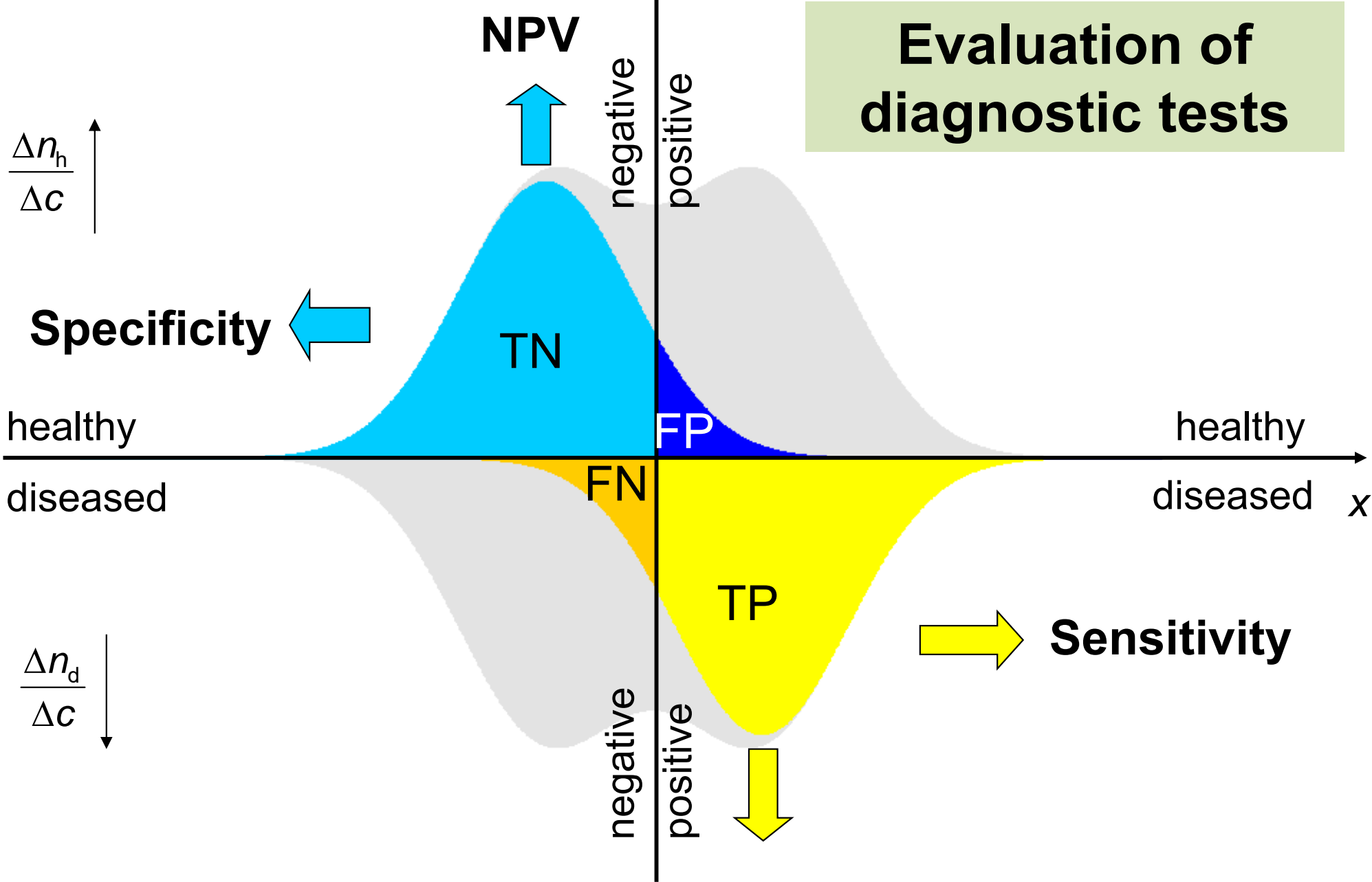
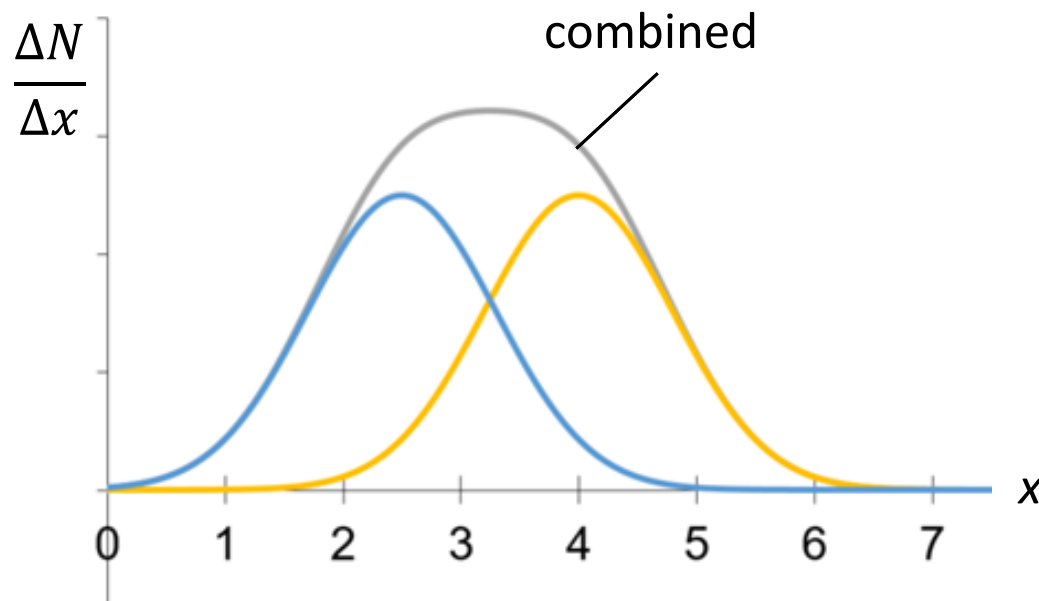
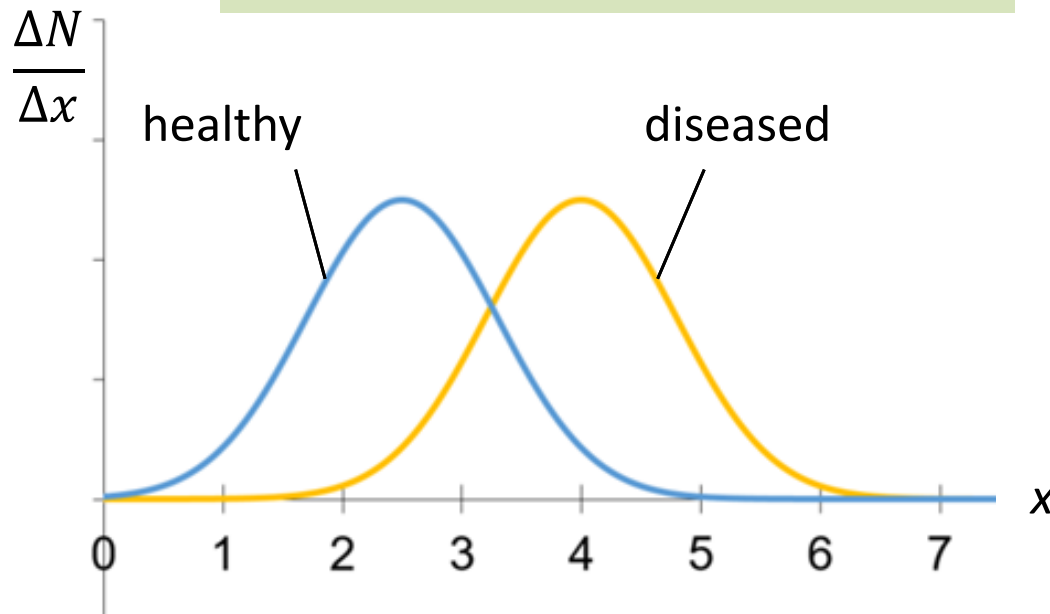


Evaluation of diagnostic tests



Overlapping distributions



everything has a **distribution** (eye color, height, cholesterol level,...)

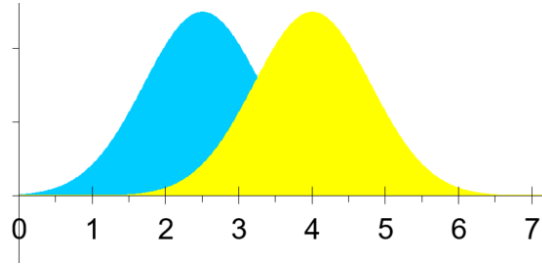
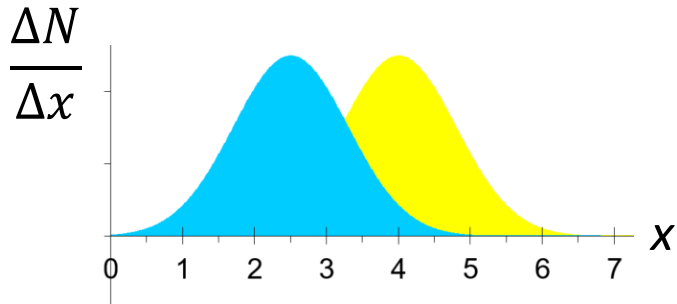
consider a continuous variable with a different distribution in the diseased population than in the healthy one

suppose that the measurable parameter is typically **larger** in the **diseased** population than in the healthy one (if smaller, the reciprocal of the original parameter can be used)

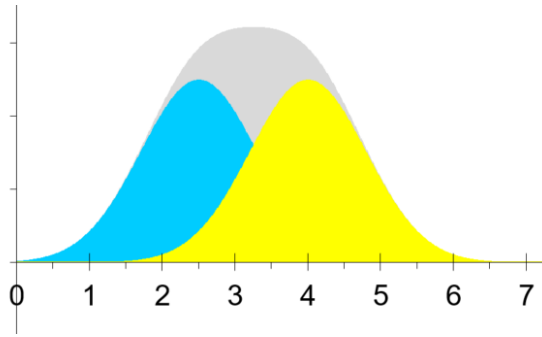
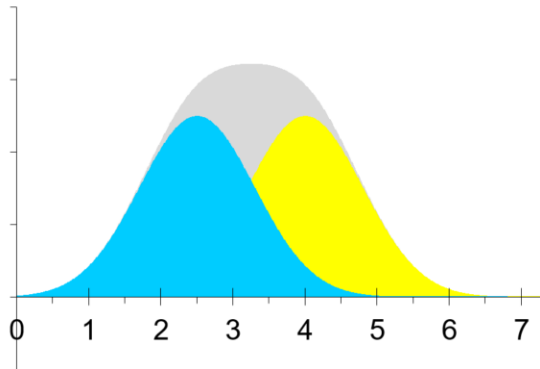
the figure shows two such density functions; the area under the curve corresponds to the number of individuals

in the present example the number of healthy and diseased is the same (the standard deviation of the parameter is the same)

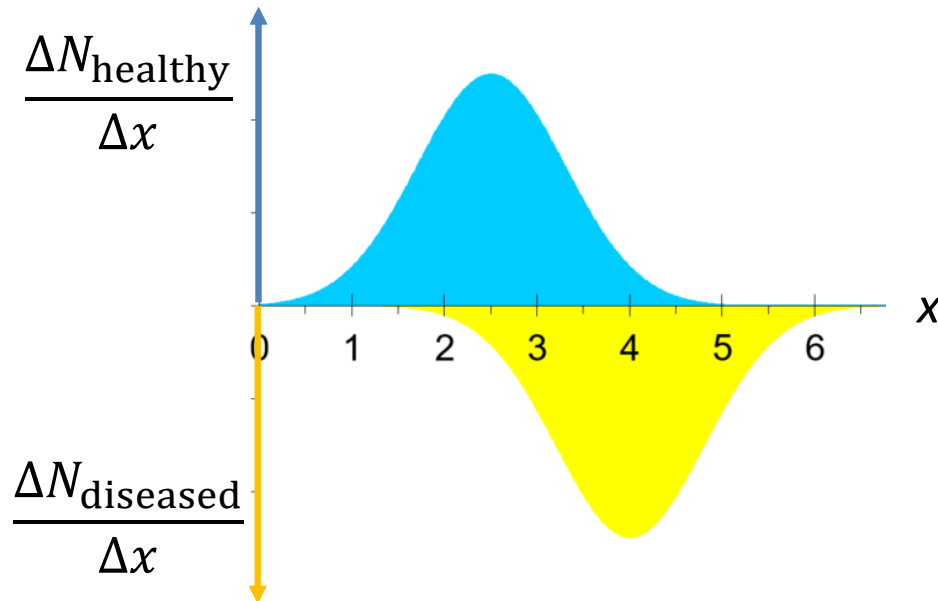
Representations



due to the great importance of the areas under the curve, we prefer an image that colors the areas instead of a line drawing



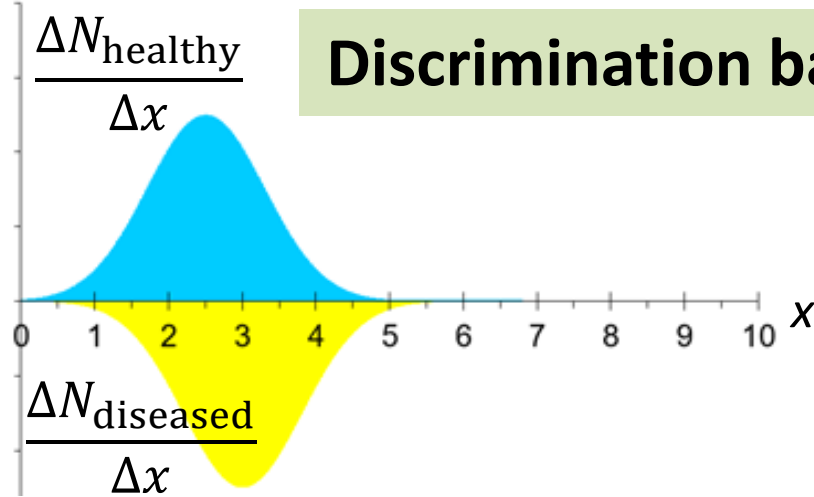
due to the overlaps, correct coloring is difficult (or impossible) in the usual representation



proposed new representation: instead of the negative axis, another positive axis, for the diseased

Discrimination based on overlap magnitude

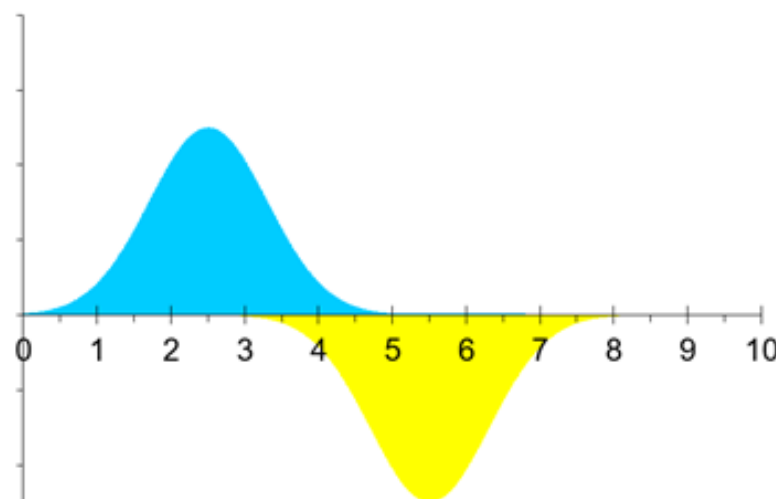
healthy
avg: 2,5
sd: 0,8
diseased
avg: 3,0
sd: 0,8



full
overlap

useless
method

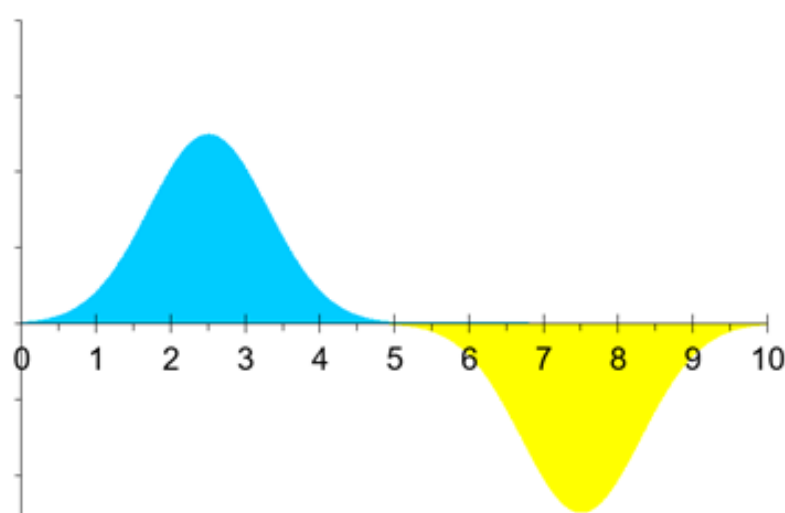
healthy
avg: 2,5
sd: 0,8
diseased
avg: 5,5
sd: 0,8



partial
overlap

(possibly) usable
method

healthy
avg: 2,5
sd: 0,8
diseased
avg: 7,5
sd: 0,8



no
overlap

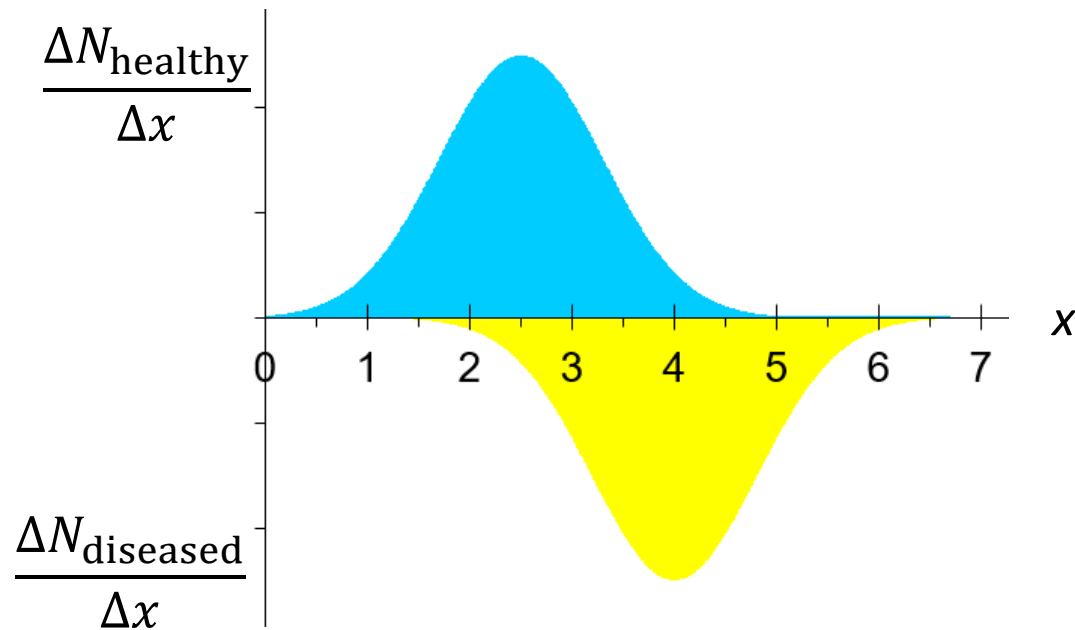
perfect
method

Prevalence

frequency of diseased in
examined population

measure of how common
the disease is

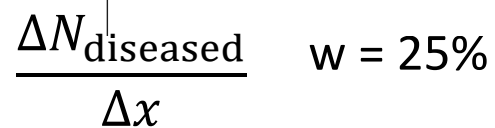
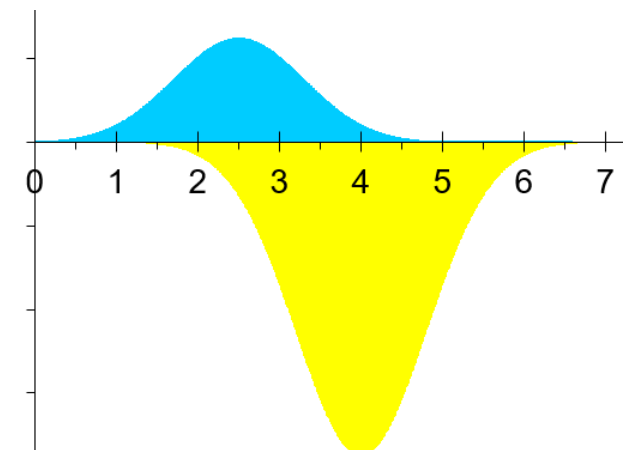
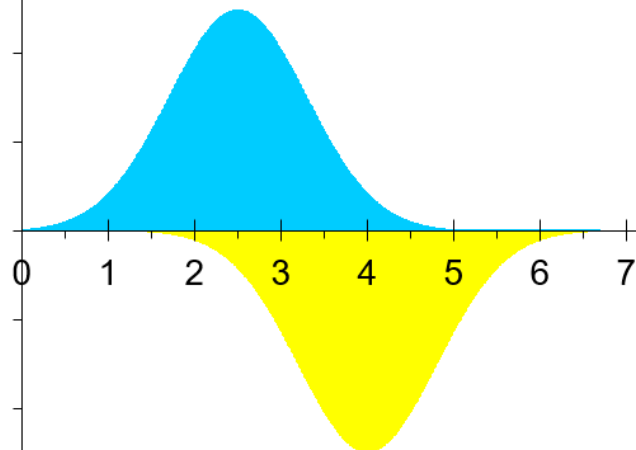
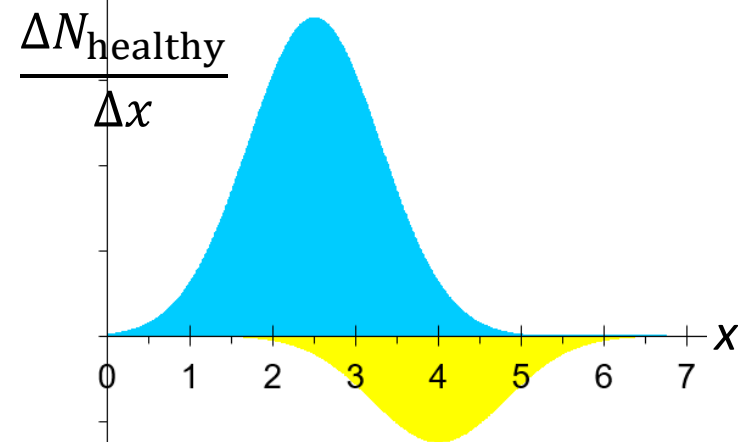
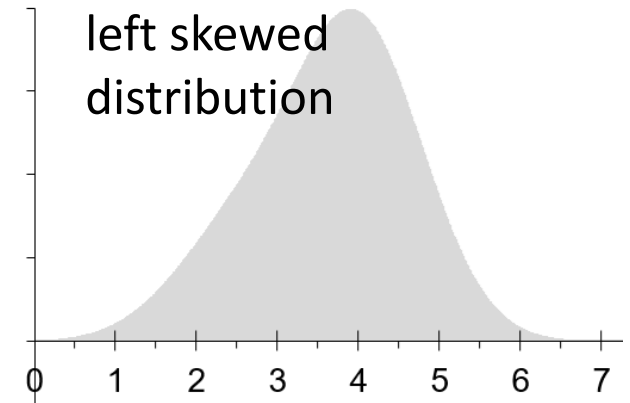
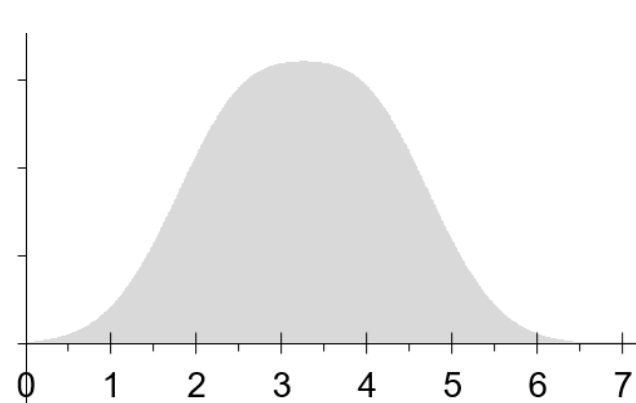
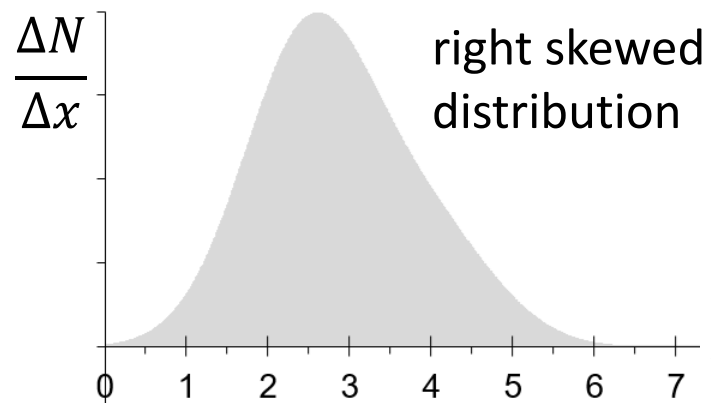
= probability prior to test
= a-priori-probability



$$w = \frac{\text{diseased}}{\text{total}} = \frac{\text{diseased}}{\text{diseased} + \text{healthy}} = \frac{\text{de} - \text{sp}}{\text{se} - \text{sp}}$$

cf: incidence = the number of new cases in a given period and in a given number of population, e.g. 29 per year, per 10 000 people

Effect of prevalence on combined distributions

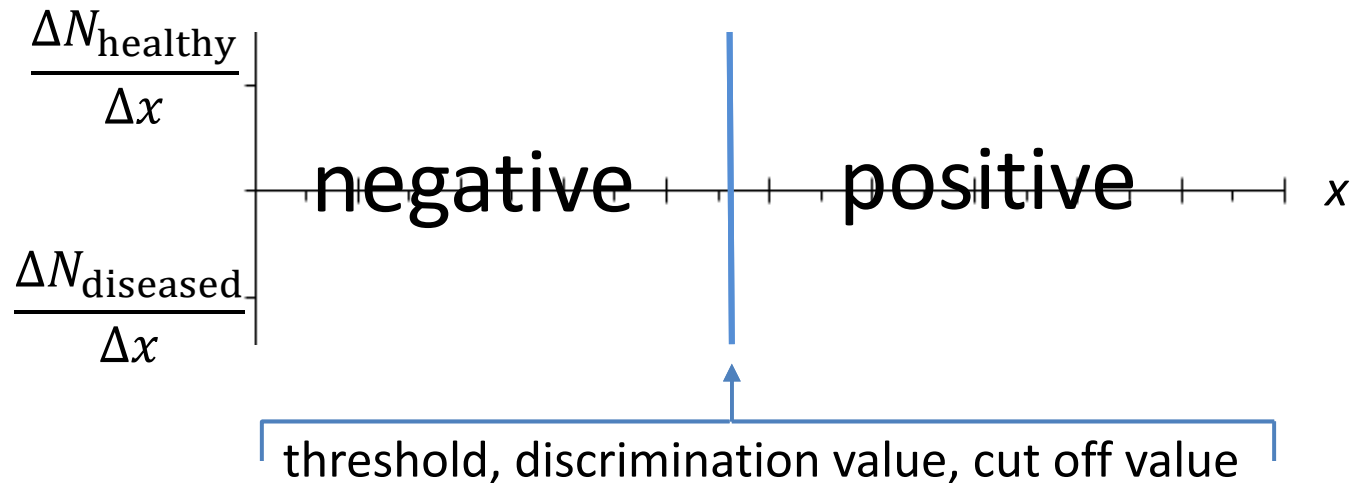


$w = 50\%$

$w = 75\%$

A negative test result below the threshold and a positive test result above it

among the possible measurement parameter values, by designating a **threshold** value, we decide which will be the positive values and which are the negative ones according to the test method



the wish/**desire**/request that the diseased and the positive, respectively healthy and negative match each other as much as possible

however, the **classification** is almost never perfect:

there will be diseased who are positive: true positive, TP ✓

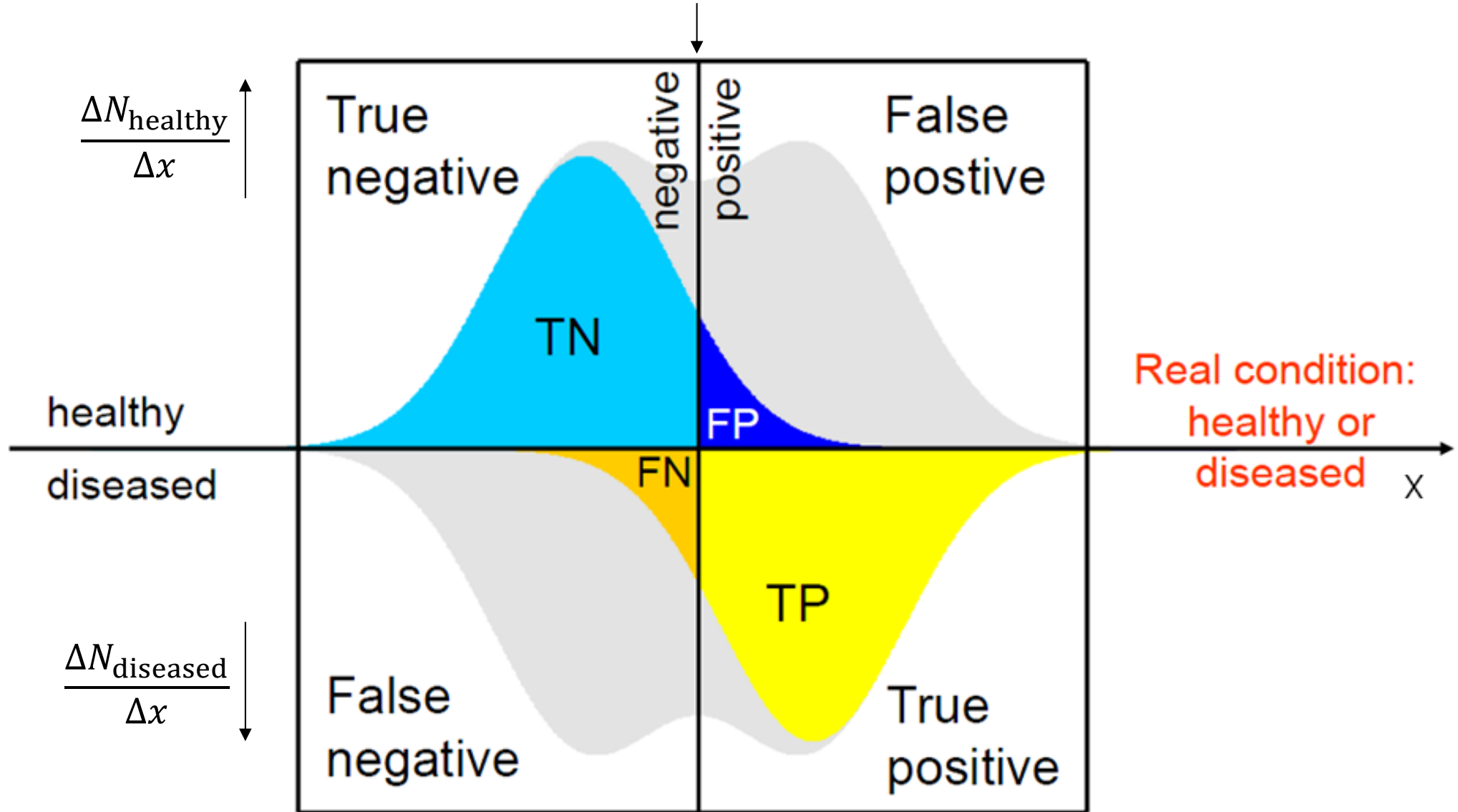
there will be diseased who are negative: false negative, FN 💣

there will be healthy who are negative: true negative, TN ✓

there will be healthy who are positive : false positive, FP 💣

Confusion matrix

cut off value



Test result (prediction):
negative or positive

Parameters of diagnostic „goodness”

based on one (or more) measured parameters diagnostic tests divide the examined into (test) **positive** and (test) **negative** groups

the “**goodness**” of grouping **cannot** be characterized by a single number

(a) how well does it catch those **to be caught**?

e.g. the probability of a COVID infected stating/determining to be positive

(b) how well does it leave those **to be left alone**?

e.g. the probability of claiming to be negative for a person not infected with a COVID

(c) how reliable is a **positive test result**?

in the case of a positive test result, how certain the patient is diseased

e.g. in the case of a positive COVID test, how certain it is that the person is infected with COVID

(d) how reliable the **negative test result** is?

in the case of a negative test result, how certain the person is healthy

e.g. in the case of a negative COVID test, how certain it is that the person is not infected with COVID

The goodness of a test can be described in terms of the following diagnostic parameters

Sensitivity

Specificity

PPV, relevance

NPV, segregation

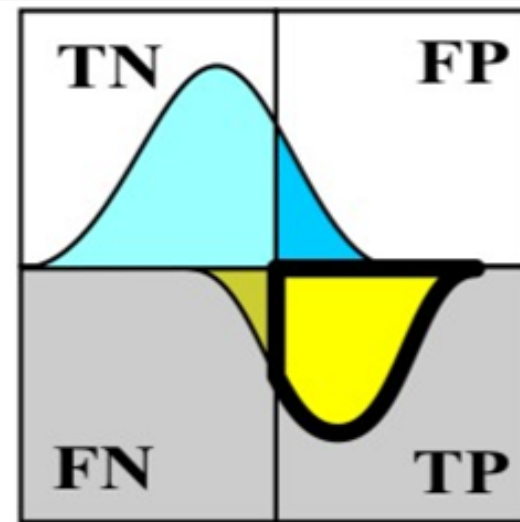
only 3 independent!

Every method must be compared with a reference-method: **gold standard**
method known to always work
(sometimes only the result of an autopsy)



Diagnostic sensitivity

= positive within diseased
 = true positive rate
 = recall rate



probability that the
 test finds the
 diseased positive

$$\frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} = \boxed{\text{se}} = \frac{\text{true positive}}{\text{diseased}} = \boxed{\frac{TP}{TP + FN}} = \underline{\underline{p(\text{positive}|\text{diseased})}}$$

discr. threshold ↓ sens. ↑

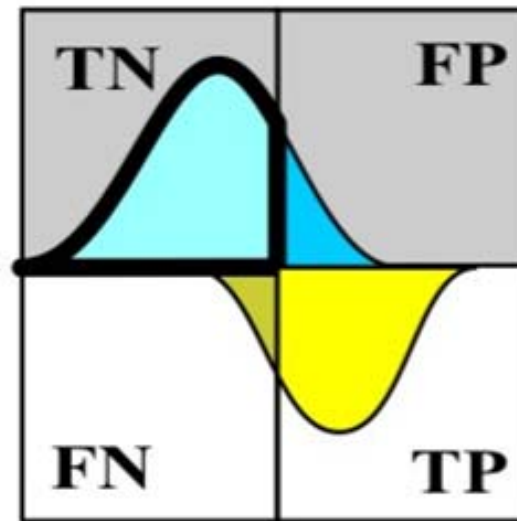
Large-sensitivity tests are required:

In early diagnosis (screening) so that few patients remain unrecognized.
 If the risk of disease is higher than the risk of treatment.

Diagnostic **specificity**

= negative among
healthy

= true negative rate



probability that
the test finds a
healthy negative

$$\frac{\text{True Negative}}{\text{True Negative} + \text{False Positive}} = \boxed{\text{sp}} = \frac{\text{true negative}}{\text{healthy}} = \frac{\boxed{\text{TN}}}{\boxed{\text{TN} + \text{FP}}} = \underline{p(\text{negative}|\text{healthy})}$$

discr. threshold \uparrow spec. \uparrow

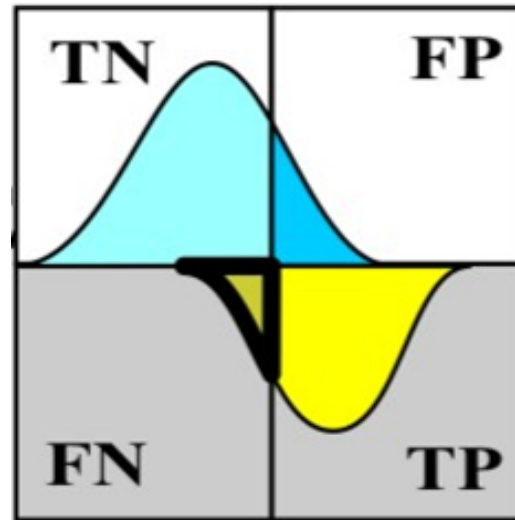
High-specificity tests are important:

When the false positive values have severe consequences (e.g. surgery).

When the risk of treatment is higher than the risk of disease.

Diagnostic False Negative Rate

Type-II error



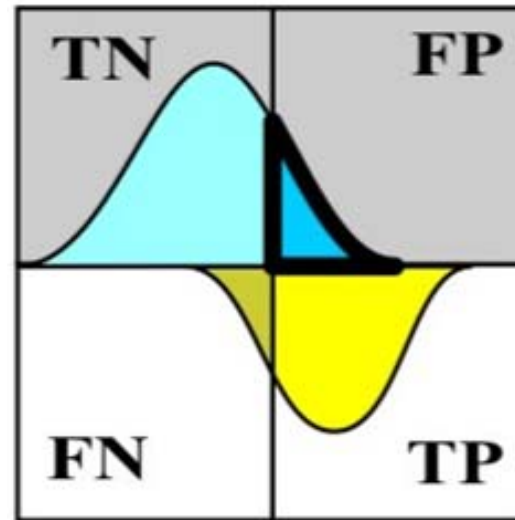
the probability that the
test will find a diseased
negative

negative among diseased

$$\frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{se} = \frac{\text{FN}}{\text{diseased}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = \underline{p(\text{negative}|\text{diseased})}$$

Diagnostic False Positive Rate

Type-I error



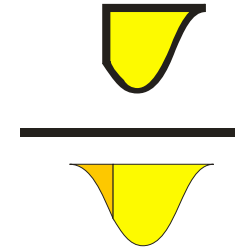
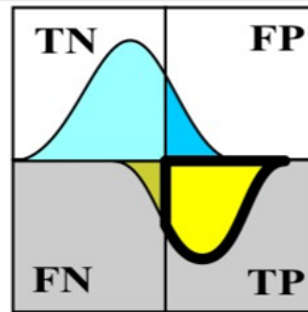
the probability that
the test will find a
healthy positive

positive among the
healthy

$$\frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{sp} = \frac{\text{FP}}{\text{healthy}} = \frac{\text{FP}}{\text{TN} + \text{FP}} = \underline{p(\text{positive}|\text{healthy})}$$

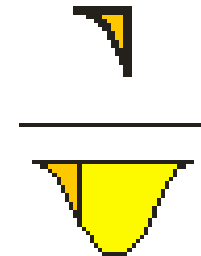
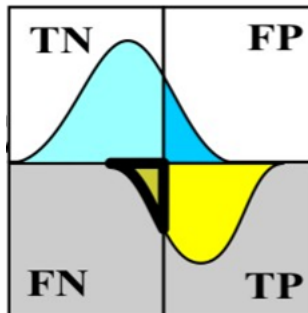
Horizontal rates are independent of prevalence

sensitivity
(se)



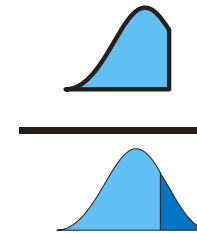
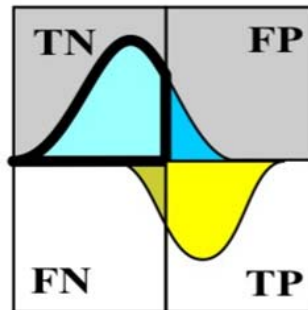
$$se = \frac{TP}{TP + FN}$$

false negative rate
(1-se)



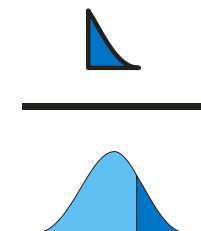
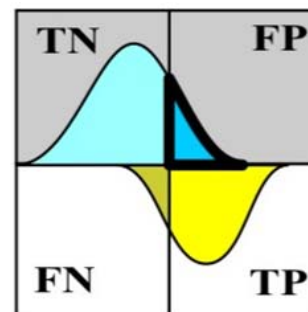
$$1 - se = \frac{FN}{FN + TP}$$

specificity
(sp)



$$sp = \frac{TN}{TN + FP}$$

false positive rate
(1-sp)



$$1 - sp = \frac{FP}{TN + FP}$$

Predictive values (vertical rates)

a-posteriori-probabilities; they depend strongly on prevalence

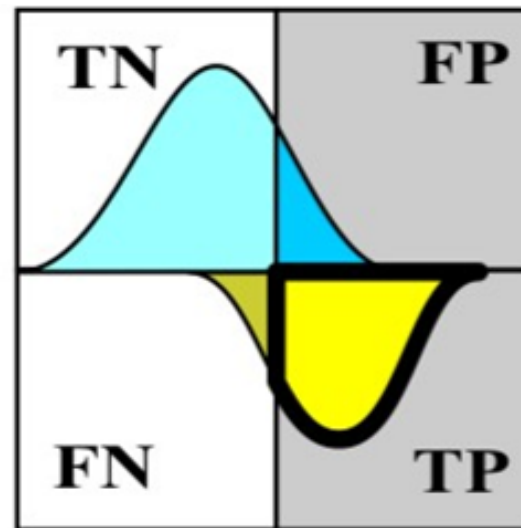
Positive predictive value

= **PPV**

= predictive value positive

= PVP

= **diagnostic relevance**



probability of
disease if test is
positive

diseased among
positives

$$\frac{\text{TP}}{\text{TP} + \text{FP}} = \text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{\text{se} \cdot w}{\text{se} \cdot w + (1 - \text{sp}) \cdot (1 - w)} = \underline{\underline{p(\text{diseased}|\text{positive})}}$$

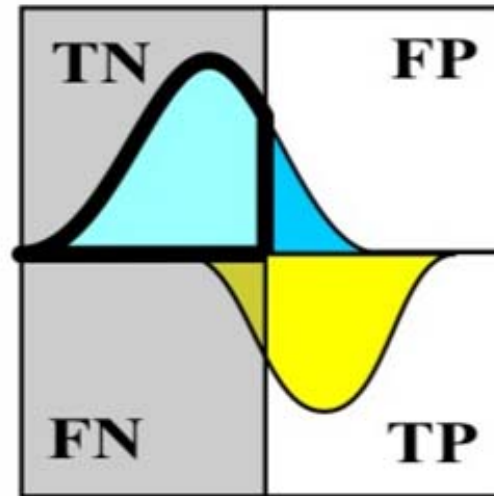
Negative predictive value

= NPV

= predictive value negative

= PVN

= diagnostic **segregation**



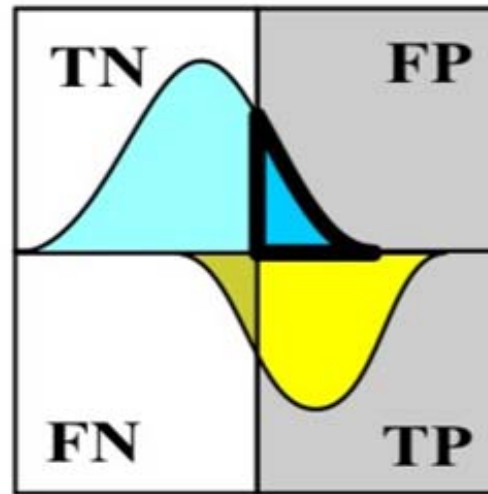
probability of health
if test is negative

healthy among
negatives

$$\frac{\text{Area under curve to the left of threshold}}{\text{Total area under curve}} = \boxed{\text{NPV}} = \frac{\text{TN}}{\text{negative}} = \frac{\boxed{\text{TN}}}{\text{TN} + \text{FN}} = \frac{\text{sp} \cdot (1 - w)}{\text{sp} \cdot (1 - w) + (1 - \text{se}) \cdot w} = \underline{\underline{p(\text{healthy}|\text{negative})}}$$

False alarm rate

=1-PPV



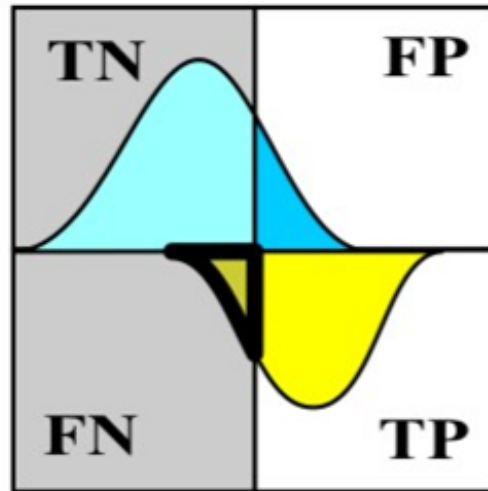
the probability of the
absence of the disease
if the test is positive

healthy among
positives

$$\frac{\text{FP}}{\text{FP} + \text{TP}} = 1 - \text{PPV} = \frac{\text{FP}}{\text{positive}} = \frac{\text{FP}}{\text{FP} + \text{TP}} = \underline{p(\text{healthy}|\text{positive})}$$

False reassurance rate

=1-NPV



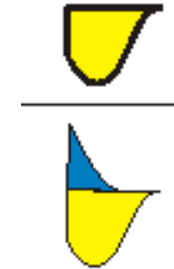
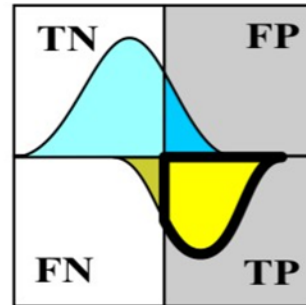
the probability of the
presence of the disease
if the test is negative

diseased among
negatives

$$\frac{\text{FN}}{\text{FN} + \text{TN}} = 1 - \text{NPV} = \frac{\text{FN}}{\text{negative}} = \frac{\text{FN}}{\text{FN} + \text{TN}} = \underline{p(\text{diseased}|\text{negative})}$$

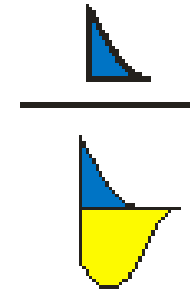
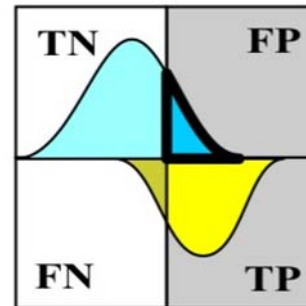
Vertical rates are dependent of prevalence

positive predictive
value
(PPV)



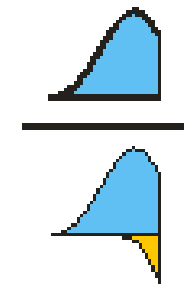
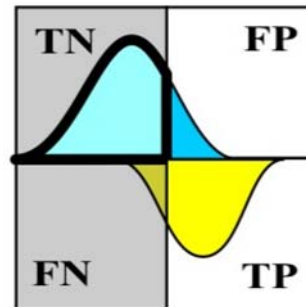
$$PPV = \frac{TP}{FP + TP}$$

false alarm rate
(1-PPV)



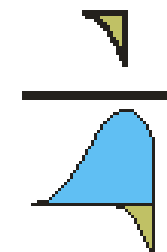
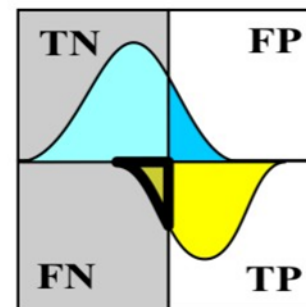
$$1 - PPV = \frac{FP}{FP + TP}$$

negative
predictive value
(NPV)



$$NPV = \frac{TN}{TN + FN}$$

false reassurance rate
(1-NPV)



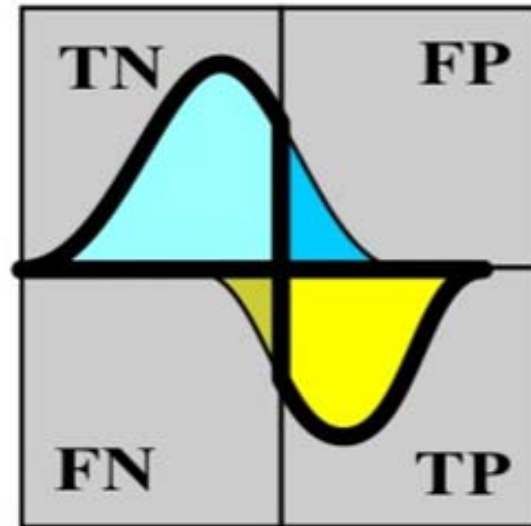
$$1 - NPV = \frac{FN}{TN + FN}$$

Diagnostic accuracy

= da=de

= efficacy/efficiency

= correct classification rate



probability of correct diagnosis

$$\frac{\text{Area under curve to the right of threshold}}{\text{Total area under curve}} = \boxed{\text{de}} = \frac{\text{TP} + \text{TN}}{\text{total}} = \boxed{\frac{\text{TP} + \text{TN}}{\text{TN} + \text{FP} + \text{FN} + \text{TP}}} = \underline{\text{se} \cdot w + \text{sp} \cdot (1 - w)}$$

often: discrimination threshold is chosen so that accuracy is maximized

Effect of prevalence

case1: $w = 50\%$

NPV = 94.7%

sp = 90%

Gold-standard		Test	
		negative	positive
	healthy	90	10
	diseased	5	95

se = 95%

(de = 92.5%)

PPV = 90.5%

NPV = 99.4%

Case 2: $w = 10\%$

		Test	
		negative	positive
Gold-standard	healthy	810	90
	diseased	5	95

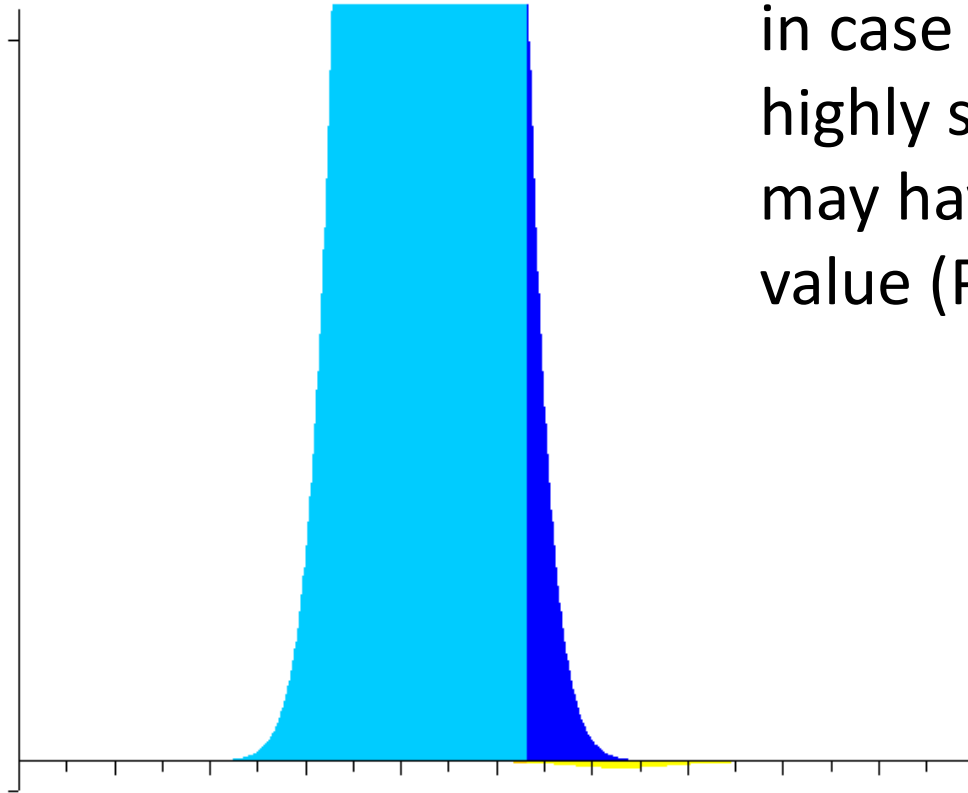
sp = 90%

se = 95%

(de = 90.5%)

PPV = 51.4%

in case of very small prevalence a highly sensitive and specific test may have low positive predictive value (PPV)



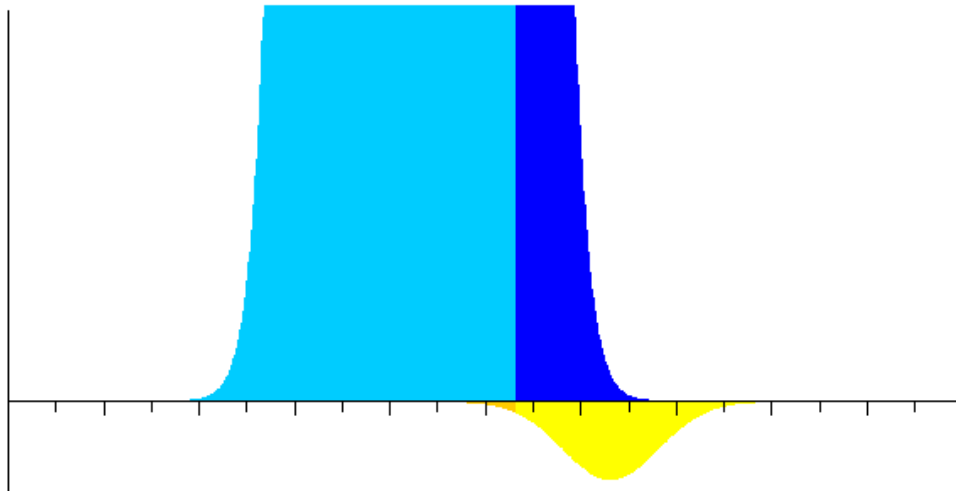
prevalence = 0.1 %

sensitivity = 98 %

specificity = 98 %



PPV = 4 %



A patient comes to your office frantic over the results of a home HIV test. The test touts 99% sensitivity and 99% specificity. On questioning, you determine that this patient is at low risk for HIV; given your assessment of his risk factors, you believe he comes from a population group that has a baseline prevalence of HIV of 1 in 100,000. He now presents to you with a positive result on his home HIV test. Given his baseline risk and the positive home test, what are the chances that this patient is actually HIV positive?

$$(PPV = 9.89 \times 10^{-4} \cong 0.001)$$

Stuart Spitalnic: Test properties I: Sensitivity, specificity, and predictive values; Hospital Physician, September 2004, 27-31

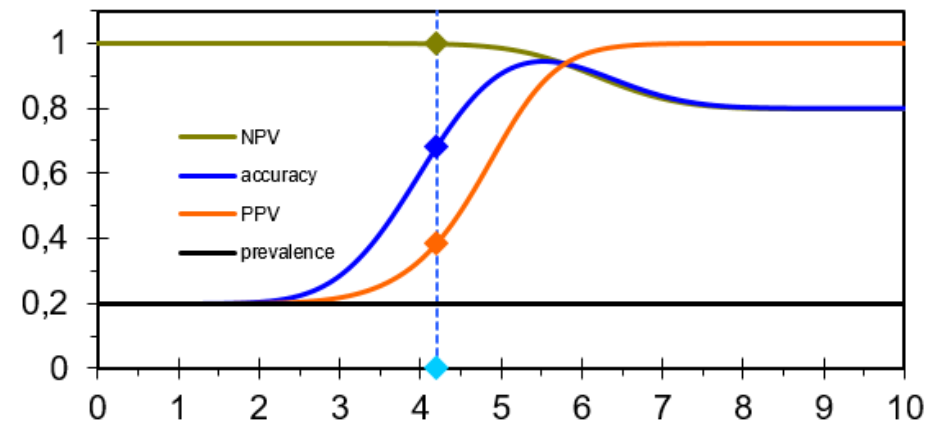
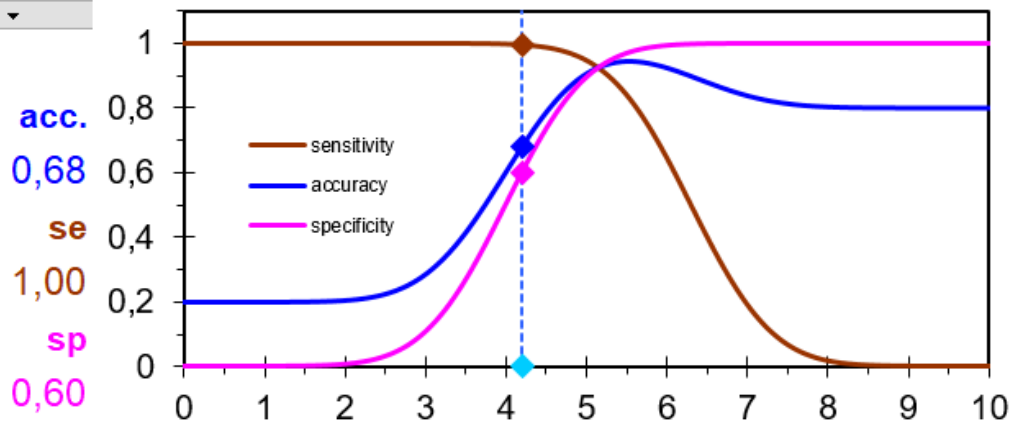
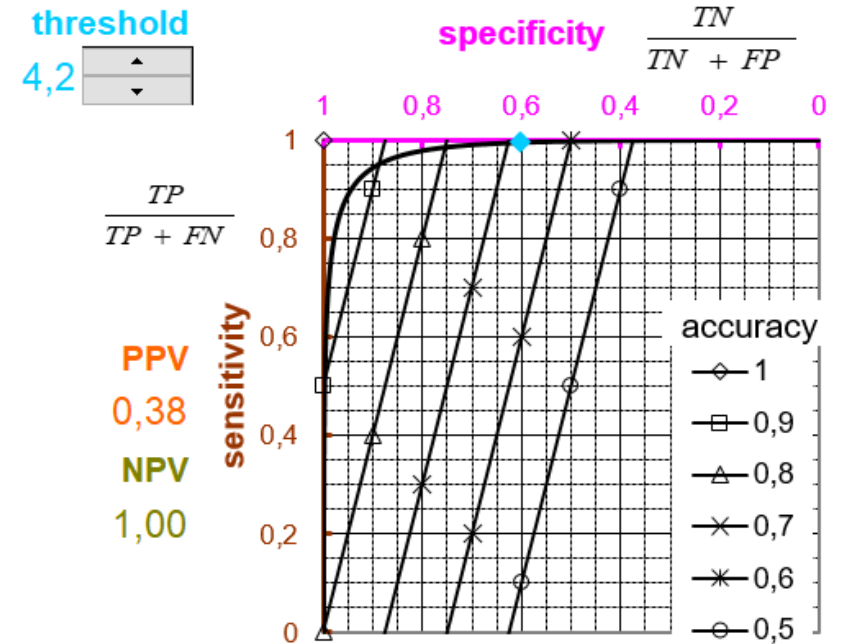
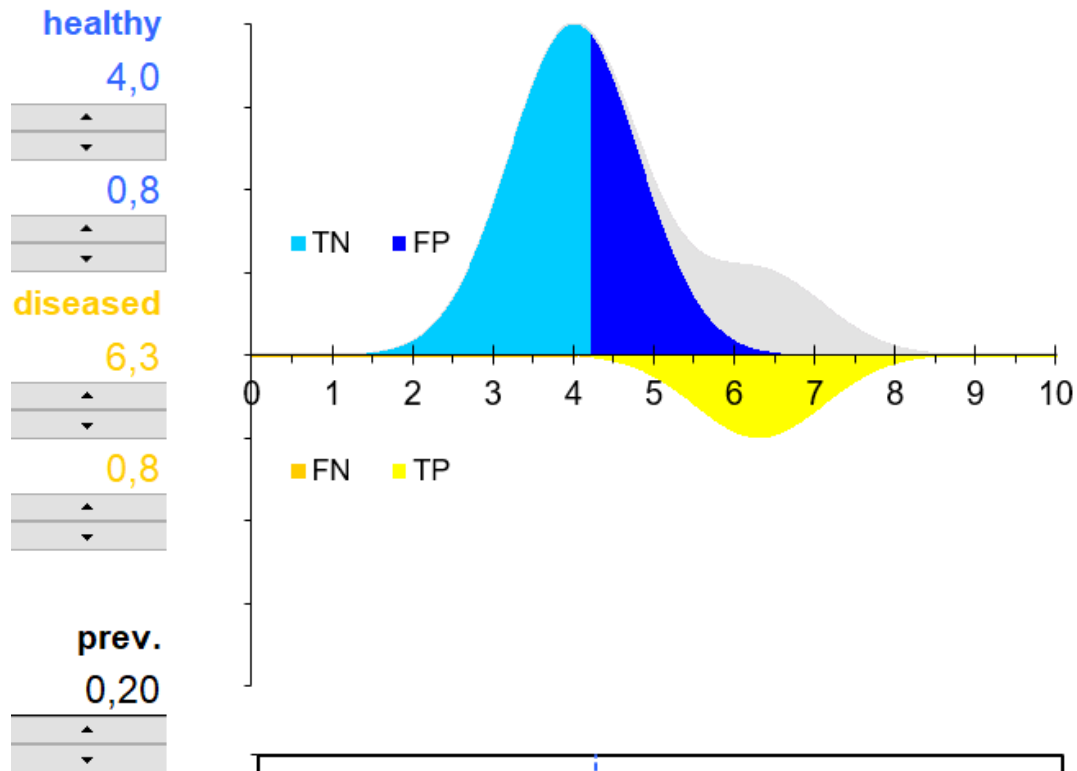
Overview

Sensitivity	se	$\frac{TP}{TP + FN}$	$p(P D)$	positive within diseased	True Positive Rate	prevalence-independent
Specificity	sp	$\frac{TN}{TN + FP}$	$p(N H)$	negative among healthy	True Negative Rate	
False Negative Rate	1-se	$\frac{FN}{TP + FN}$	$p(N D)$	negative among diseased		
False Positive Rate	1-sp	$\frac{FP}{TN + FP}$	$p(P H)$	positive among the healthy		
Positive Predictive Value	PPV	$\frac{TP}{TP + FP}$	$p(D P)$	diseased among positives	Relevance	prevalence-dependent
Negative Predictive Value	NPV	$\frac{TN}{TN + FN}$	$p(H N)$	healthy among negatives	Segregation	
False alarm rate	1-PPV	$\frac{FP}{TP + FP}$	$p(H P)$	healthy among positives		
False reassurance rate	1-NPV	$\frac{FN}{TN + FN}$	$p(D N)$	diseased among negatives		

conditional probability (Bayes)

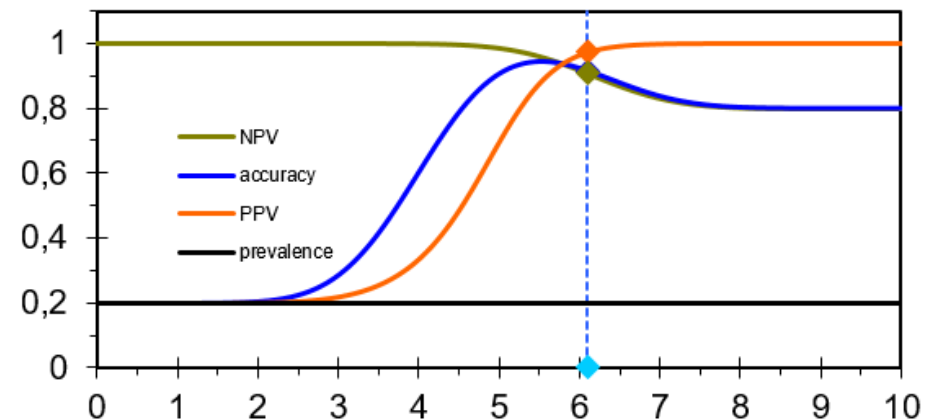
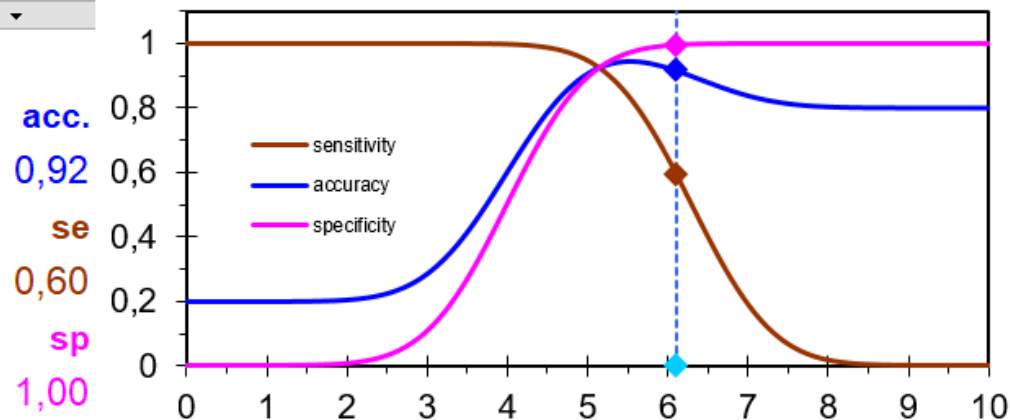
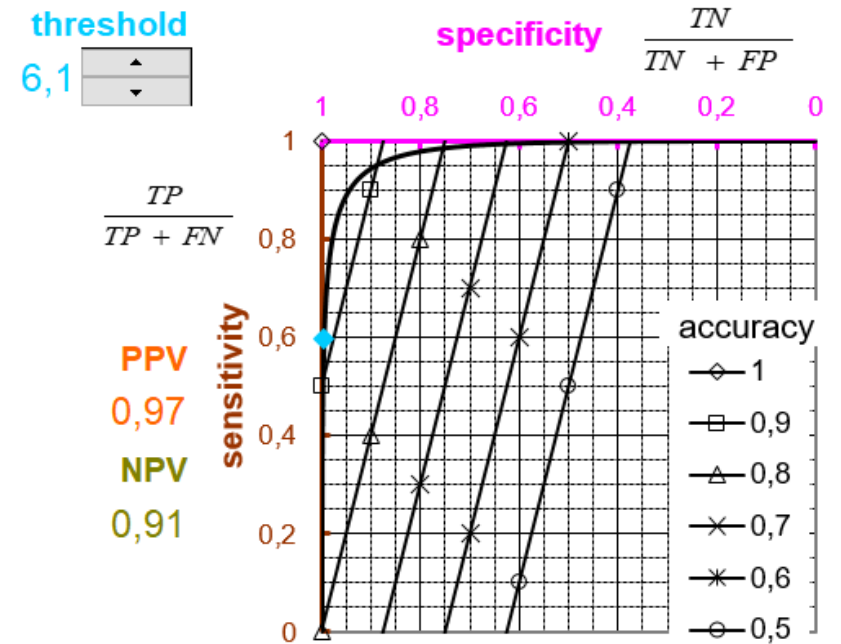
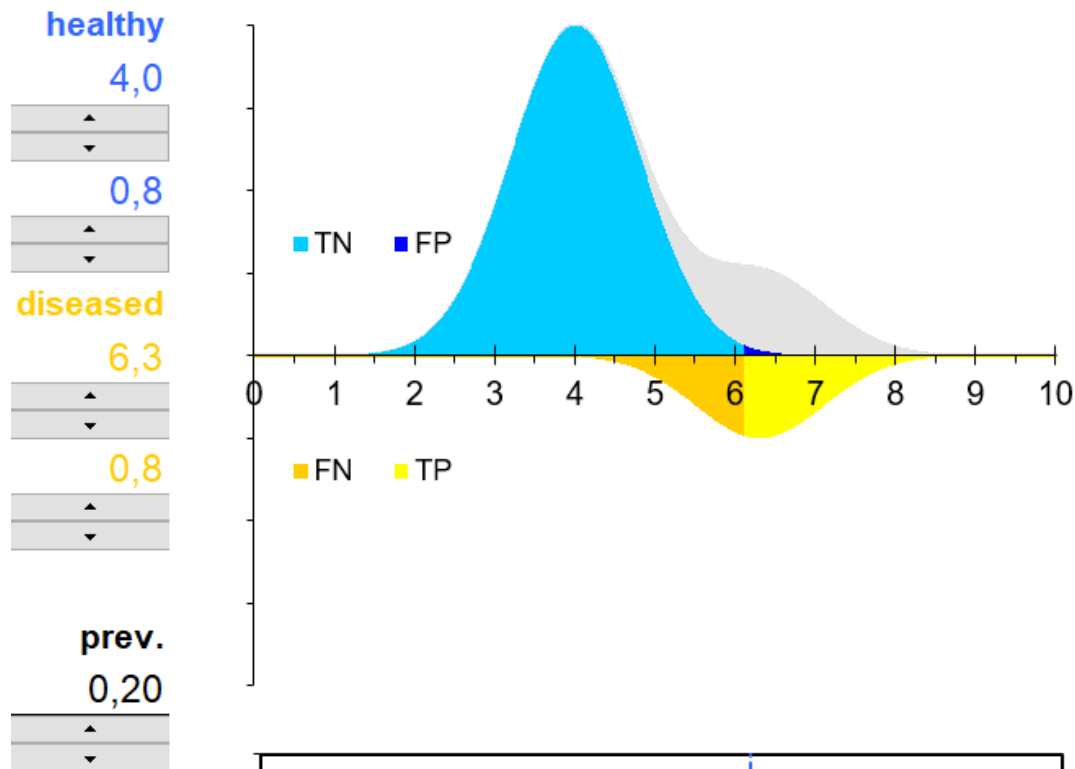
how well does it catch
those to be caught?

Maximize diagnostic sensitivity



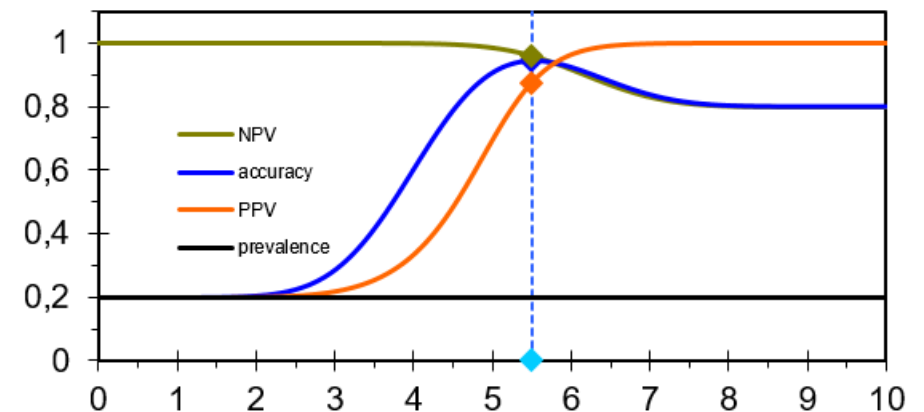
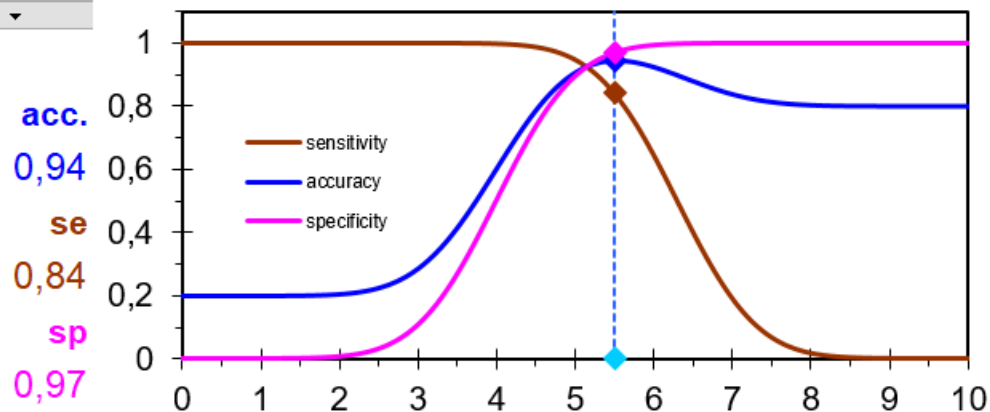
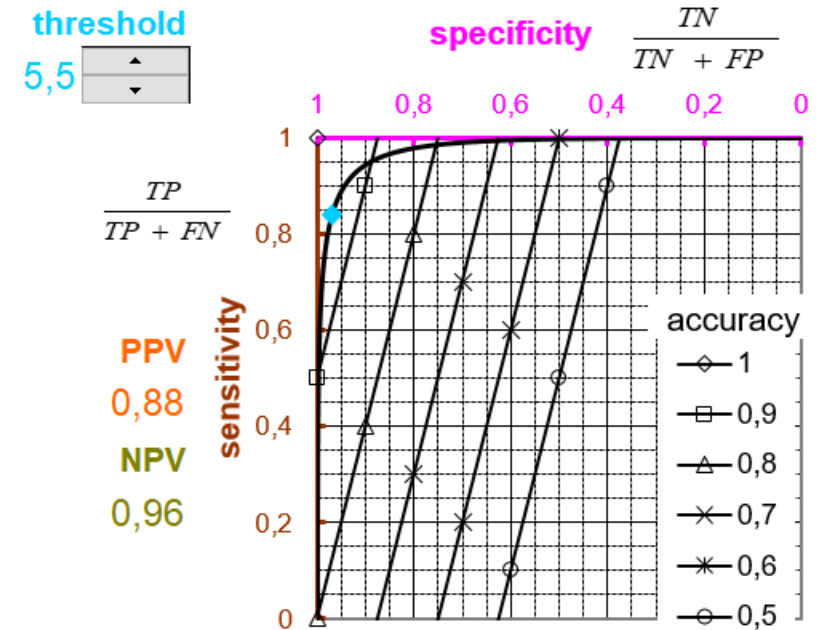
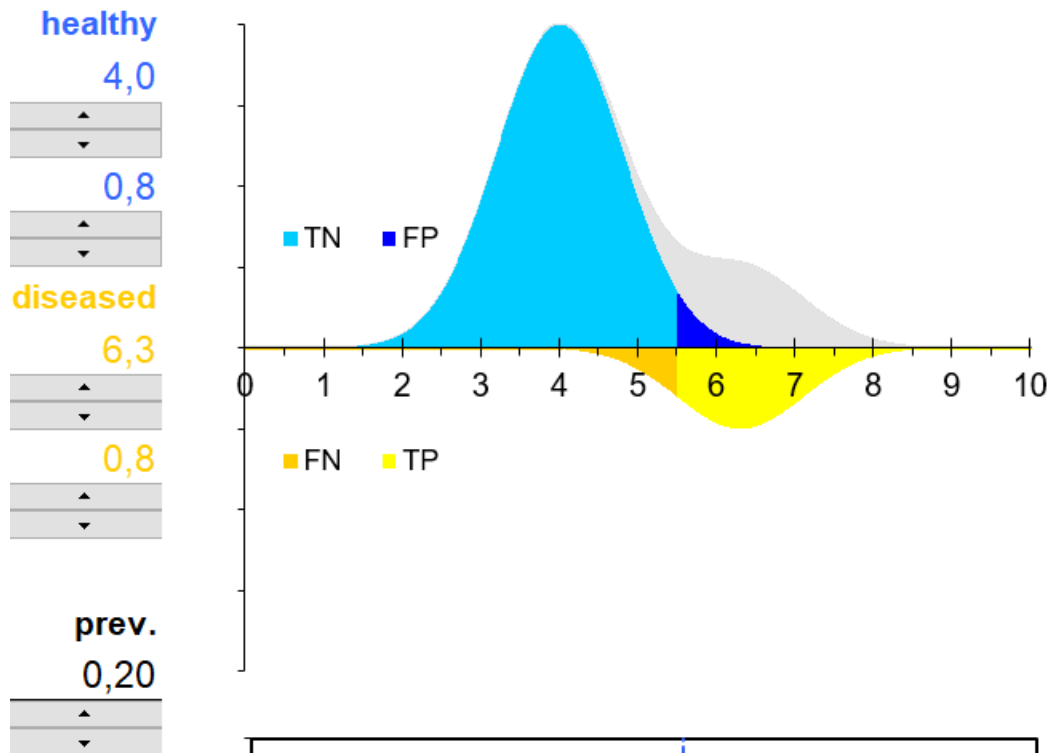
how well does it leave
those **to be left alone?**

Maximize diagnostic specificity



it is equally important to catch those **to be caught** and to leave those **to be left alone**?

Maximize diagnostic accuracy



Take-home message

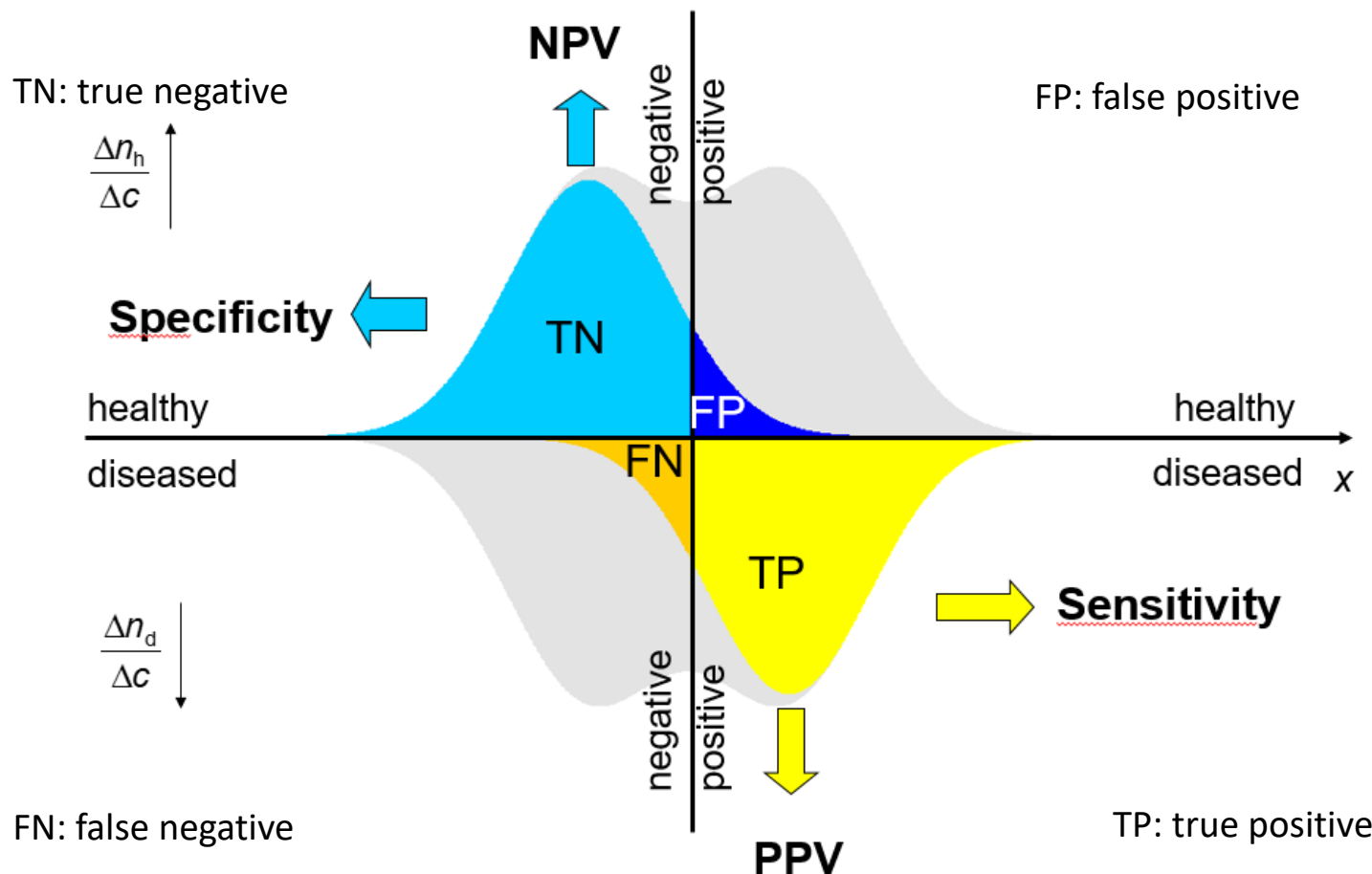
everything has a **distribution**; the distribution of sick and healthy values **overlap**

whether it is possible to decide which is more important :

to detect the disease in as many patients as possible in order to receive treatment (**maximizing sensitivity**), or

to assume a false positive value (minimizing false-positive ratio or **maximizing specificity**) in as few healthy people as possible so that they do not receive unnecessary therapy

if they cannot be decided, they are equally important: **maximizing accuracy**



diagnostic accuracy:

$$de = \frac{TP + TN}{\text{total}} = \frac{TP + TN}{TN + FP + FN + TP}$$

in case of several possible methods:
ROC

