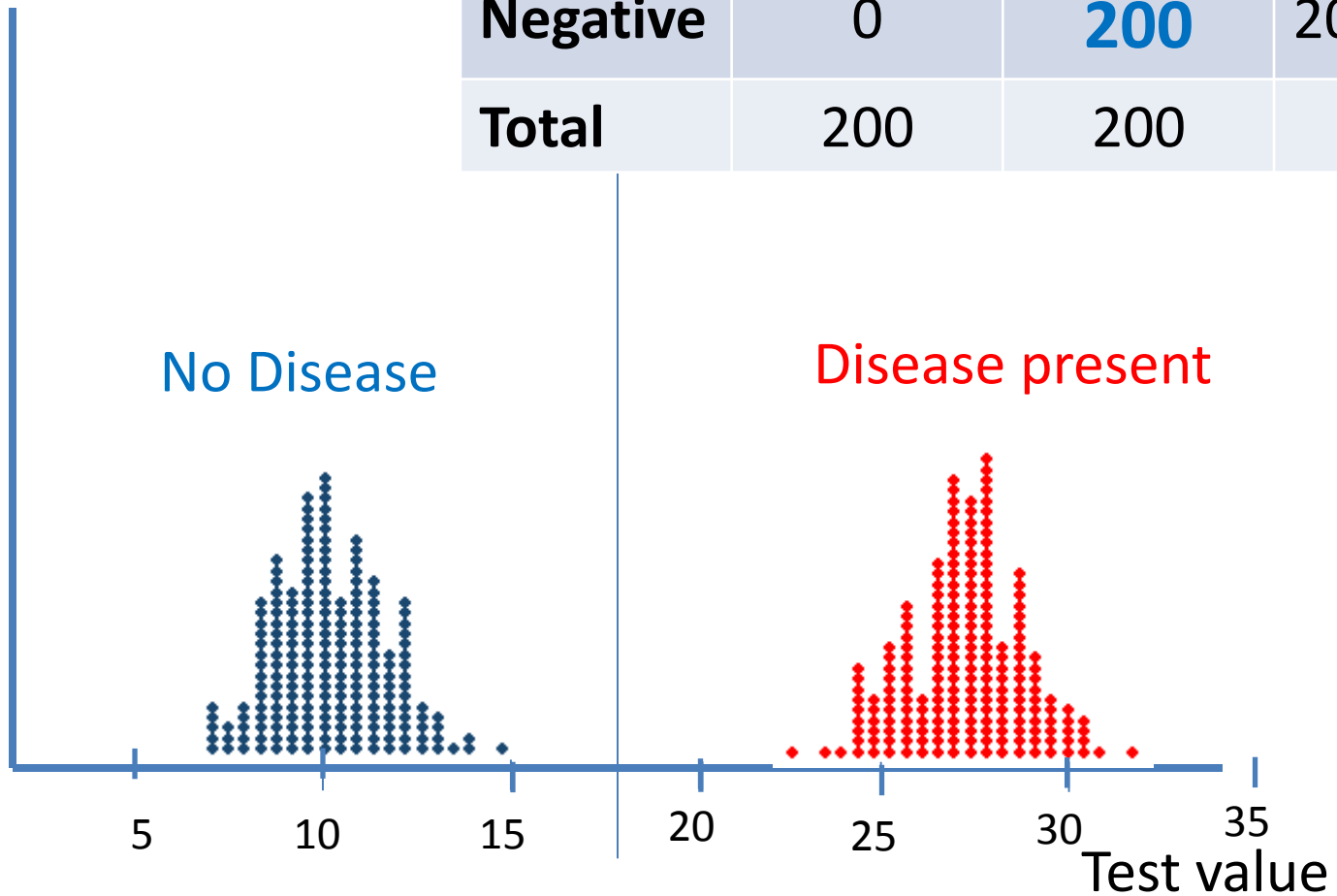| Societal Impact |
| Cost effectiveness |
| Clinical effectiveness |
| Clinical performance |
| Analytical performance |

# Case:

Your father has just returned from his annual physical.  His doctor suggested that he consider a prostatic specific antigen (PSA) test to screen for prostate cancer.  He is unsure what to do and asks your advice.  Should he take the test?
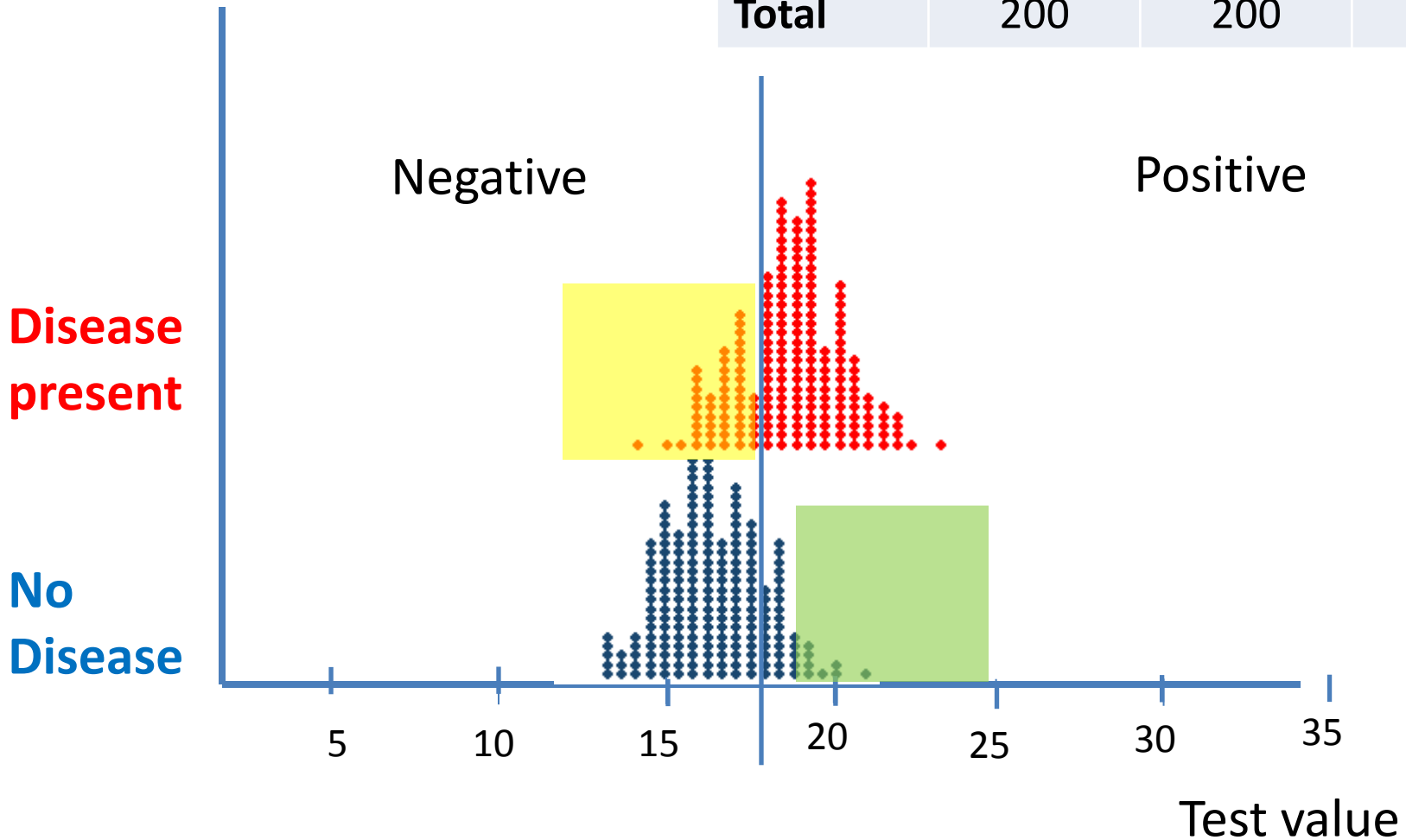
# The basic task: classification

**Perfect Test:**

| Test | Disease | | |
| --- | --- | --- | --- |
| | **Present** | **Absent** | Total |
| **Positive** | **200** | 0 | 200 |
| **Negative** | 0 | **200** | 200 |
| **Total** | 200 | 200 | |

No Disease

Disease present

5    10    15    20    25    30    35

Test value

# Not so Perfect Test:

| Test | Disease | | Total |
|---|---|---|---|
| | **Present** | **Absent** | **Total** |
| **Positive** | 148 | **40** | 188 |
| **Negative** | **52** | 160 | 212 |
| **Total** | 200 | 200 | |

| | Disease | | |
|---|---|---|---|
| **Test** | **Present** | **Absent** | **Total** |
| **Positive** | **148** | 40 | 188 |
| **Negative** | 52 | 160 | 212 |
| **Total** | 200 | 200 | |

**Negative**  **Positive**

**Disease present**

**No Disease**

**True Positives**

5   10   15   20   25   30   35

**Test value**

| Test | Disease | | Total |
| --- | --- | --- | --- |
| | Present | Absent | |
| Positive | 148 | 40 | 188 |
| Negative | **52** | 160 | 212 |
| Total | 200 | 200 | |

Negative     Positive

Disease present

No Disease

False negatives

5     10     15     20     25     30     35

Test value

| Test | Disease | | |
| | Present | Absent | Total |
|---|---|---|---|
| Positive | 148 | 40 | 188 |
| Negative | 52 | **160** | 212 |
| Total | 200 | 200 | |

|  | Disease | | |
|---|---|---|---|
| Test | Present | Absent | Total |
| Positive | 148 | **40** | 188 |
| Negative | 52 | 160 | 212 |
| Total | 200 | 200 | |

Negative   Positive

Disease present

No Disease

False positives

5    10    15    20    25    30    35

Test value

# How well did we classify those <u>with</u> disease?

Sensitivity = 148/200 =  74%



| Test | Disease | | |
| | Present | Absent | Total |
|---|---|---|---|
| **Positive** | **148** | 40 | 188 |
| **Negative** | **52** | 160 | 212 |
| **Total** | **200** | 200 | |

Negative   Positive

**Disease present**

**No Disease**

5    10    15    20    25    30    35

Test value

# How well did we classify those __without__ disease?

Specificity = 160/200 =  80%

| Test | Disease Present | Absent | Total |
|------|---------|--------|-------|
| **Positive** | 148 | **40** | 188 |
| **Negative** | 52 | **160** | 212 |
| **Total** | 200 | **200** | |

**Disease present**

**No Disease**

Negative    Positive

5    10    15    20    25    30    35

Test value

# Sensitivity = accuracy in the **diseased** group
# Specificity = accuracy in the **nondiseased** group

| Test | Disease Present | Disease Absent | Total |
|------|---------|--------|-------|
| **Positive** | 148 | 40 | 188 |
| **Negative** | 52 | 160 | 212 |
| **Total** | 200 | 200 | |

# Two useful mnemonics*

**SnNout:**

High **_Sen_**sitivity Test with a **_N_**egative result rules **_out_**

**SpPin:**

High **_Sp_**ecificity Test with a **_p_**ositive result rules **_in_**

# PSA test

- Sensitivity 90%
- Specificity 20%

**How might this test be useful?**
- **SnNout?**
- **SpPin?**

# Results May Vary……

**Sensitivity of PSA Tests**

# Specificity of PSA Test

# Sensitivity and Specificity Depend on Cutoff Values

|  | Cutoff A | Cutoff B |
|---|---|---|
| **Sensitivity** | 74% | 22% |
| **Specificity** | 80% | 98% |

| Test | Disease Present | Absent | Total |
|------|---------|--------|-------|
| **Positive** | 148 | 40 | 188 |
| **Negative** | 52 | 160 | 212 |
| **Total** | 200 | 200 | |

**Disease present**

**No Disease**

Sensitivity = 74%
Specificity = 80%

5  10  15  20  25  30  35

# Threshold Effects on Test Performance

| Test | Disease | | Total |
|---|---|---|---|
| | **Present** | **Absent** | **Total** |
| **Positive** | 44 | 4 | 88 |
| **Negative** | 156 | 196 | 312 |
| **Total** | 200 | 200 | |

Sensitivity = 22%
Specificity = 98%

**Disease present**

**No Disease**

5   10   15   20   25   30   35

# Threshold Effects on Test Performance

| Test | Disease | | |
|------|---------|--------|-------|
| | **Present** | **Absent** | **Total** |
| **Positive** | 196 | 156 | 312 |
| **Negative** | 4 | 44 | 88 |
| **Total** | 200 | 200 | |

**Disease present**

**No Disease**

Sensitivity = 98%
Specificity = 22%

5   10   15   20   25   30   35

# Tradeoff: Specificity vs Sensitivity

| Threshold | Sensitivity | Specificity |
|-----------|-------------|-------------|
| 15 | 98 | 22 |
| 18 | 74 | 80 |
| 20 | 22 | 98 |

# How to Compare Tests

|  | Sensitivity | |
| --- | --- | --- |
| Specificity | Test A | Test B |
| 0.1 | 0.98 | 0.95 |
| 0.2 | 0.97 | 0.91 |
| 0.3 | 0.95 | 0.88 |
| 0.4 | 0.91 | 0.78 |
| 0.5 | 0.87 | 0.70 |
| 0.6 | 0.80 | 0.56 |
| 0.7 | 0.70 | 0.40 |
| 0.8 | 0.55 | 0.20 |
| 0.9 | 0.38 | 0.10 |
| 1.0 | 0.10 | 0.01 |

# Receiver Operating Characteristic (ROC) Curve



Test A

Test B

| Specificity | Sensitivity | |
| --- | --- | --- |
| | Test A | Test B |
| 0.1 | 0.98 | 0.95 |
| 0.2 | 0.97 | 0.91 |
| 0.3 | 0.95 | 0.88 |
| 0.5 | 0.87 | 0.70 |
| 0.6 | 0.80 | 0.56 |
| 0.7 | 0.70 | 0.40 |
| 0.8 | 0.55 | 0.20 |
| 1.0 | 0.10 | 0.01 |

Sensitivity = TPR

1 – specificity = FPR

1.0

0.5

0.0

0.5

1.0

ROC Curve for The Perfect Test

# Test Performance is Related to Area Under the Curve (AUC)

# Area Under the Curve (AUC)



**Perfect Test**
**AUC = 1.0**

**Real Test**

**Useless Test**
**AUC = 0.5**

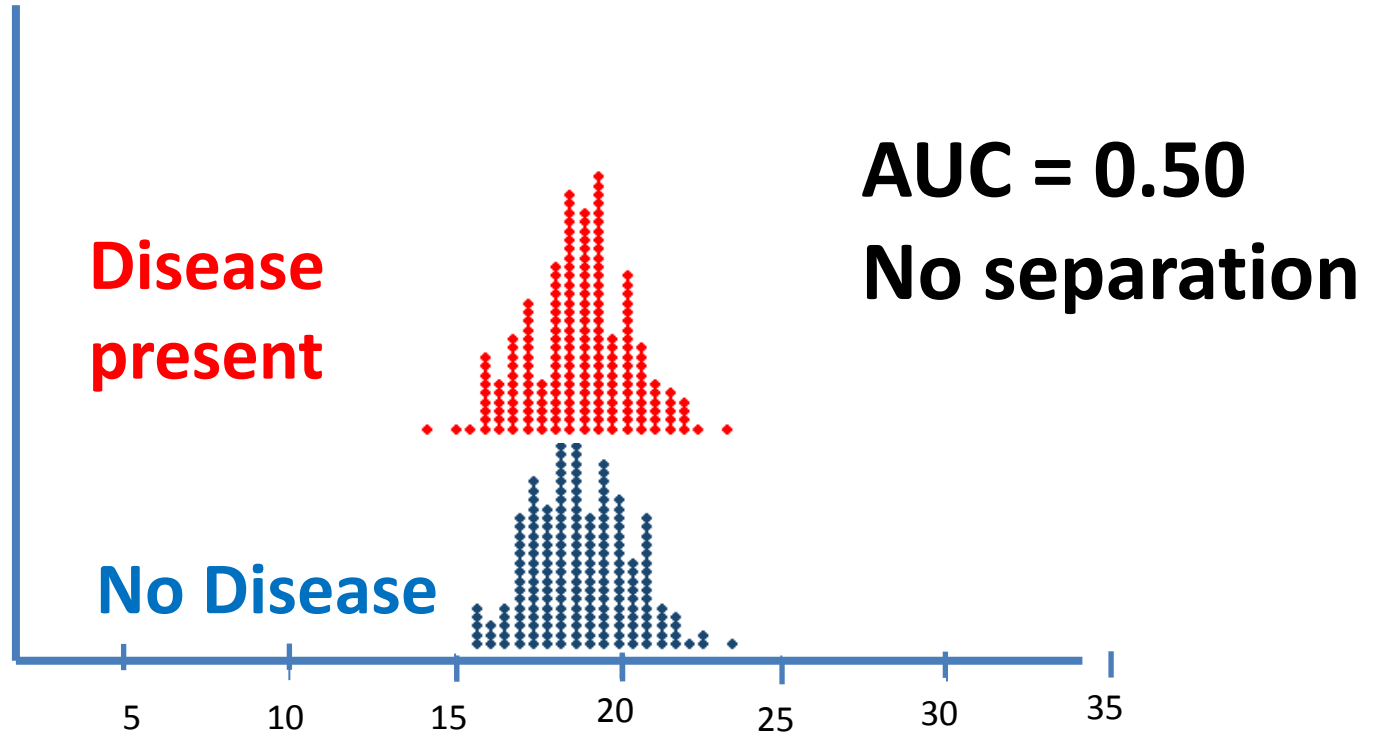# Does the AUC mean anything?

Does the AUC Mean Anything?

# Meaning of AUC

Prob ($T_{Disease} > T_{No\ Disease}$)

**AUC = 1.0 → perfect separation**



**No Disease**   **Disease**

Test value, T

# Meaning of AUC

$$\text{Prob} (T_{\text{Disease}} > T_{\text{No Disease}})$$



**Disease present**

**No Disease**

**AUC = 0.50**
**No separation**

# PSA vs PSA velocity

| Study | Area Under the Curve (AUC) | |
| --- | --- | --- |
| | PSA velocity | PSA |
| Eggener, 2005 | 0.91 | 0.88 |
| Ciatto, 2004 | 0.74 | 0.67 |
| Berger, 2007 | 0.87 | 0.65 |

# Take home message:

| Sensitivity Specificity | → | ROC Curves AUC |
|---|---|---|

Threshold Effects                    No Threshold Effects

# How do I know if a test is useful?

# Is this test useful?

# Usefulness is defined by the customer



Will this test tell me whether my patient has Prostate Cancer?
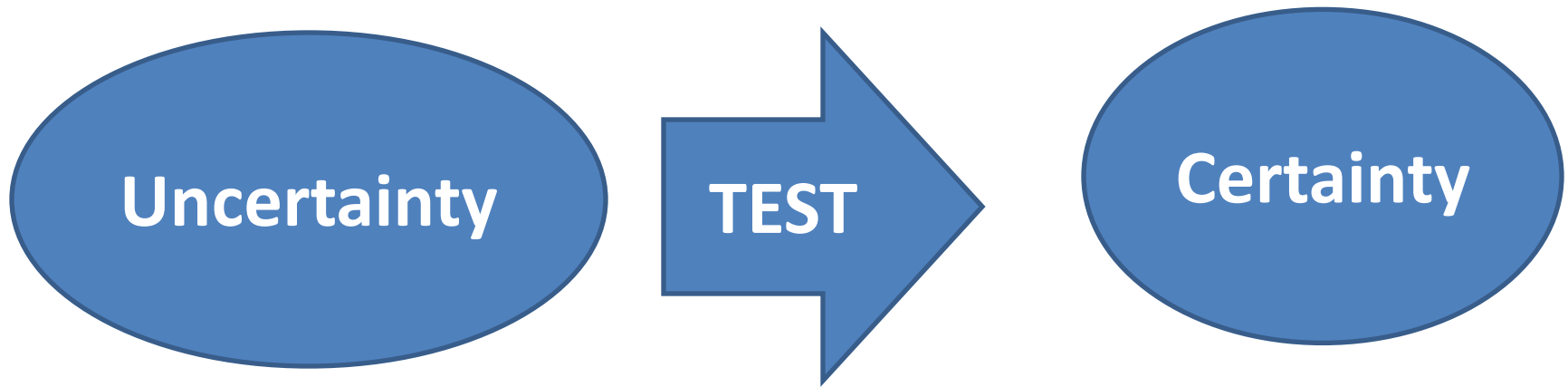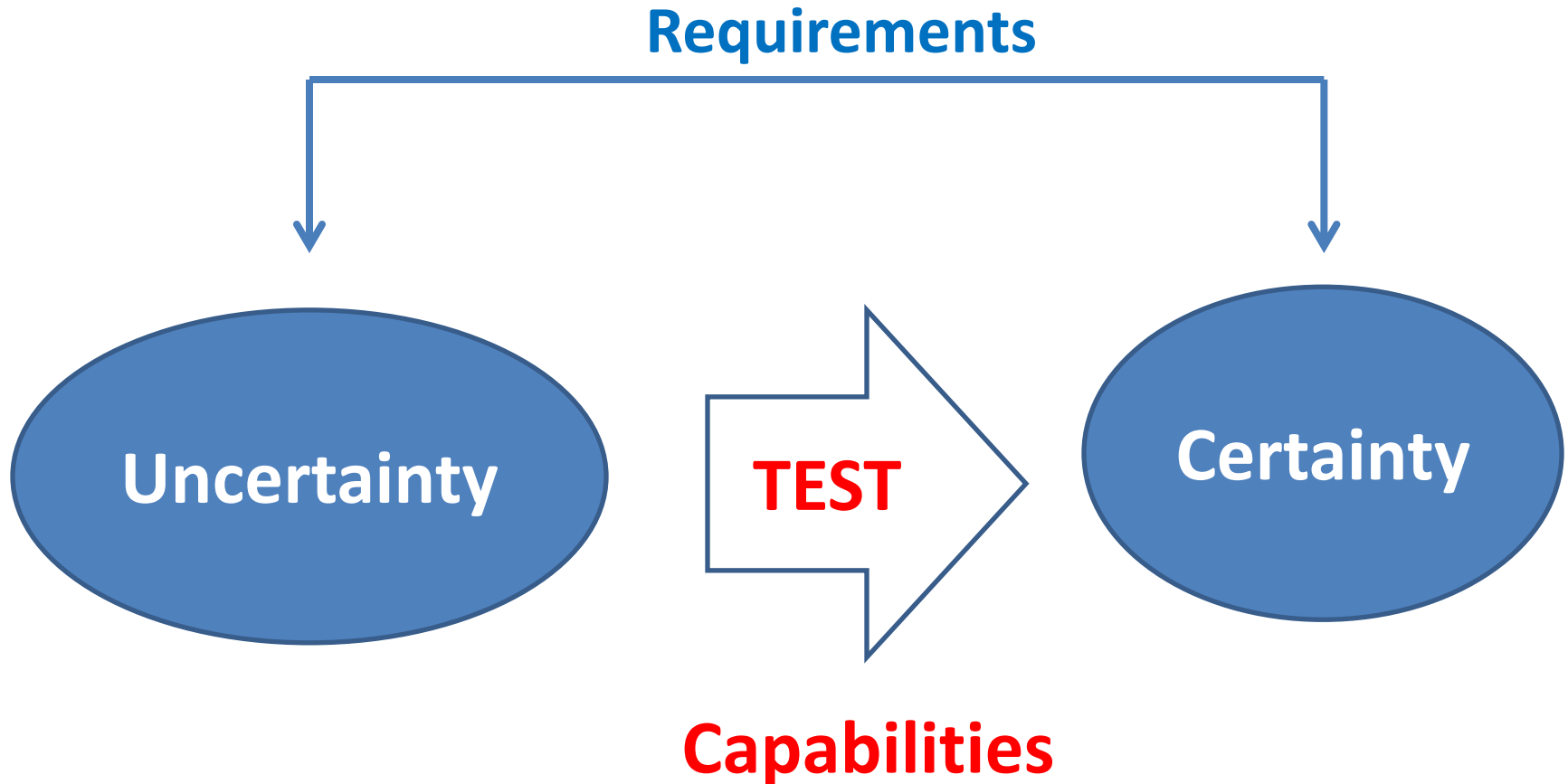
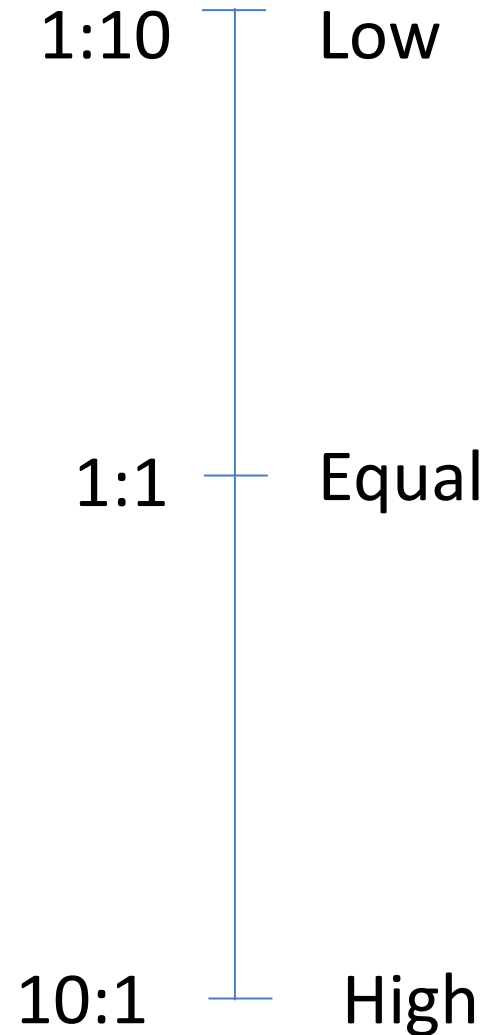# Usefulness is defined by the customer

# The customer's problem:

# Usefulness = Capabilities - Requirements
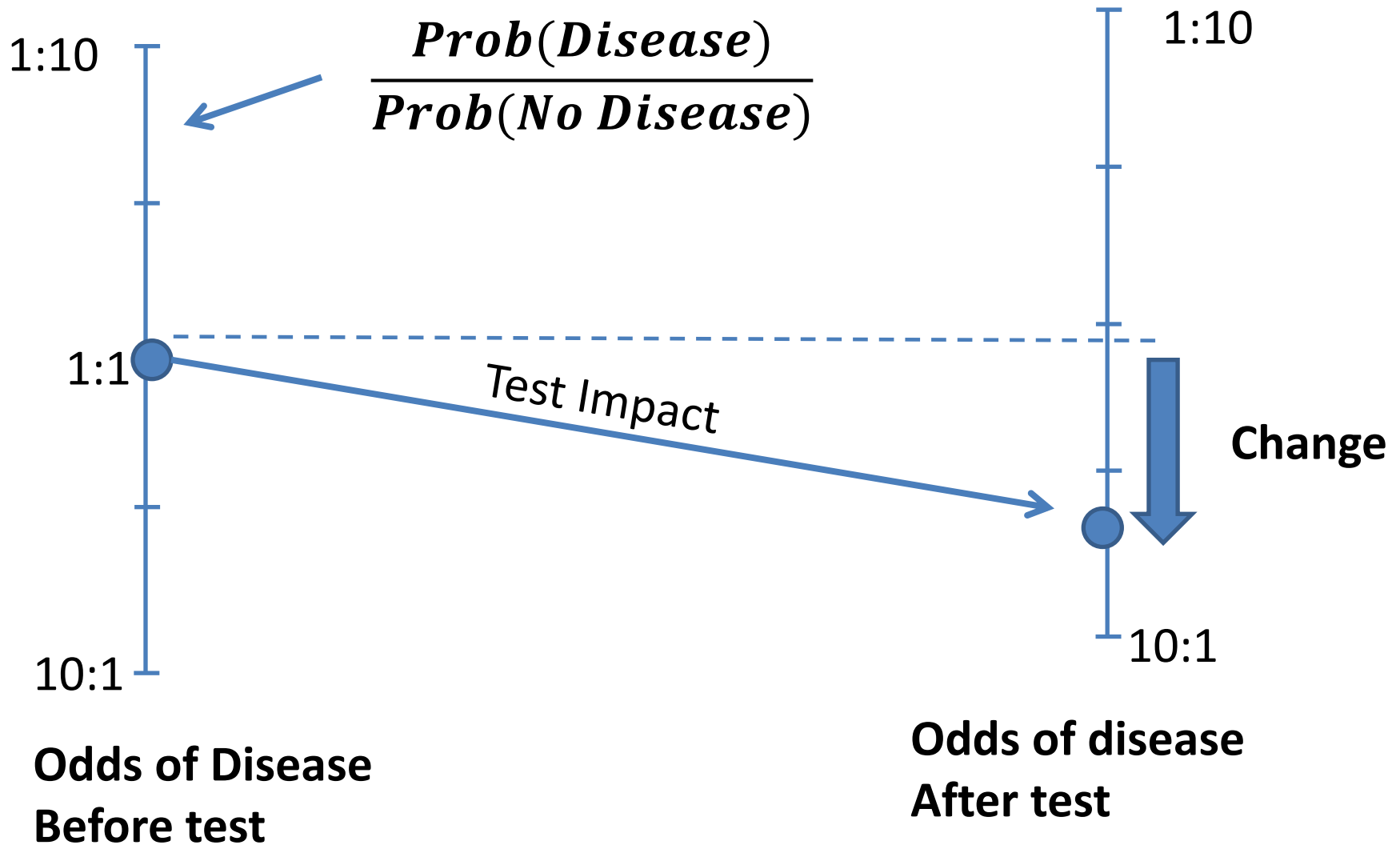
**Requirements**

**Uncertainty**

**TEST**

**Certainty**

**Capabilities**

# How to think about certainty

**The Odds Scale**

$$\textbf{Odds} = \frac{Prob(Disease)}{Prob(No\ Disease)} = \frac{P}{1-P}$$

1:10 — Low

1:1 — Equal
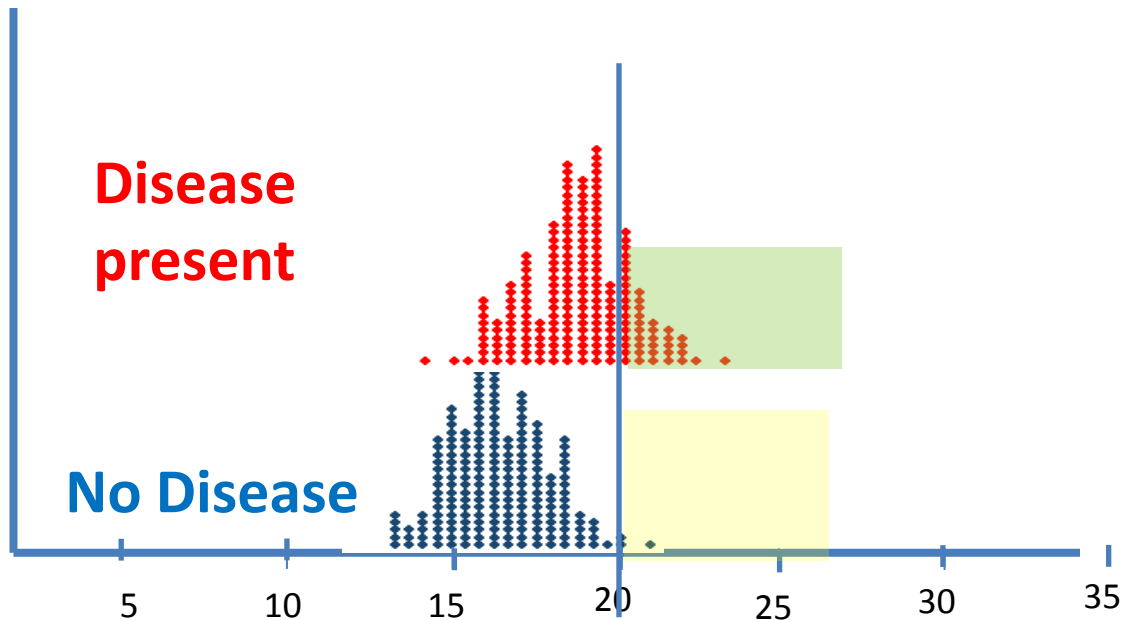
10:1 — High

# Capability = Change in Certainty



$$\frac{Prob(Disease)}{Prob(No\ Disease)}$$

1:10

1:1

10:1

**Odds of Disease**
**Before test**

Test Impact

1:10

10:1

**Change**

**Odds of disease**
**After test**

# What is the impact of a _positive_ result?

## Positive Likelihood ratio, LR+

$$\text{LR+} = \frac{Prob\ (Positive|disease)}{Prob(Positive|no\ disease)} = \frac{40/200}{4/200} = 10.0$$
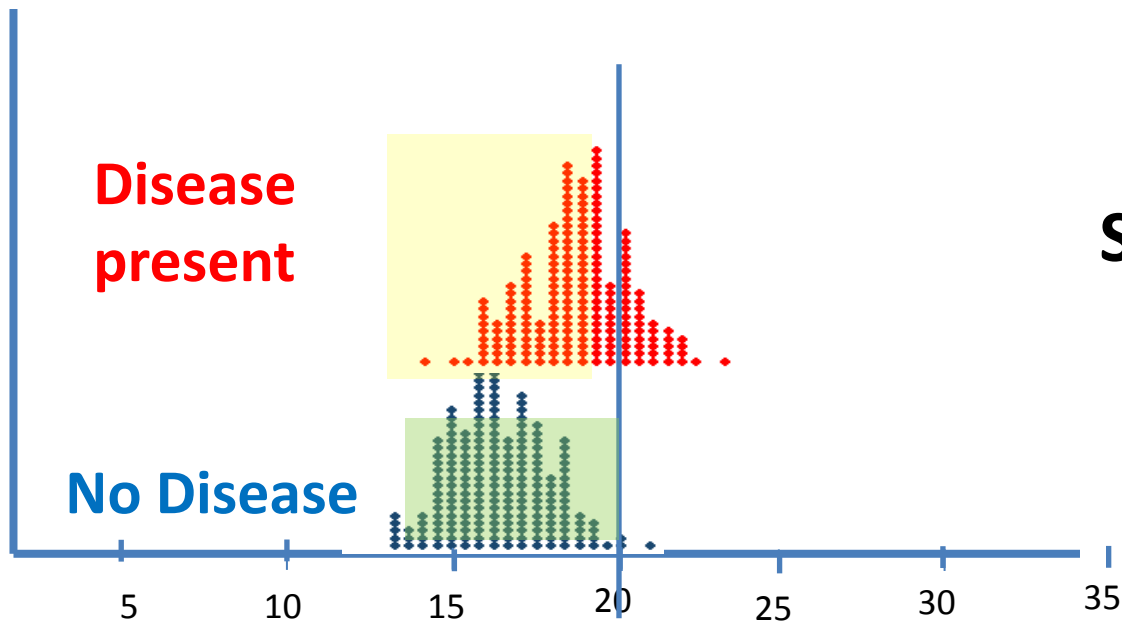


Disease present

No Disease

**Bigger LR+  is better**

# What is the impact of a _negative_ result?

## Negative Likelihood ratio, LR-

$$\text{LR-} = \frac{\textbf{\textit{Prob (Negative }|disease)}}{\textbf{\textit{Prob(Negative}|no\ disease)}} = \frac{160/200}{196/200} = 0.8$$
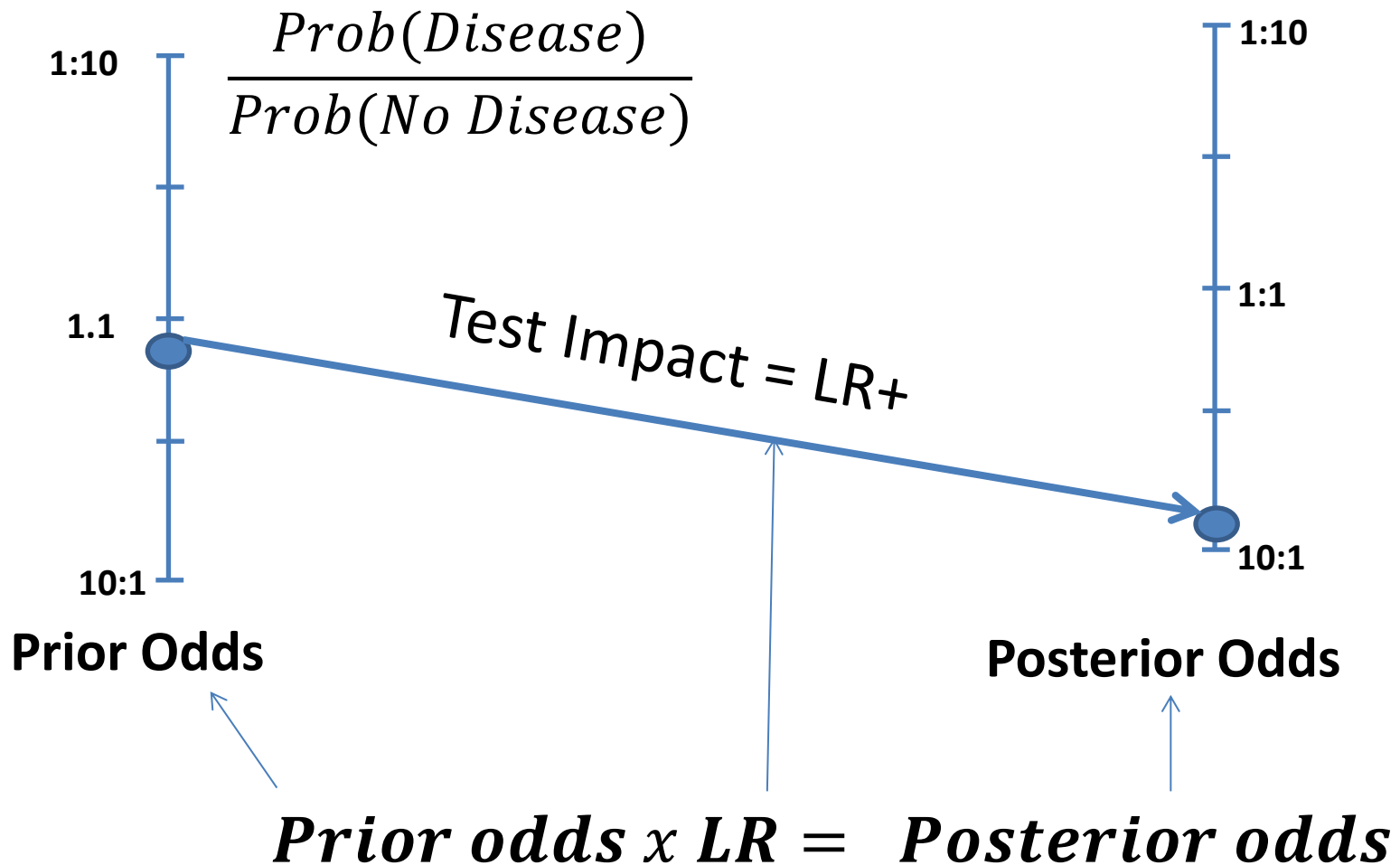


**Disease present**

**No Disease**

## Smaller LR-  is Better

# Key Relationship

$$\text{Posterior Odds} = \text{LR} \times \text{Prior Odds}$$

# Likelihood = Impact Factor

$$\frac{Prob(Disease)}{Prob(No\ Disease)}$$

1:10

1.1

10:1

**Prior Odds**

Test Impact = LR+

1:10

1:1

10:1

**Posterior Odds**

$Prior\ odds\ x\ LR =\ \ Posterior\ odds$

# Two ways to be certain:



$$\frac{Prob(Disease)}{Prob(No\ Disease)}$$

1:10

1:1

10:1

**Prior Odds**

LR-
Negative result

LR+
Positive result

1:10    Exclude

1:1

10:1

Confirm

**Posterior  Odds**
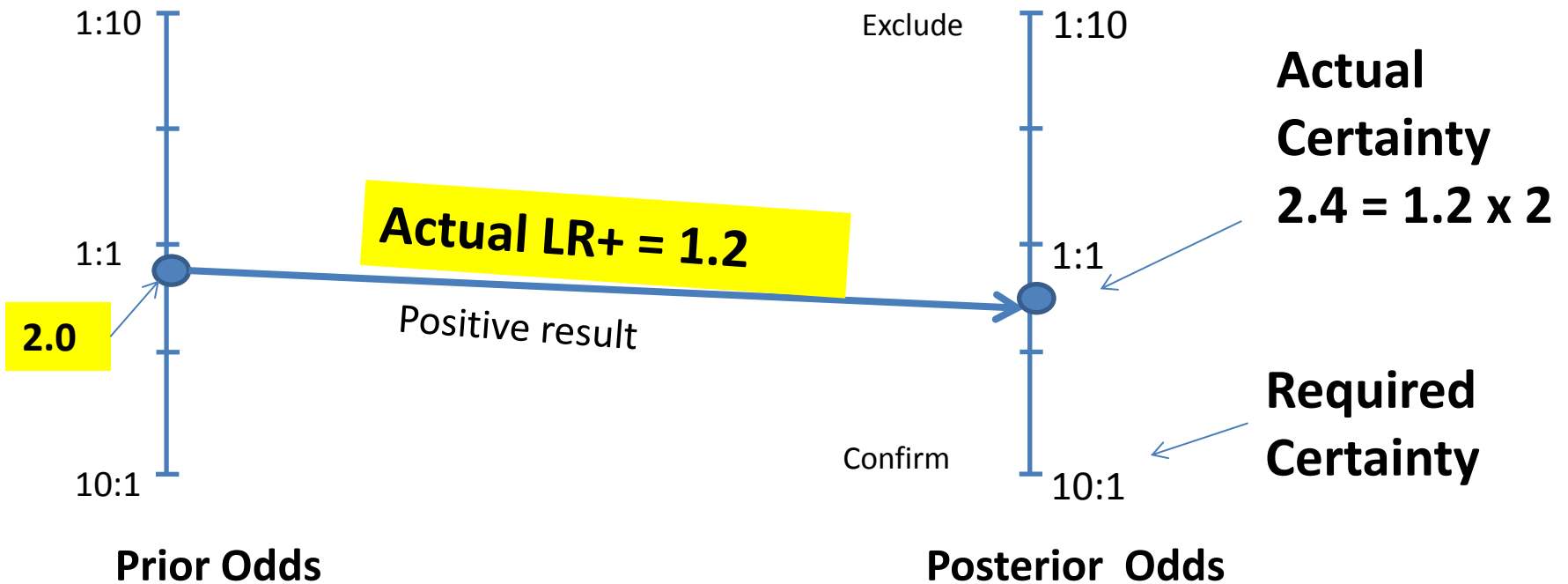
*Prior odds x Likelihood Ratio =  Posterior odds*

# A test can solve the problem if:

$$LR_{Actual} > LR_{Required}$$

$$LR^+ > \frac{PosteriorOdds}{PriorOdds} \quad \text{Confirm}$$

Or:

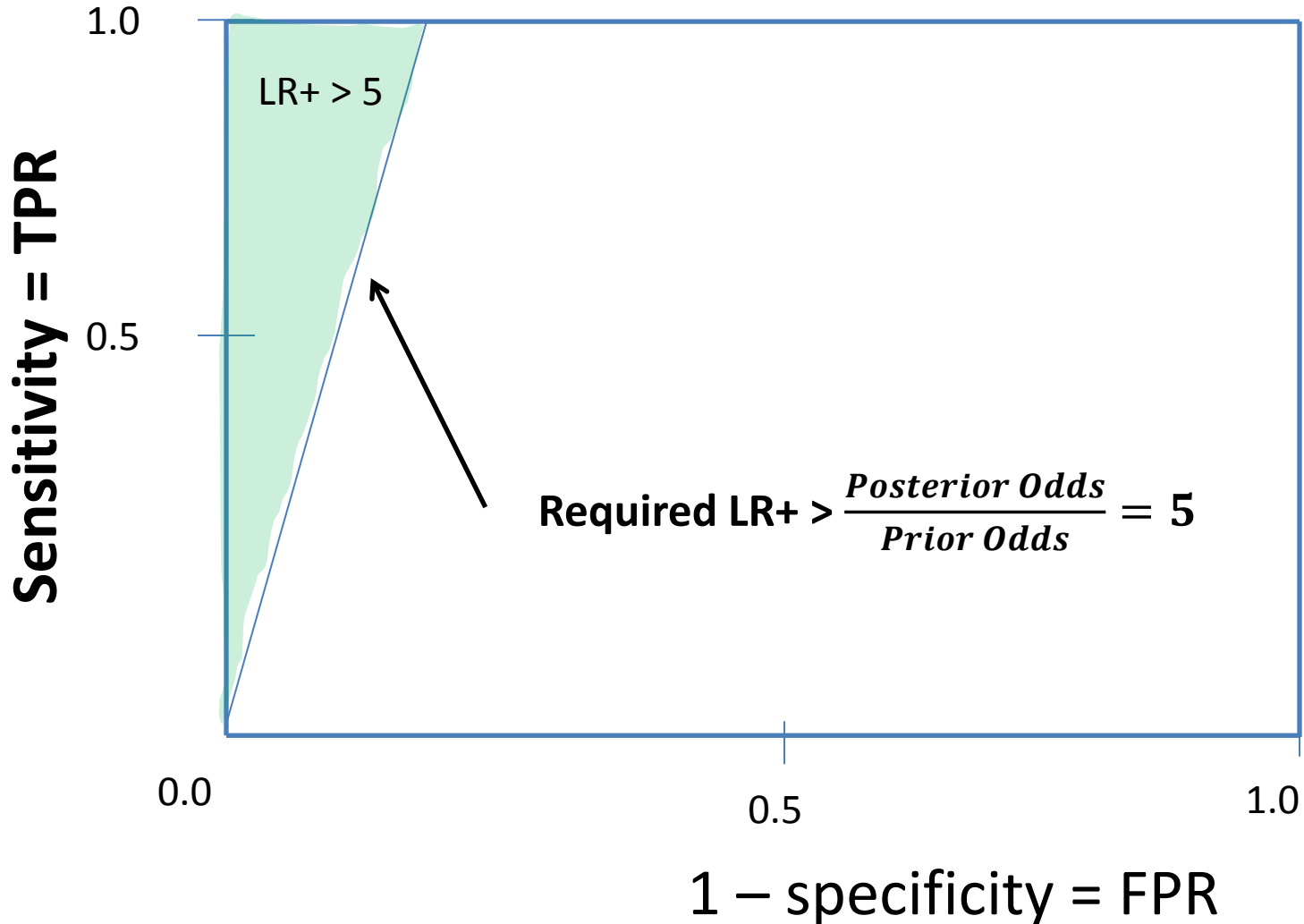$$LR^- < \frac{PosteriorOdds}{PriorOdds} \quad \text{Exclude}$$

# Rule-in, Confirmation Zone
## High LR+ is better

LR+ > 5

Required LR+ > $\frac{Posterior\ Odds}{Prior\ Odds}$ = 5

Sensitivity = TPR

1 − specificity = FPR

# Rule-out, Exclusion Zone:
## Low LR- is better

LR- < 0.2

$$\text{Required LR-} < \frac{\textit{Posterior Odds}}{\textit{Prior Odds}} = 0.2$$

Sensitivity = TPR

1 – specificity = FPR

# Capabilities > Requirements?



Sensitivity = TPR

1.0

SnNout
SpPin

SnNout

SpPin

0.5

*Required* (LR),
determined by
the customer

*actual* LR, capabilities

0.0

0.5

1.0

1 – specificity = FPR
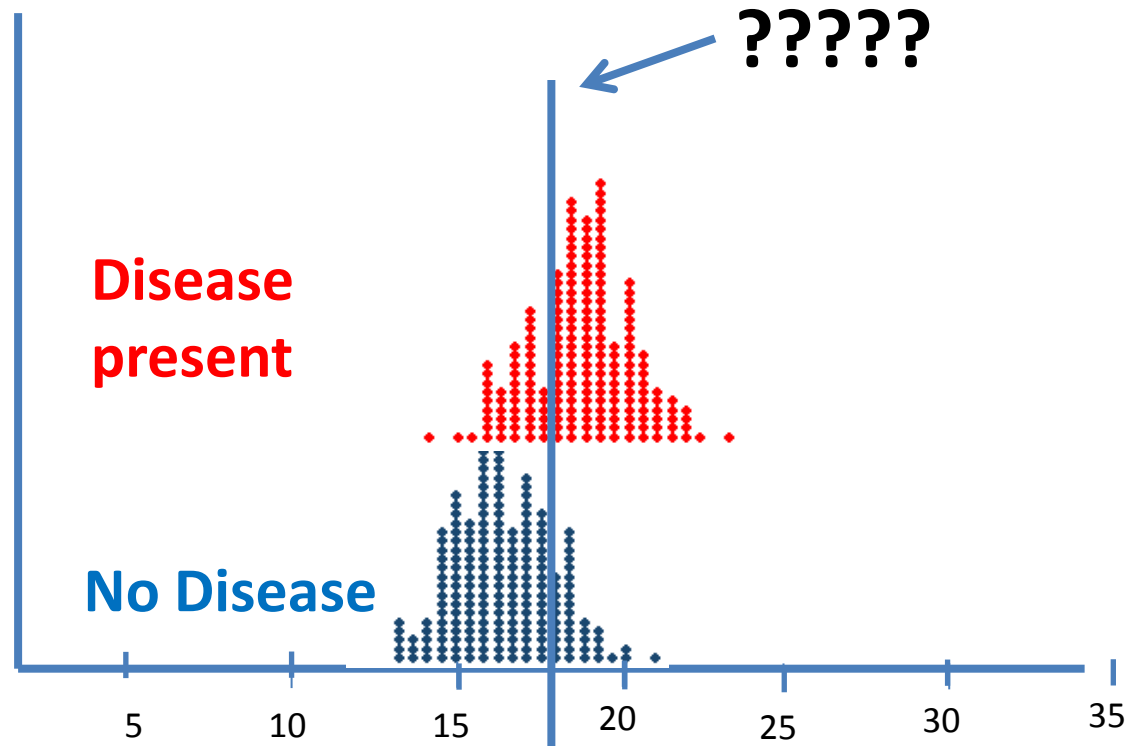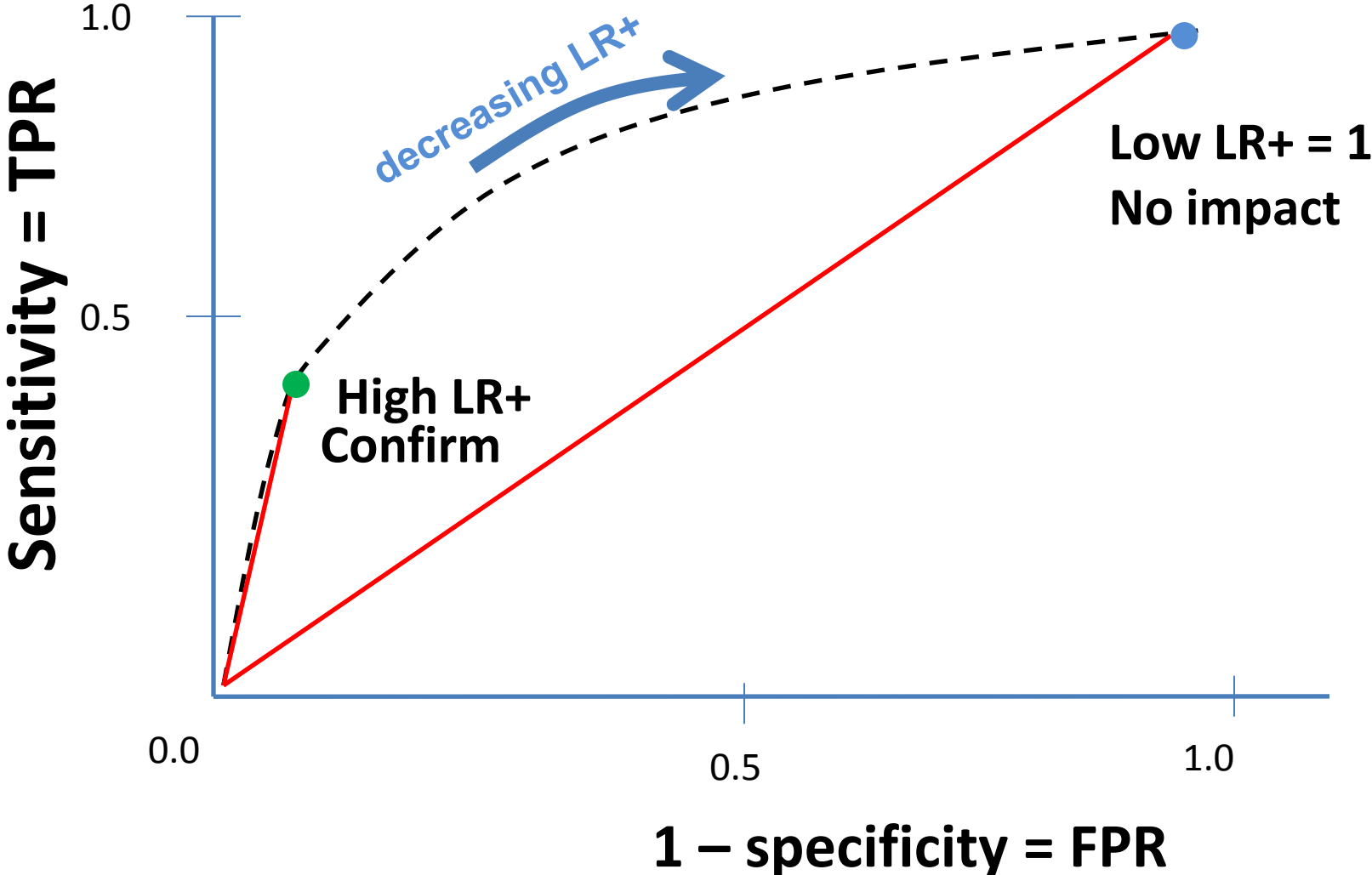
# Key Points:

- **Accuracy ≠ Usefulness**
- **Potential Usefulness = LR = f(Sn,Sp)**
- **Usefulness = Capabilities - Requirements:**
  - The objective (exclude, confirm)
  - Prior uncertainty
  - Required certainty
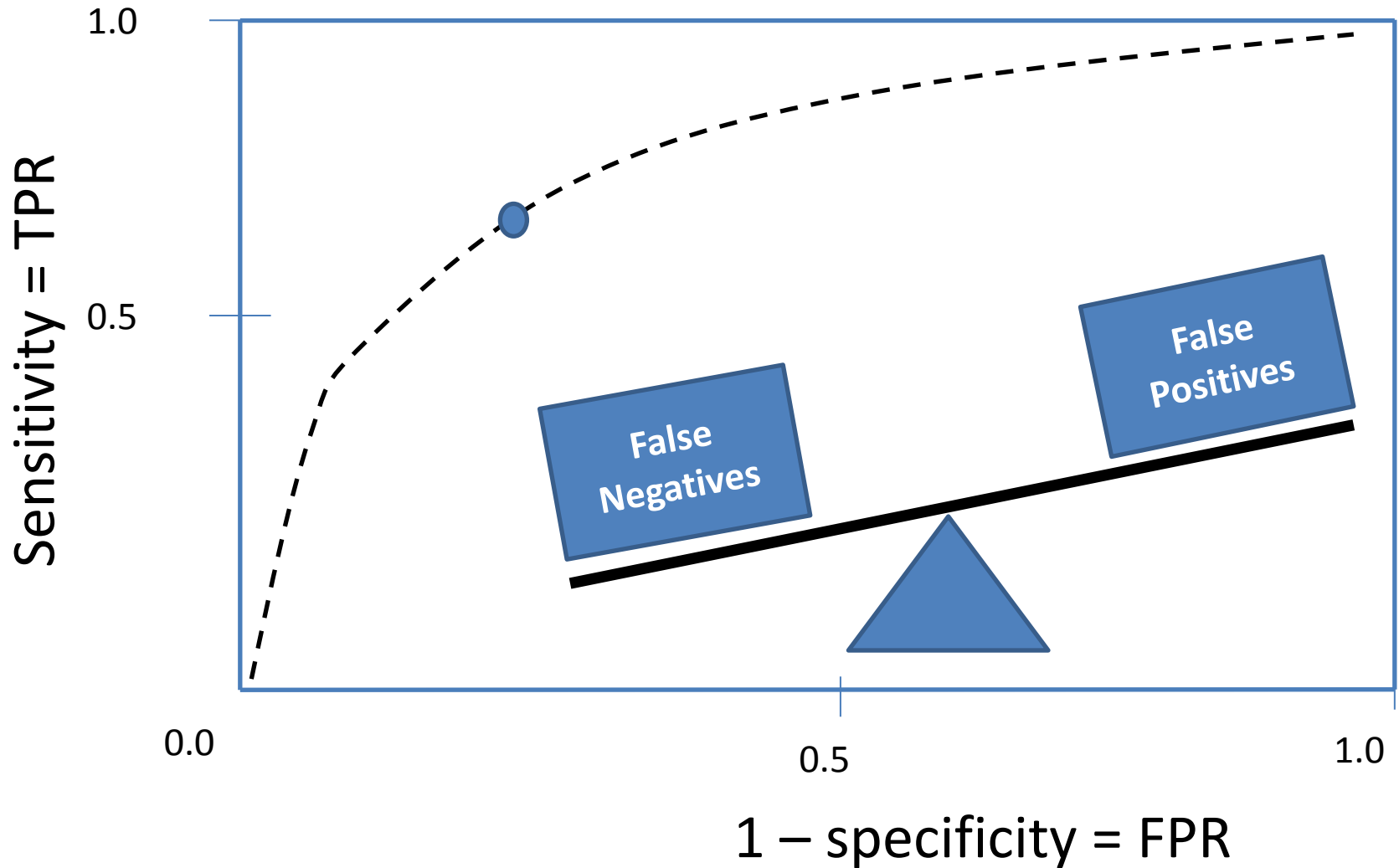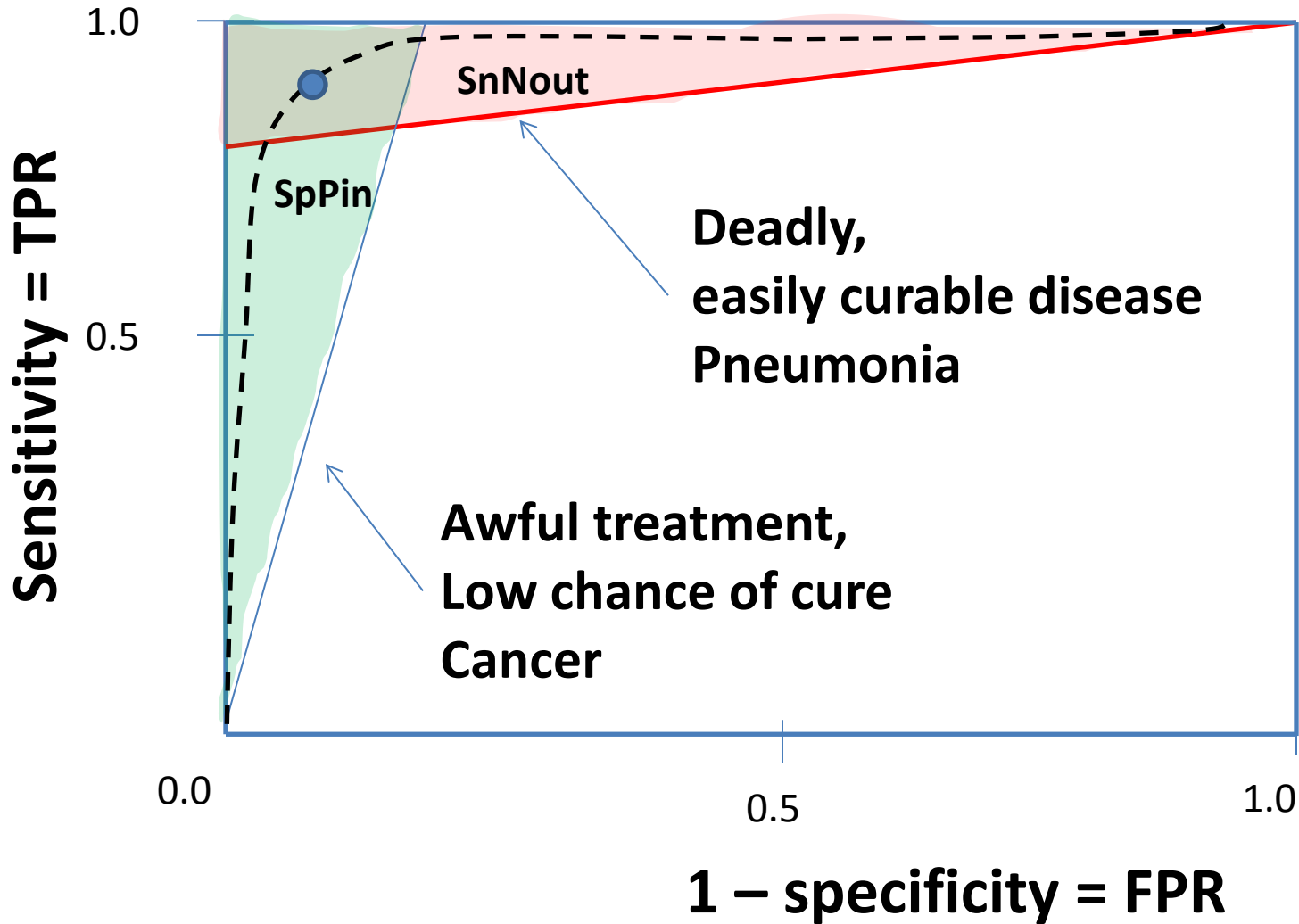  - The Test Impact (*actual* LR vs *required* LR)
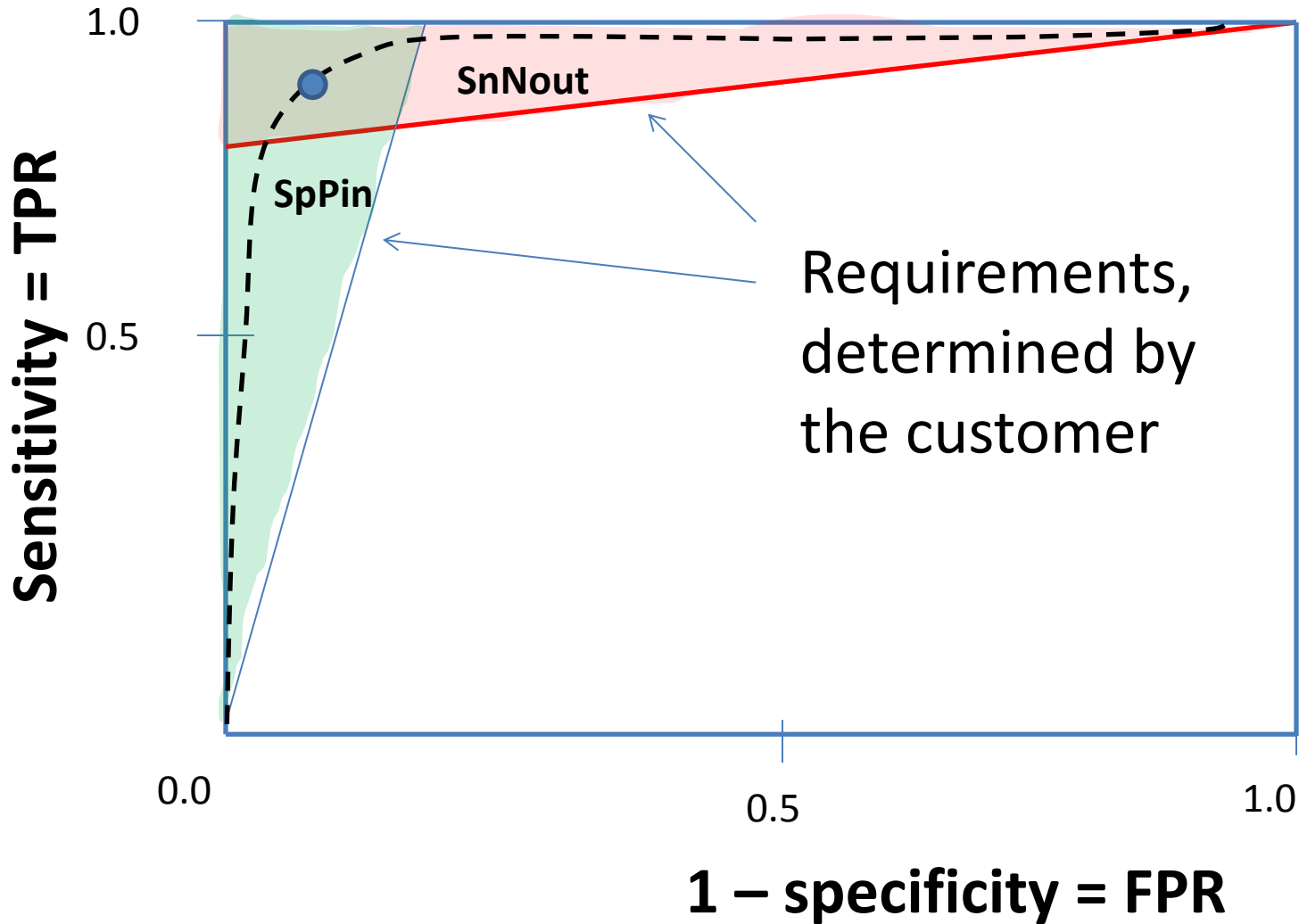
# ROC curve = set of available Likelihood Ratios

# How to select a threshold value

# Diagnostic Zones for Thresholds

Accuracy , Usefulness and Optimality

Sensitivity = TPR

1 − specificity = FPR

SnNout

SpPin

Requirements, determined by the customer

# Key Points – Setting Thresholds

- **Comparing Tests**
  - Thresholds are a nuisance
  - ROC/AUC facilitates comparisons of *diagnostic* accuracy

- **Using Tests**
  - Thresholds are required
    - Define a test
    - Link capabilities and requirements
    - Can be set to optimize performance
      - Optimum is context dependent
      - Depends on error costs

# Comparing Test Performance
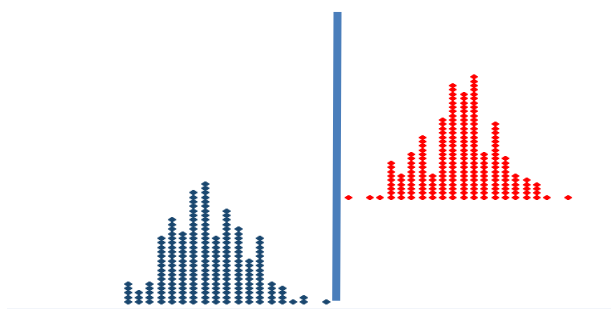
**Why do test results differ?**
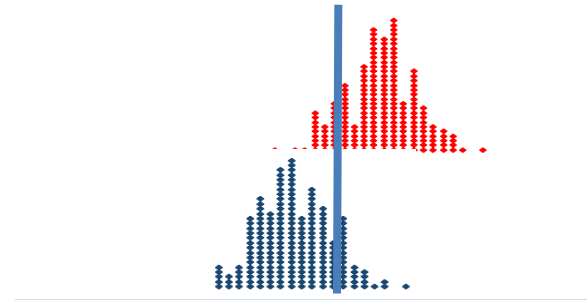
1. **True differences**
2. False differences
    - (bias)
    - thresholds
3. Random variation (imprecision)

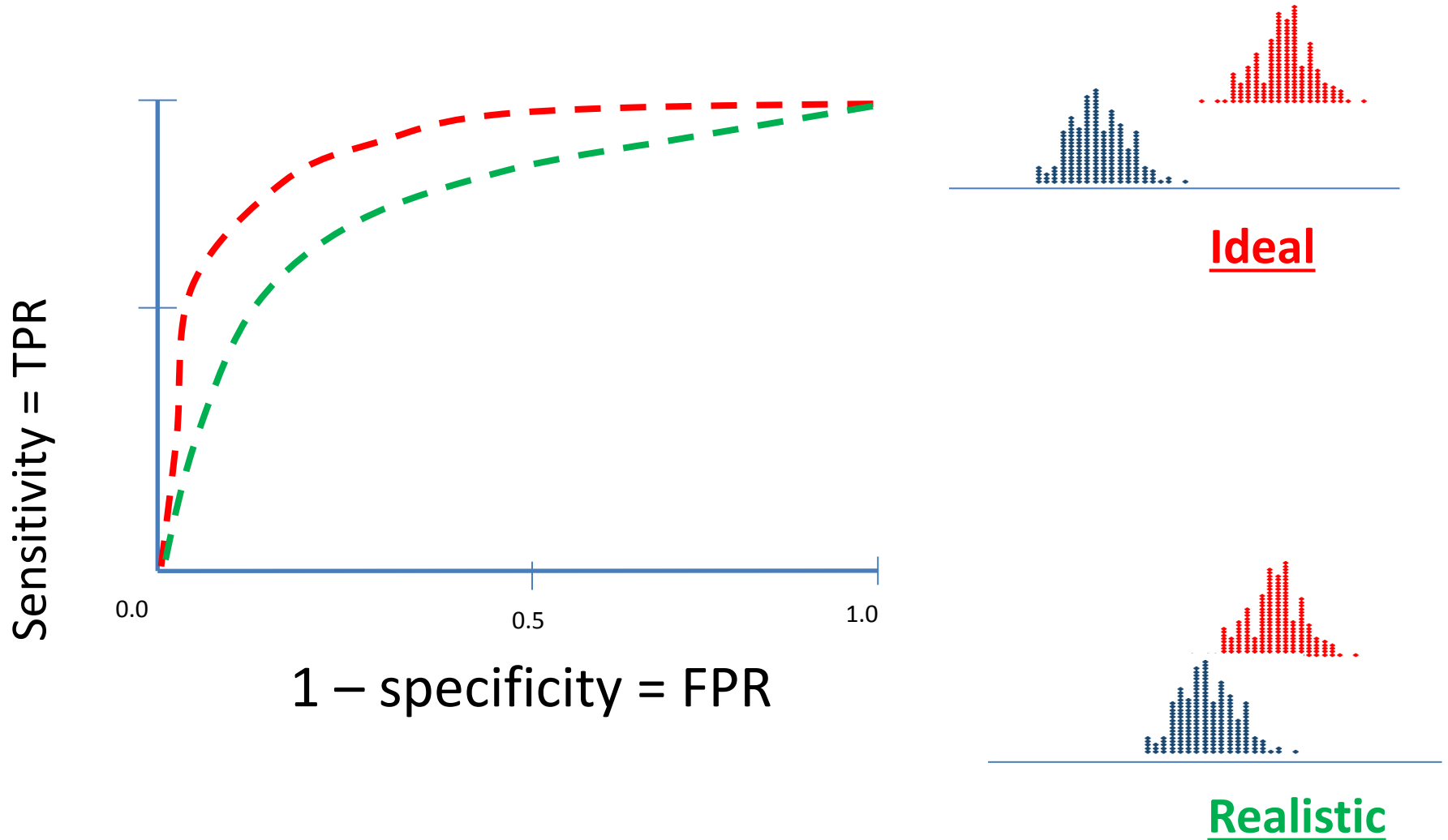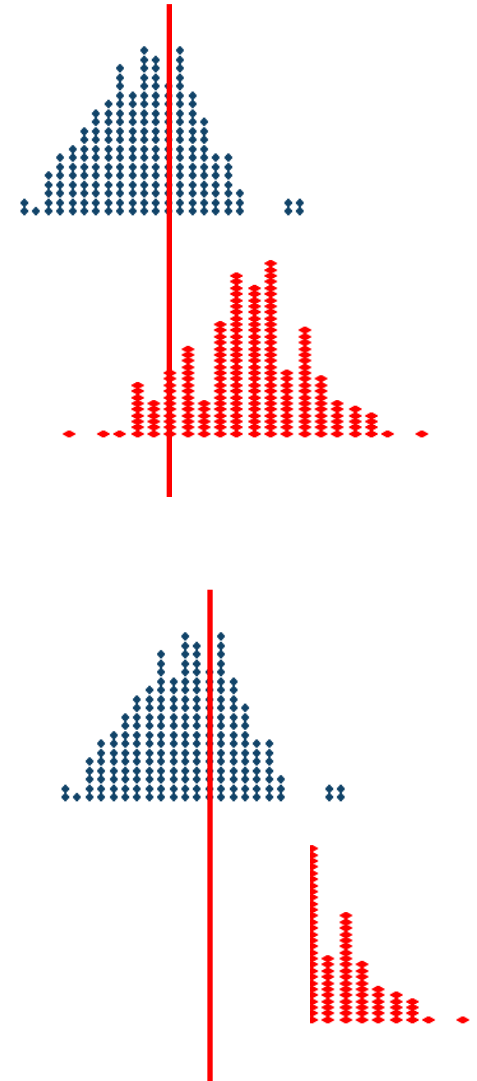# Evaluating Test Accuracy: Ideal vs Realistic Conditions

# Test Performance
## Ideal vs Realistic Conditions

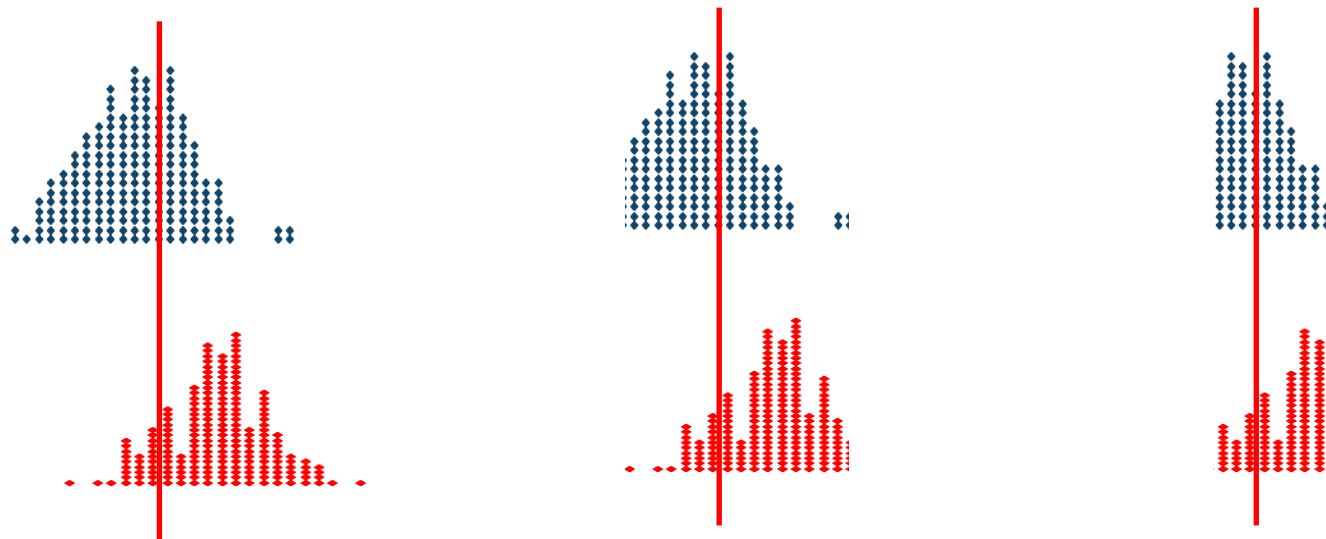# Referral Pattern & Disease Spectrum

Patient with Complaint
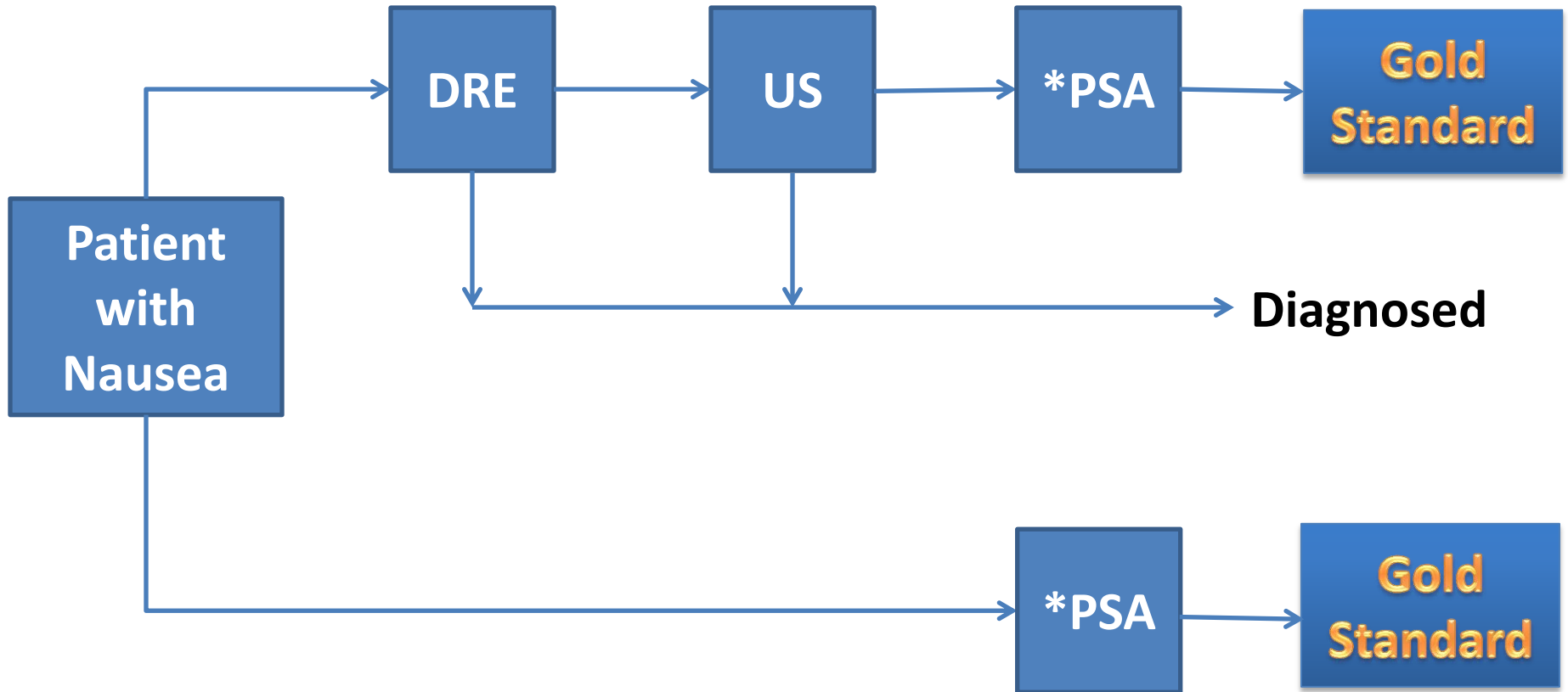
General Practitioner

Emergency Dept

# Referral Pattern & Disease Spectrum

# Effect of Prior Testing
## Will the index test perform differently?



*Index Test = the test of interest

# Defining a Test: PICCO
## Sources of <u>Real</u> Differences: Context is Everything

| | | |
|---|---|---|
| **P** | **Population** | **Setting**<br>**Exclusion/Inclusion criteria**<br>**Referral pattern**<br>**Comorbidities**<br>**Age, Gender** |
| **I** | **Index Test** | **Method (in detail)**<br>**Cutoff**<br>**Skill level** |
| **C** | **Condition** | **Disease of interest** |
| **C** | **Comparator (reference test)** | **Definition of disease** |
| **O** | **Outcome measure** | **Diagnostic accuracy**<br>**Discomfort, adverse events**<br>**Operational (TAT, Availability, cost, etc)** |

# Comparing Test Performance
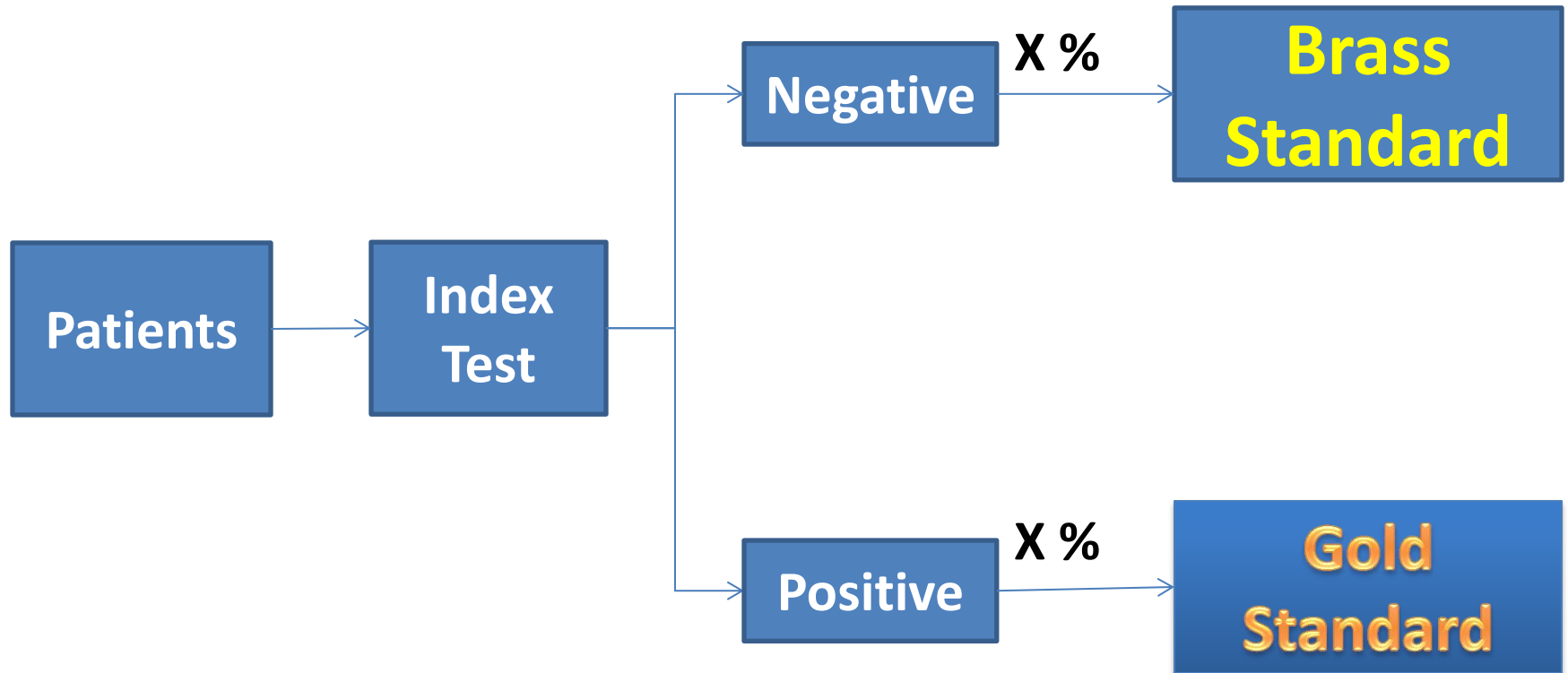
**Why do test results differ?**
1. True differences
2. False differences
   - **bias**
   - thresholds
3. Random variation (imprecision)

# Sources of Bias
## (Phantom Differences)

- Imperfect gold standard

- Verification Bias

- Indeterminate results

- Others….

# Imperfect Gold Standard
## (Differential verification)

Patients → Index Test

Index Test → Negative → **X %** → **Brass Standard**

Index Test → Positive → **X %** → **Gold Standard**

# Verification bias
(Differential sampling)

Patients → Index Test

Negative — X % → Gold Standard

Positive — Y % → Gold Standard

# Bias due to indeterminates

**Evaluator A:**
No indeterminates

Low sensitivity
Low specificity

**Evaluator B:**
Many indeterminates
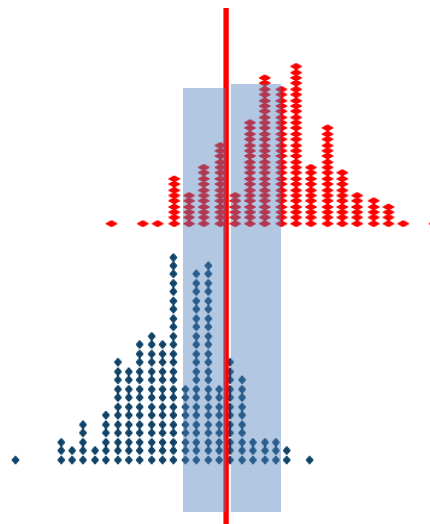
High sensitivity
High specificity

# Indeterminates:
## Where do these values go?

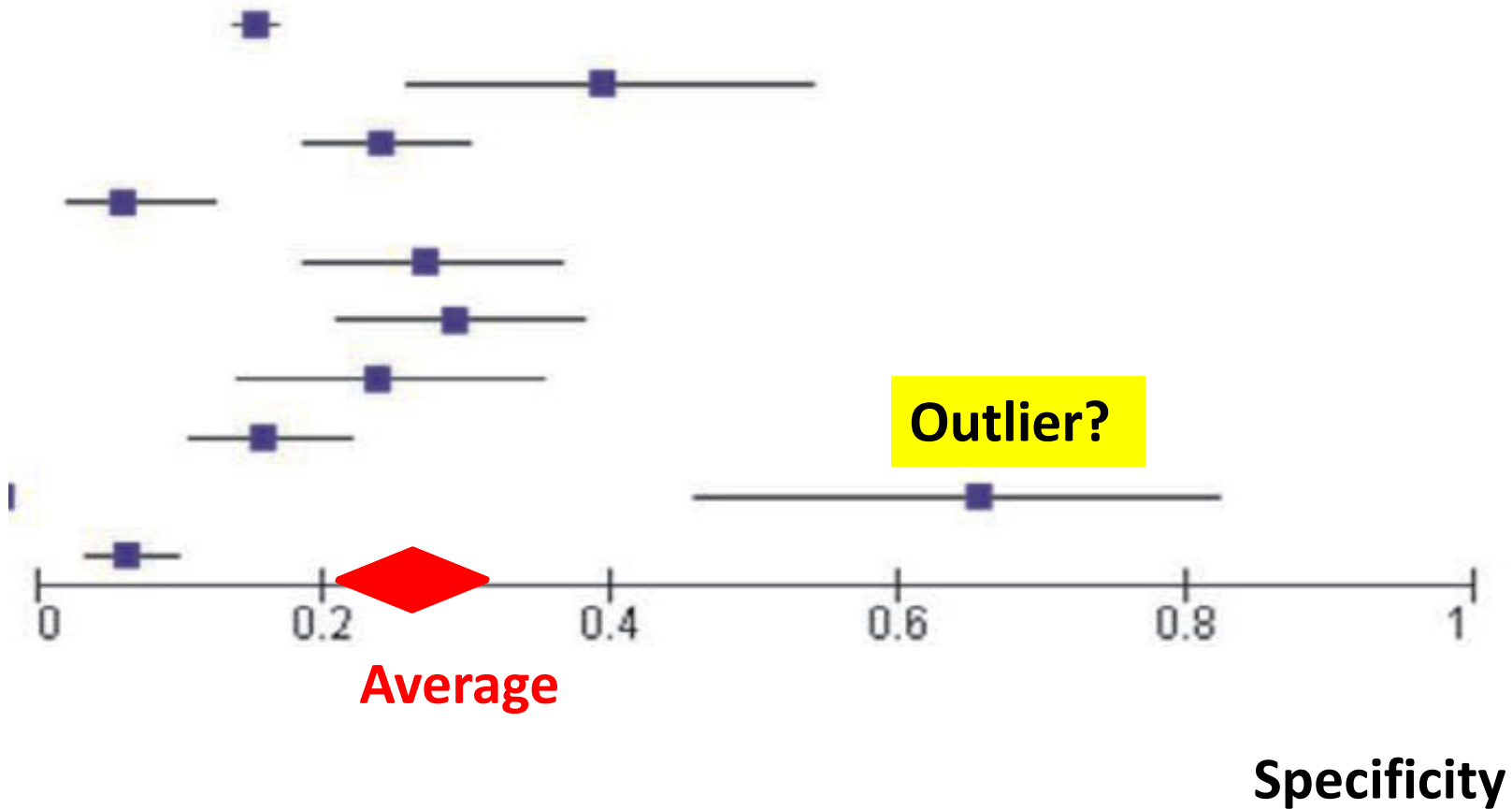| Index Test | Gold Standard | | |
|---|---|---|---|
| | Disease Present | Indeterminate | Disease Absent |
| Positive | | X | |
| Indeterminate | U | Y | V |
| Negative | | Z | |

# Comparing Test Performance

**Why do test results differ?**

1. True differences

2. False differences (bias)

3. **Random variation (imprecision)**

# Understanding Statistical Variation in Studies Meta-Analysis

# Comparing tests

| Source of Difference | Countermeasures |
|---|---|
| **True Differences** | Complete Reporting<br>PICCO<br>Meta-analysis |
| **False Differences**<br>  **bias**<br>  **thresholds** | Improved Study Design<br>ROC Curves |
| **Random variation** | Study design<br>Meta-analysis |

# Higher Levels of Test Evaluation

Societal Impact

Cost effectiveness

Clinical effectiveness

Clinical performance

Analytical performance

# Problems with Test Evaluation

- Potentially useful ≠ Clinically useful
- Potential problems
  - Tests are not used properly
  - Tests do not change diagnosis
  - Tests do not change management
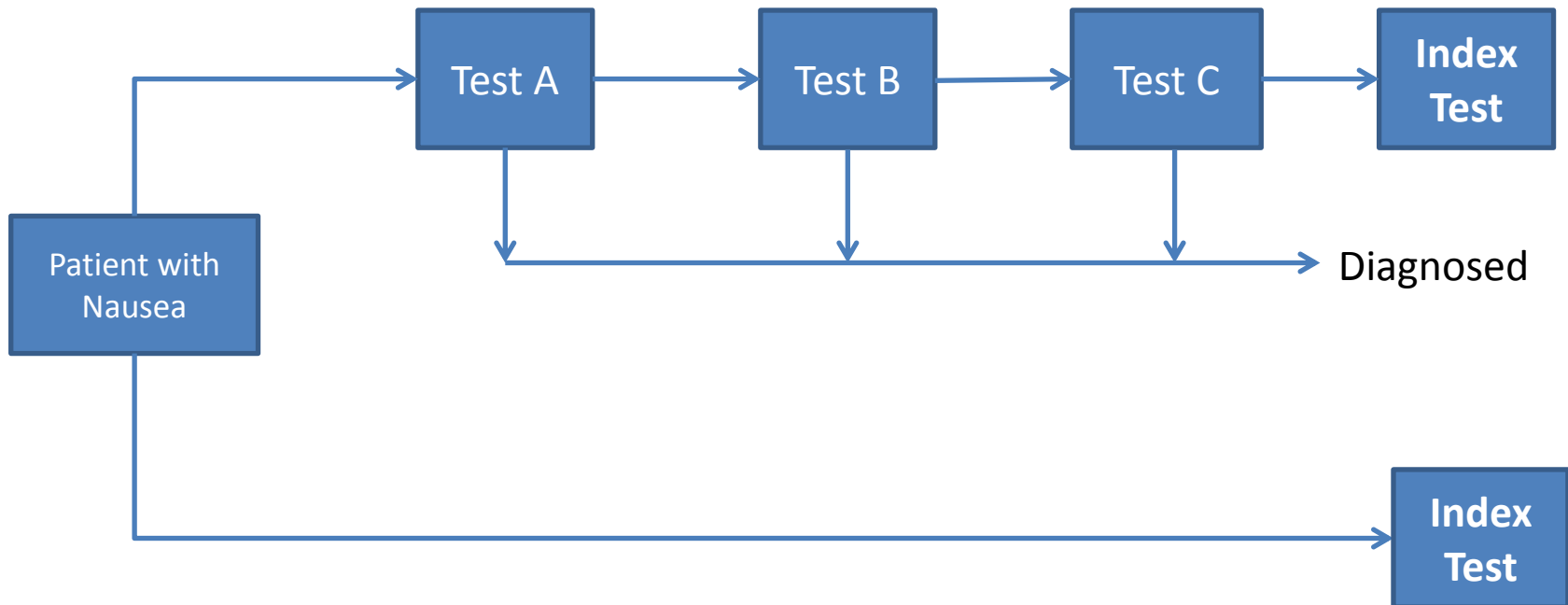- Tests are not used in isolation
  - Incremental value

# Clinical Trial Evaluation of Tests
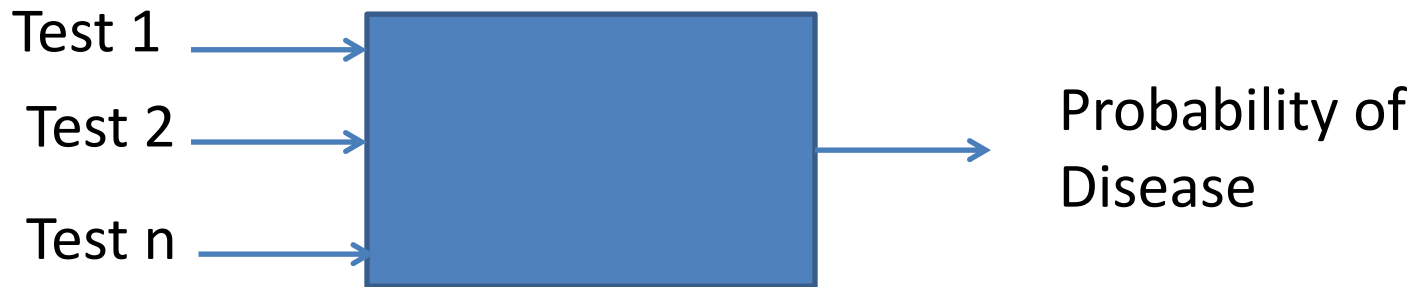
**Key Question:**

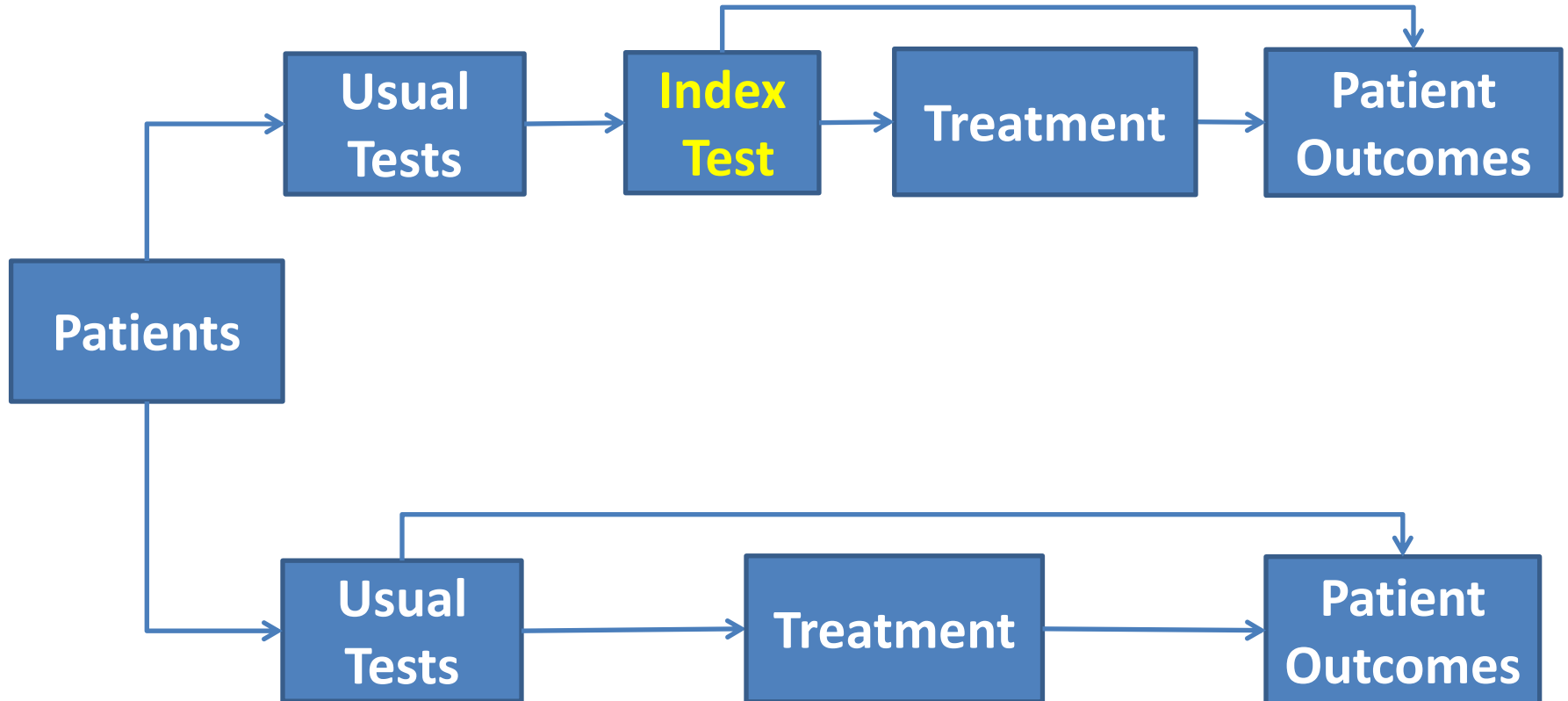*Do patients who receive this test have better outcomes?*

# Tests don't exist in isolation

- Test Research vs Diagnostic Research
- What is the *incremental* value of a test?

# Tests are often combined

Test 1

Test 2

Test n

Probability of Disease

# The acid test

# Clinical Trial Evaluation Prostate Screening (PSA)

| Outcomes | Event Rate per 1000 | | Relative Risk |
|---|---|---|---|
| | No Screen | Screen | |
| All cause mortality | 200 | 198 | 0.99 [0.97-1.01] |
| Death from prostate CA | 8 | 7 | 0.88 [0.71-1.09] |
| Prostate CA diagnosis | 44 | 64 | 1.46 [1.21-1.77] |

# Levels of Evaluation

## Therapeutics

### Phase II/III Trial – Explanatory Trial

Scientific Perspective

Hypothesis: Does this drug affect outcomes?

As-Treated Analysis

Carefully controlled population, setting

Carefully controlled administration and monitoring

### Phase III Trial – Pragmatic Trial

Policy Perspective
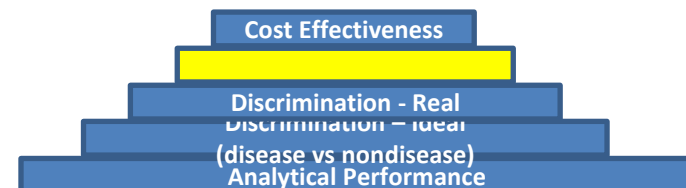
Hypothesis: Does *prescribing* this drug affect outcomes?

Intention-to-Treat Analysis

Patients seeking treatment for condition

Usual conditions

## Diagnostics

### Scientific Test Evaluation

Single test

Idealized population

Expert administration

Expert interpretation

### Pragmatic Test Evaluation

Multiple tests

Actual population

Usual conditions

Cost Effectiveness

Discrimination - Real

Discrimination – Ideal
(disease vs nondisease)

Analytical Performance

# Cost-Effectiveness Plane

Higher Cost
Less
Effective

Higher Cost
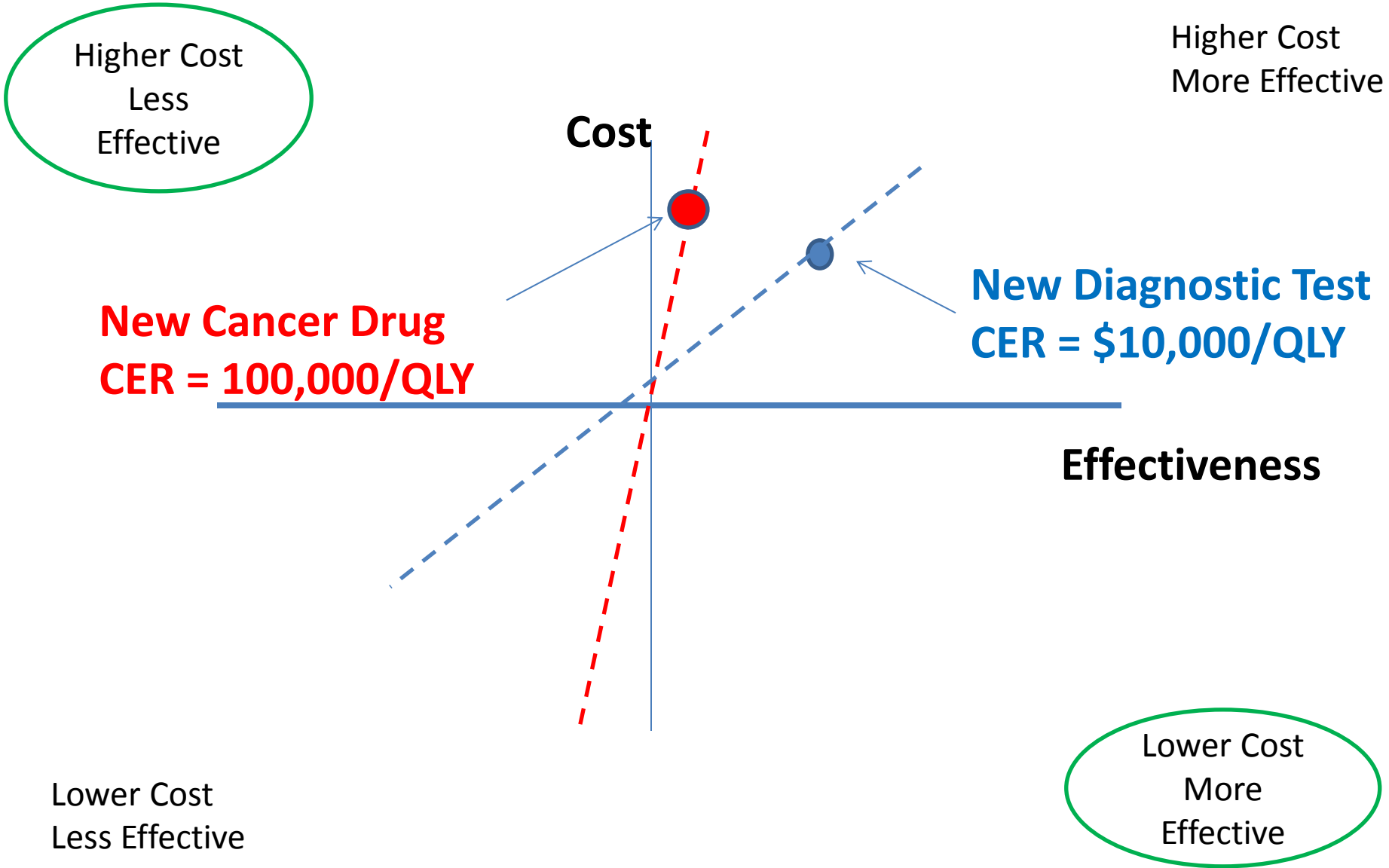More Effective

**Cost**

**New Cancer Drug**
**CER = 100,000/QLY**

**New Diagnostic Test**
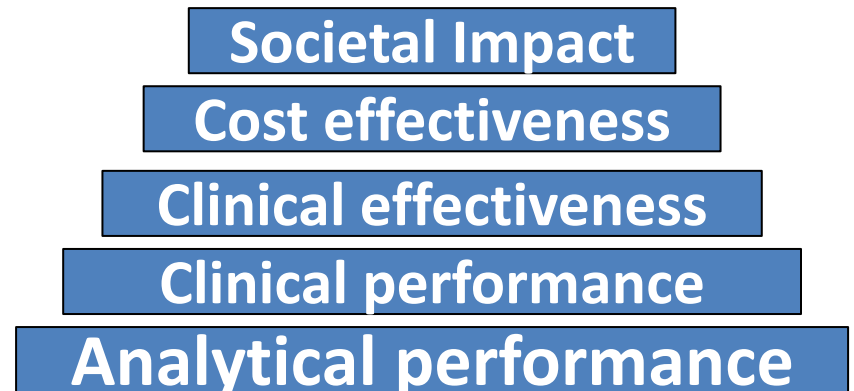**CER = $10,000/QLY**

**Effectiveness**

Lower Cost
Less Effective

Lower Cost
More
Effective

# Summary

- Many ways to assess performance
- Many reasons why studies differ
  - Real differences (PICCO)
  - False differences
    - Thresholds
    - Bias
  - Statistical variation
- Progress in Performance Evaluation
  - Quality of Reporting
  - Quality of studies
  - Types of studies
- Educating Clinicians

**Societal Impact**

**Cost effectiveness**

**Clinical effectiveness**

**Clinical performance**

**Analytical performance**

# Testing A Test:
# Beyond Sensitivity and Specificity