

Biostatistics

Week 7

➤ **Diagnostic tests as “patient classifier”**

- How can we describe the quality of a diagnostic test with binary outcome:

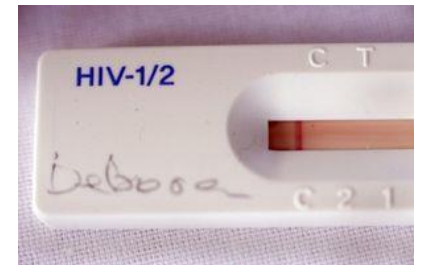
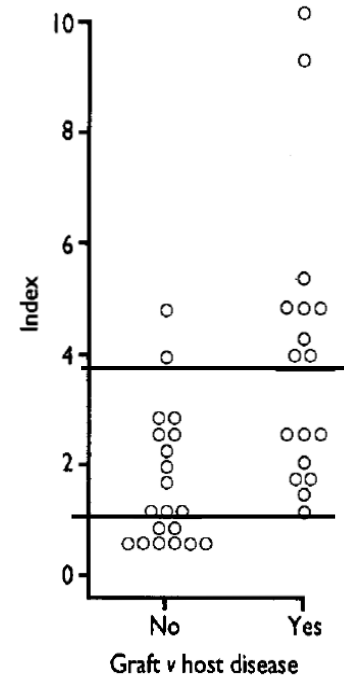
→ Sensitivity, Specificity

- How can we describe the predictive value of a binary diagnostic test:

→ PPV, NPV or positive and negative predictive value

- How to evaluate a diagnostic test with continuous score outcome:

→ ROC curve analysis and its AUC



How to quantify the performance of a test?

1. Performance characteristics of a diagnostic test in a lab setting

Sensitivity

Specificity

Choice of a threshold

2. Performance of a diagnostic test in a population application

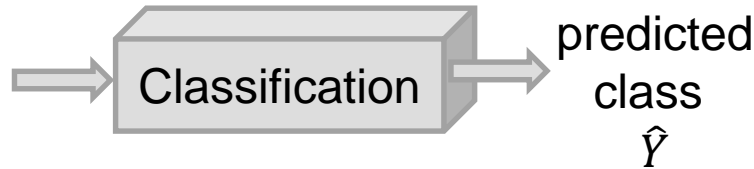
Positive predictive value of a test (PPV)

Negative predictive value of a test (NPV)

Impact of disease prevalence, sensitivity and specificity on predictive values

Binary test ore binary classification rule

Explanatory
variable \mathbf{X}
(e.g.blood sample)



Target Variable Y

2 classes:

Positive or **Negative**

1 or **0**

Yes or **No**

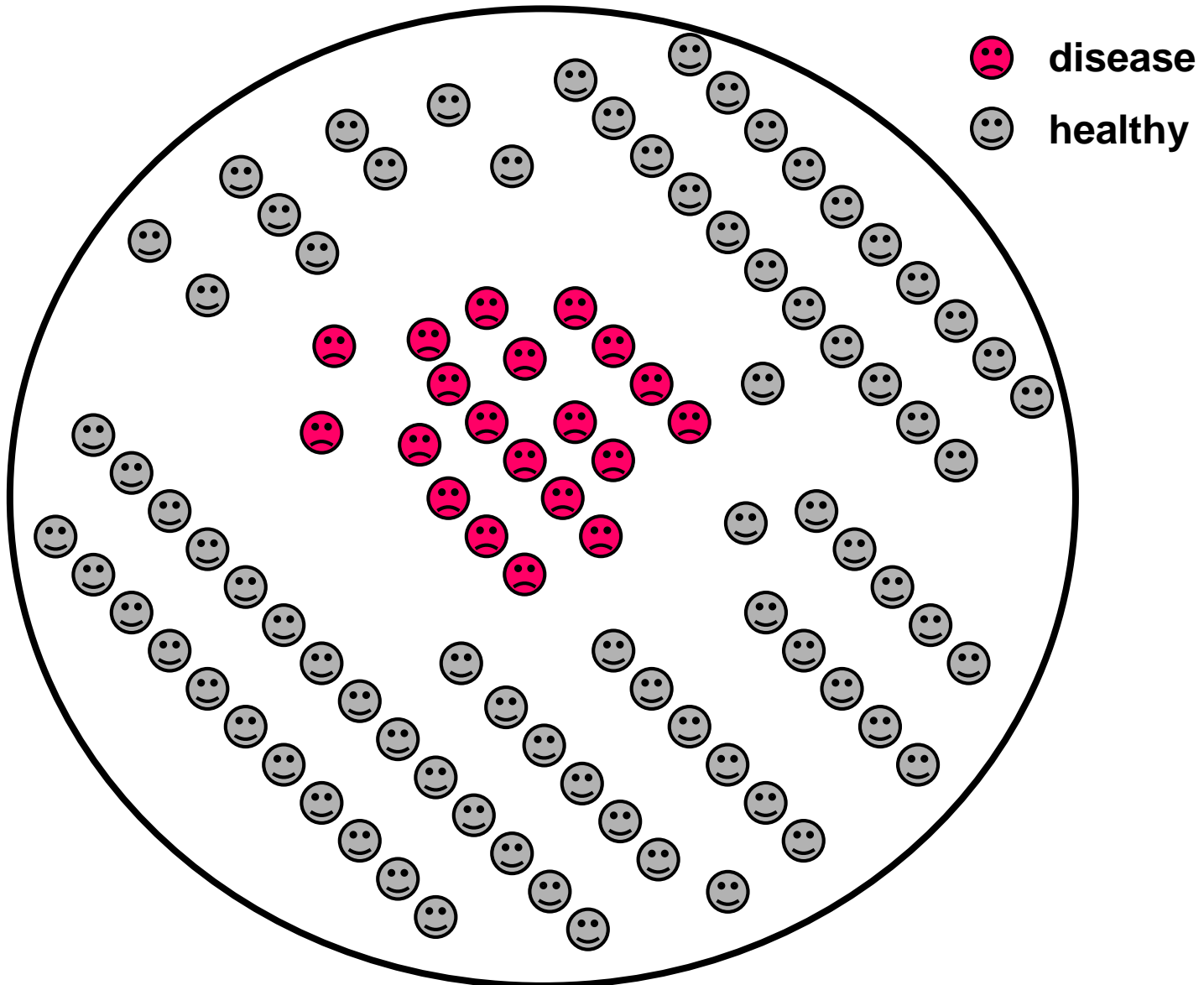
Diseased or **Healthy**

Each observation unit described by input \mathbf{x} , belongs to one of two classes.

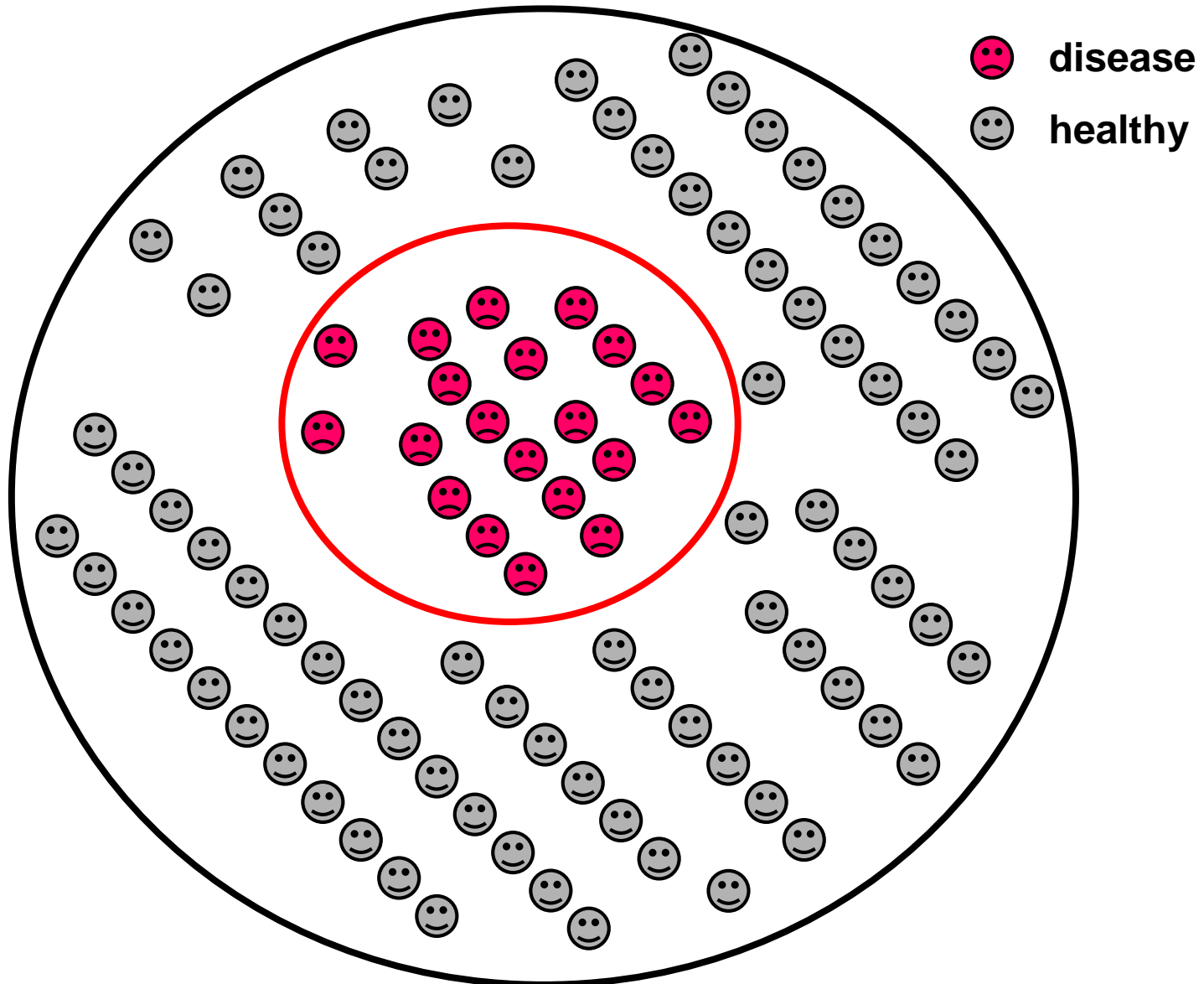
Y : true class

\hat{Y} : predicted class

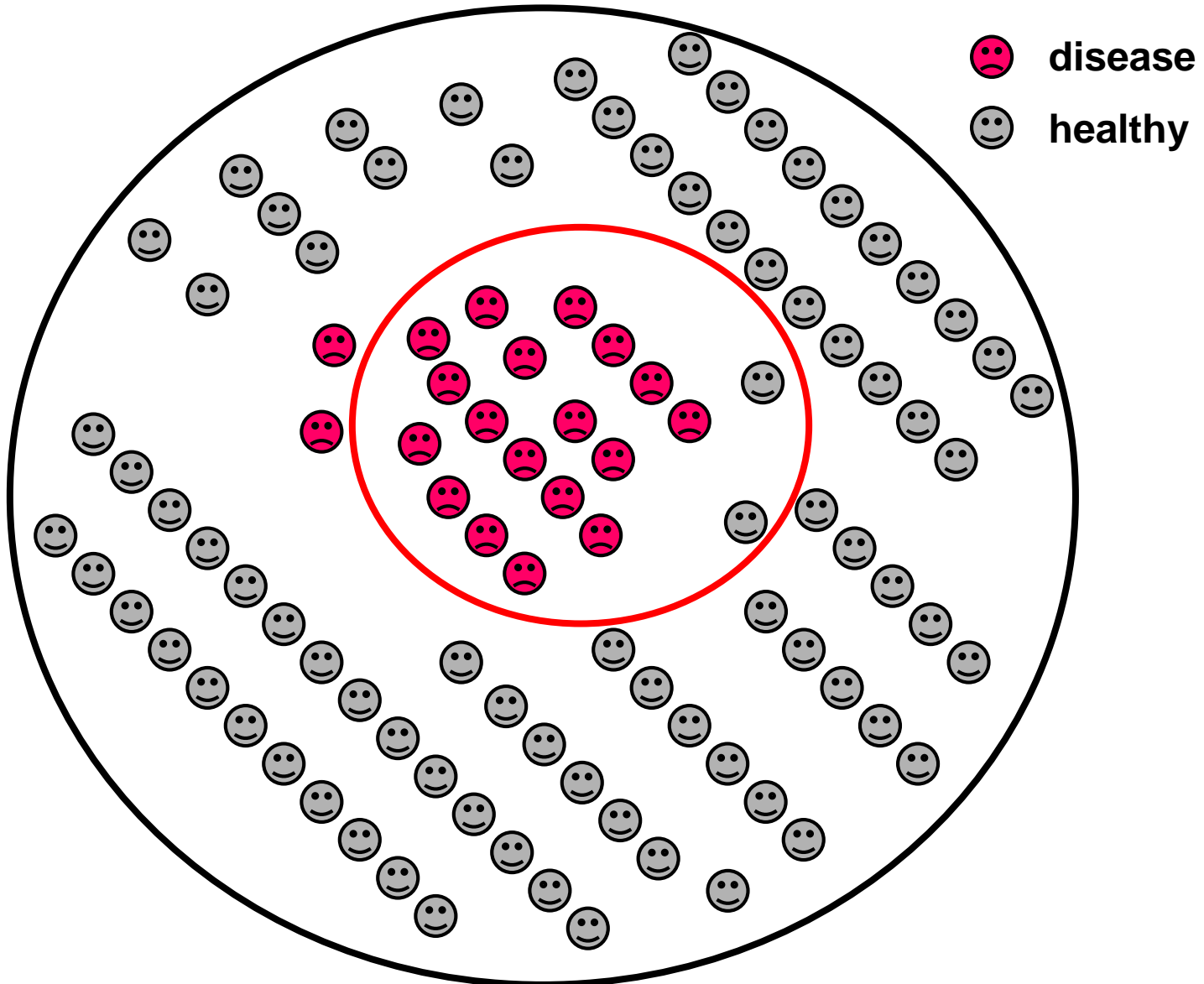
Population with diseased and healthy individuals



**A perfect diagnostic test
turns out positive for the diseased individuals only**



Real tests are not perfect



Confusion matrix: Evaluate a performed classification

Evaluation is done on a test set with known true class y and the predicted class \hat{y} .



id	true_class	pred_class
1	P	P
2	N	P
3	N	N
4	P	P
5	N	N
6	N	N

		True class	
		Positive	Negative
Predicted class	Positive	TP=2	FP=1
	Negative	FN=0	TN=3

Sensitivity and Specificity derived from a confusion matrix

Evaluation is done on a test set with known true class labels y and the predicted class label \hat{y} .

Predicted class	True class	
	Positive	Negative
Positive	TP	FP
Negative	FN	TN
		$sens = \frac{TP}{TP + FN}$ $spec = \frac{TN}{FP + TN}$

The **sensitivity** is derived from the positive examples and the **specificity** from the negative examples → both do not depend on the ratio of positive and negative classes in the test sample.

The **sensitivity** (recall) of a binary classifier is its **ability to identify correctly the positive class**.

Also called true positive rate (TPR) since it corresponds to the proportion of “Positive” instances that were classified as “Positive”

The **specificity** of a binary classifier is its **ability to identify correctly the negative class**.

Also called true negative rate (TNR) since it corresponds to the proportion of “Negative” instances that were classified as “Negative”

Positive predictive value (PPV) and negative predictive value (NPV)

Evaluation is done on a test set with known true class labels y and the predicted class label \hat{y} .

Predicted class	True class		
	Positive	Negative	
	Positive	Negative	
Positive	TP	FP	$PPV = \frac{TP}{TP + FP}$
Negative	FN	TN	$NPV = \frac{TN}{TN + FN}$
	$sens = \frac{TP}{TP + FN}$	$spec = \frac{TN}{FP + TN}$	

The **PPV** gives the probability that a instance, that was as “positive” predicted, is indeed “positive”.

The **NPV** gives the probability that a instance, that was as “negative” predicted, is indeed “negative”

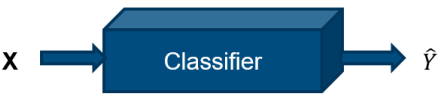
The **PPV** is derived from all as positive classified examples and the **NPV** from all as negative classified examples → both **depend on the ratio of positive and negative classes in the two prediction groups**.

Performance measures expressed as (conditional) probabilities

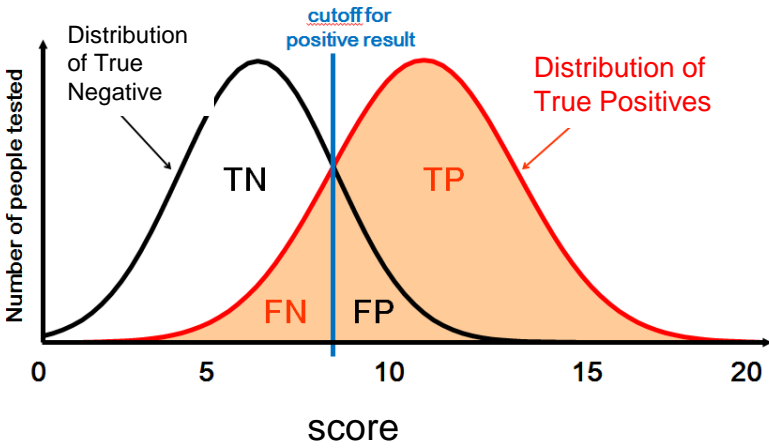
- $P(\hat{Y} = Y) = \text{acc}$: accuracy
- $P(\hat{Y} = 1 \mid Y = 1) = \text{Sens}$: true positive rate or sensitivity or recall
- $P(\hat{Y} = 0 \mid Y = 0) = \text{Spec}$: true negative rate or specificity
- $P(Y = 1 \mid \hat{Y} = 1) = \text{PPV}$: positive predictive value or precision
- $P(Y = 0 \mid \hat{Y} = 0) = \text{NPV}$: negative predictive value

		True class		
		Positive	Negative	
Predicted class	Positive	TP	FP	$PPV = \frac{TP}{TP + FP}$
	Negative	FN	TN	$NPV = \frac{TN}{TN + FN}$
		$sens = \frac{TP}{TP + FN}$	$spec = \frac{TN}{FP + TN}$	

Score based classifier

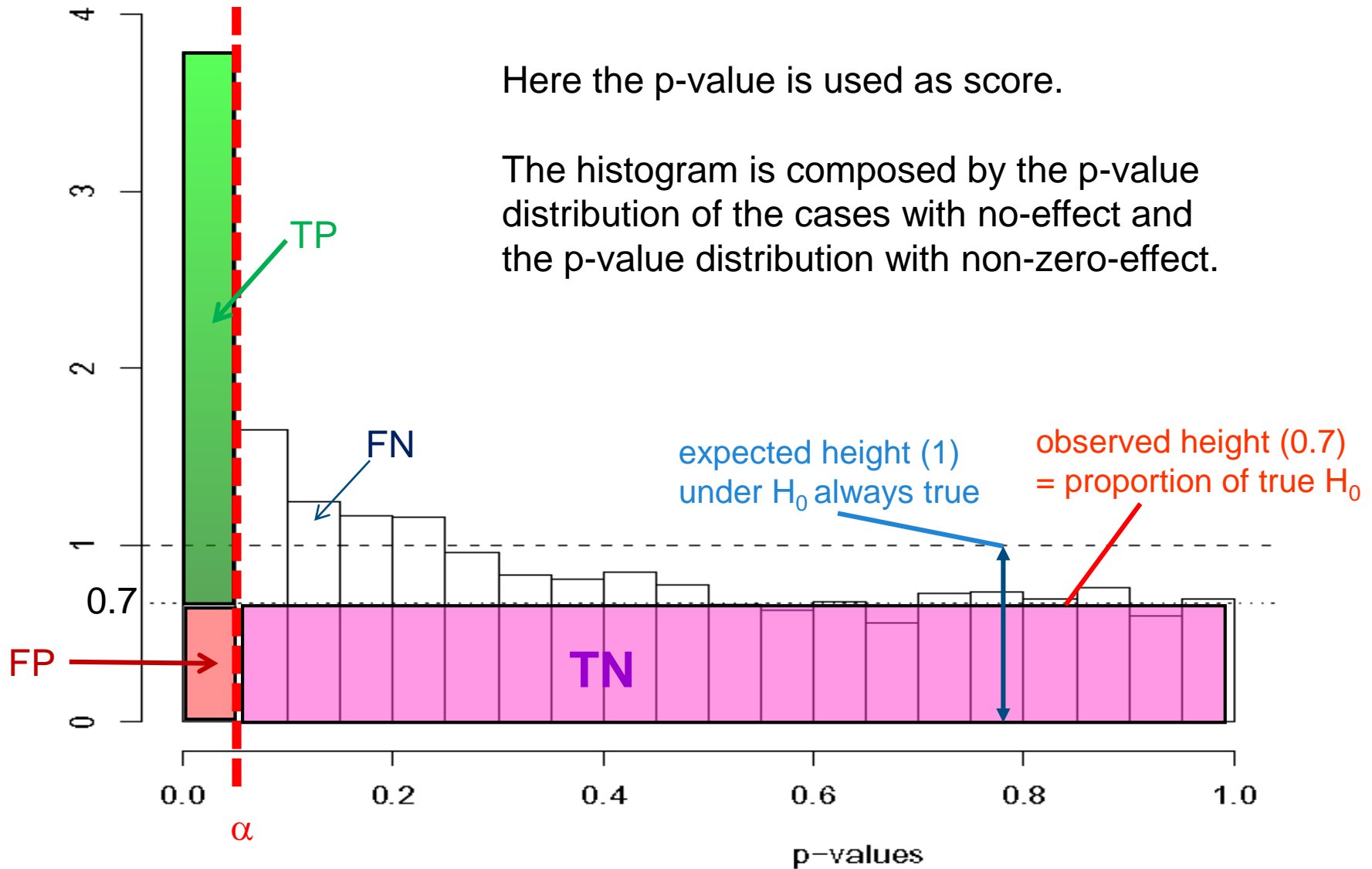


- Output: continuous score $\hat{Y}(x)$ (instead of actual class prediction)
- Discretized by choosing a cut-off
 - $\text{score} \geq c \rightarrow$ class «positive» or 1
 - $\text{score} < c \rightarrow$ class «negative» or 0

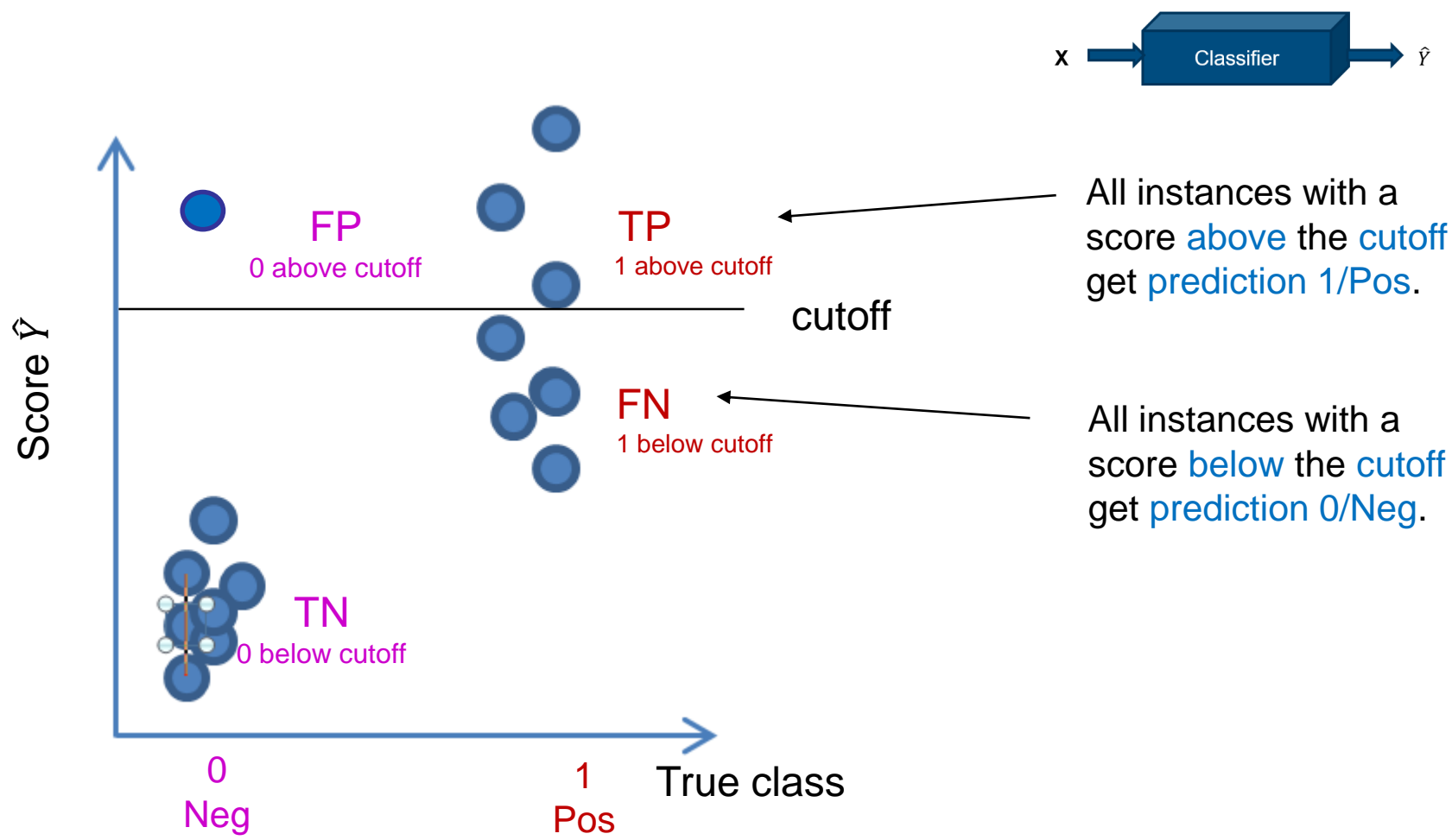


		True class		
		Positive	Negative	
Predicted class	Positive	TP	FP	$PPV = \frac{TP}{TP + FP}$
	Negative	FN	TN	$NPV = \frac{TN}{TN + FN}$
		$sens = \frac{TP}{TP + FN}$	$spec = \frac{TN}{FP + TN}$	

Recall the p-value histogram we can read off the content of the confusion matrix



Score based classifier



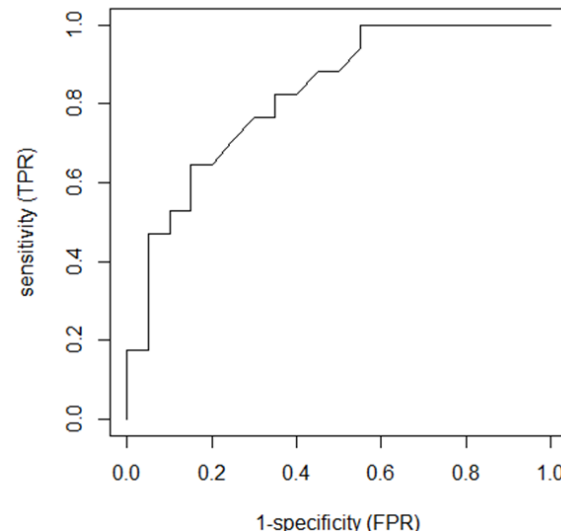
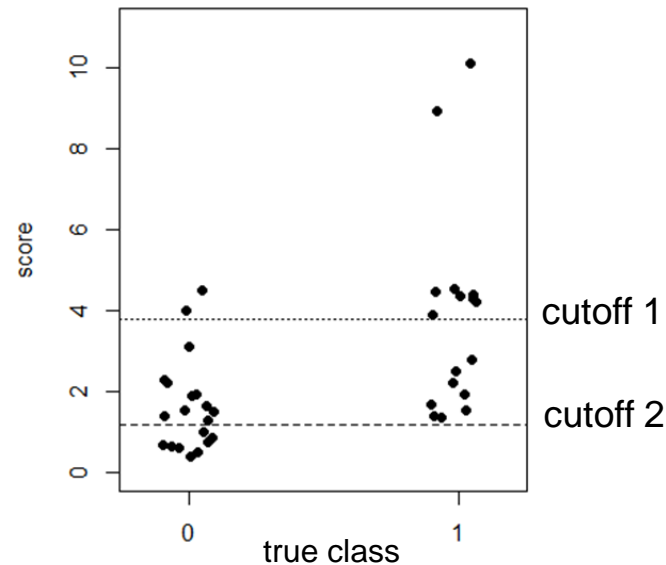
	True class	
	Positive	Negative
Positive	TP	FP
Negative	FN	TN

We can use a continuous score such as probability to construct a ROC curve

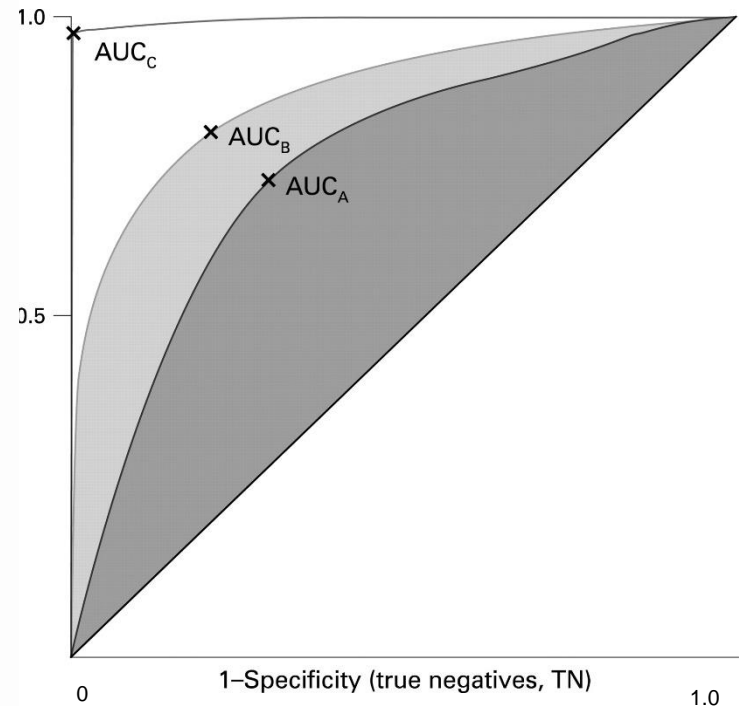
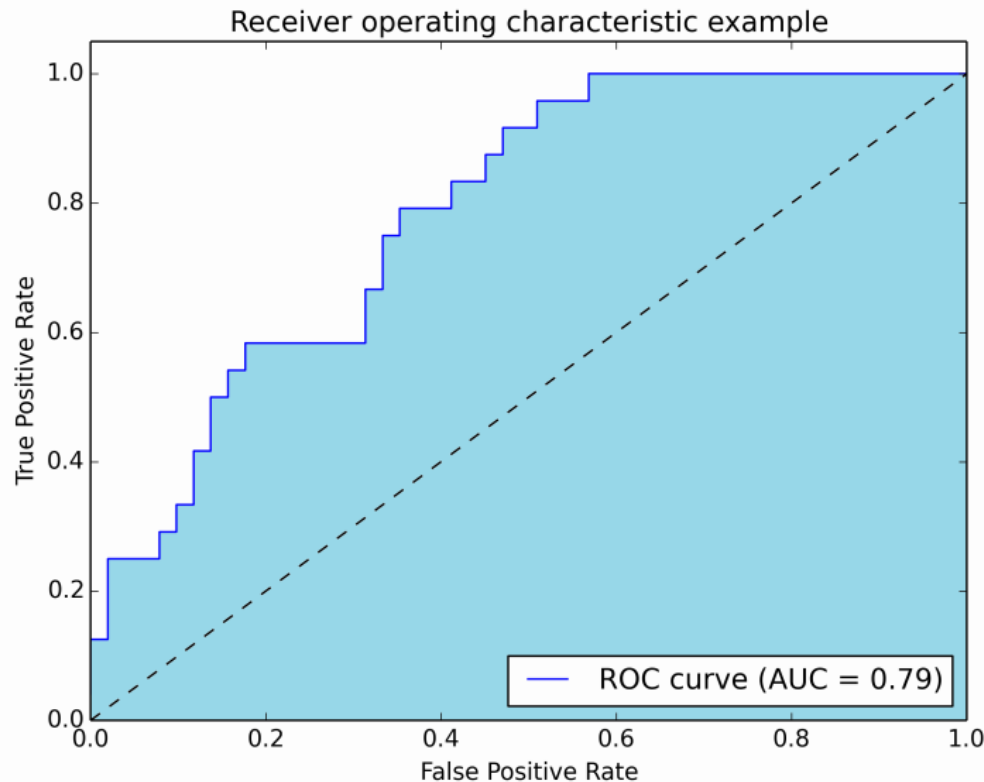
For each cutoff we get a classification rule (classify each observation with $\text{score} > \text{cutoff}$ as class 1) and a corresponding confusion matrix and can determine sensitivity and specificity

Determine the Sensitivity (true positive rate) and Specificity (true negative rate) for the indicated 2 cut-offs.

Do inn-class exercise



Use the ROC curve as performance measure by quantifying the area under the curve (AUC)

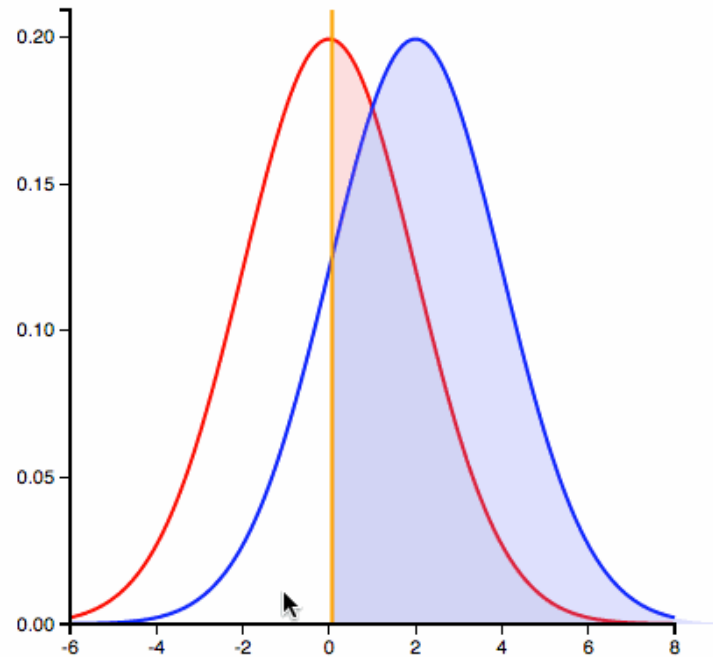
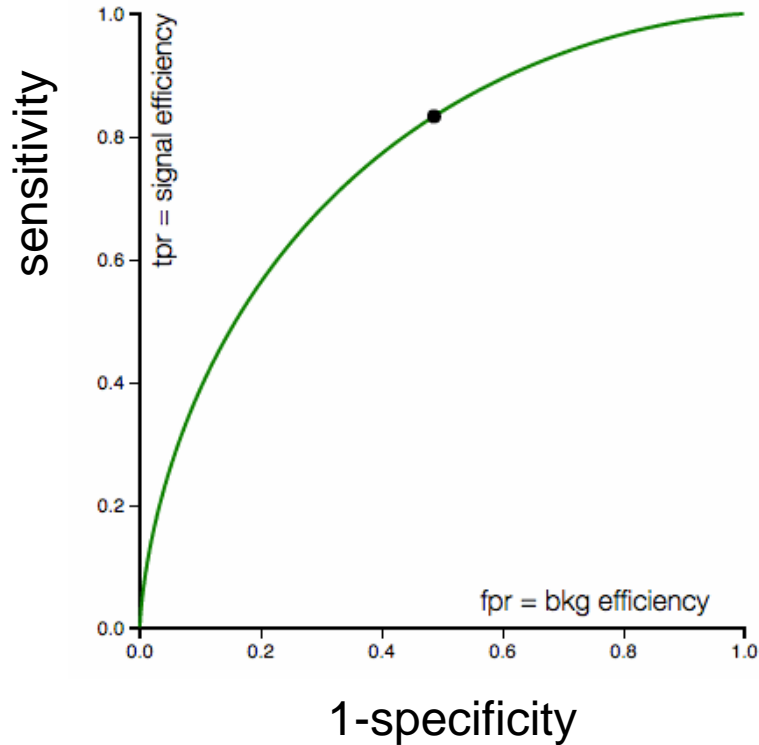


The larger the AUC the better is the performance of the diagnostic test.
A useless test has an $AUC = 0.5$.
A perfect test has an $AUC = 1$.

Nice online demos

ROC curve demo

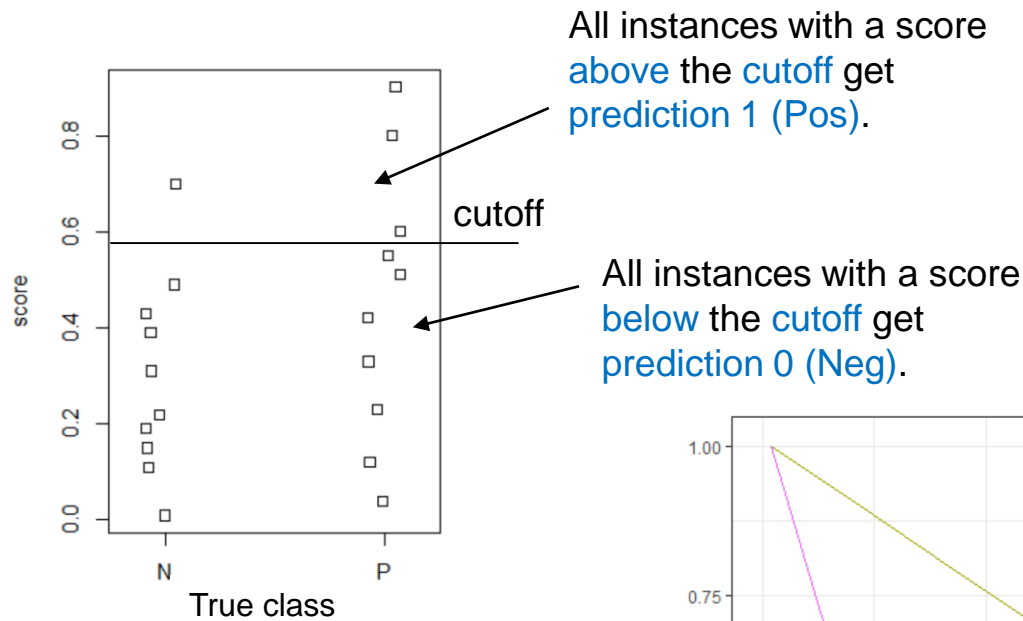
mean #1: mean #2: variance #1: variance #2:



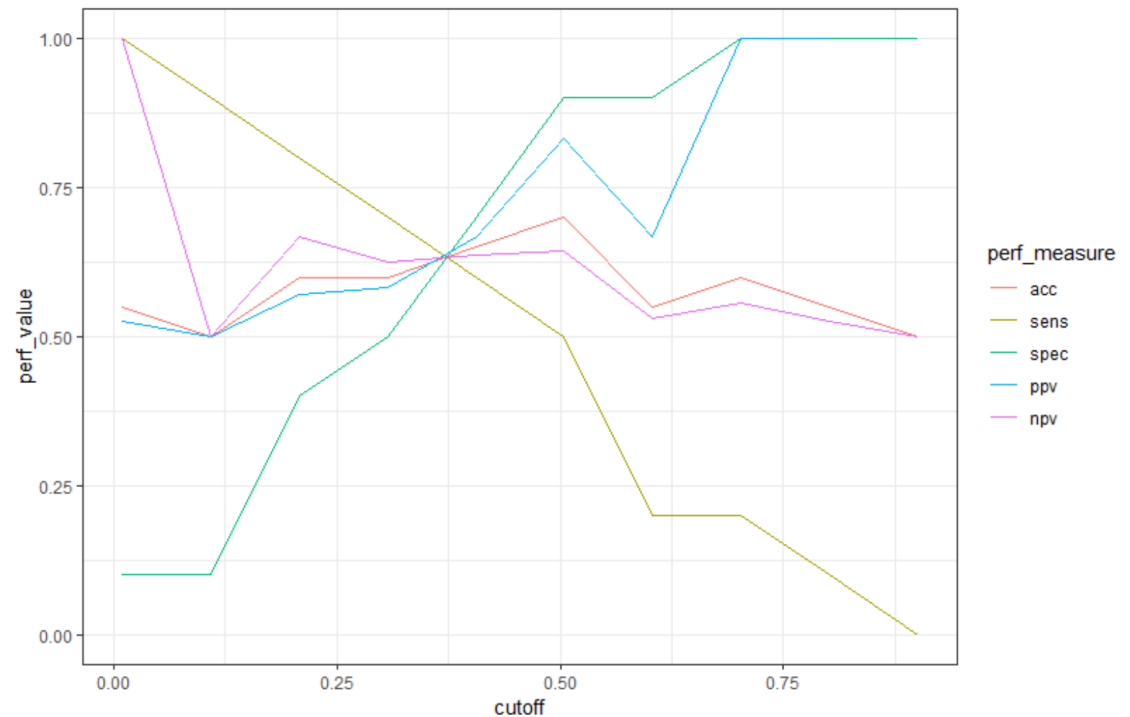
<http://arogozhnikov.github.io/2015/10/05/roc-curve.html>

Check out: <http://www.navan.name/roc/>
http://mlwiki.org/index.php/ROC_Analysis

Let's move the cutoff in scoring classifier and determine performance of resulting classification rule



Predicted class	True class		
	Positive	Negative	
Positive	TP	FP	$PPV = \frac{TP}{TP + FP}$
Negative	FN	TN	$NPV = \frac{TN}{TN + FN}$
	$sens = \frac{TP}{TP + FN}$	$spec = \frac{TN}{FP + TN}$	



How reliable is the result of a Aids-Test?

Ozzy Osbourne 'was told he could be HIV positive by doctors'

Rocker Ozzy Osbourne has revealed he was once told by doctors he could be HIV positive before a second test for the disease came back negative.



Ozzy Osbourne 'was told by doctors he could be HIV positive' Photo: AP

Prevalence, Sensitivity and Specificity

The probability that a randomly selected person has AIDS in Switzerland:
0.004

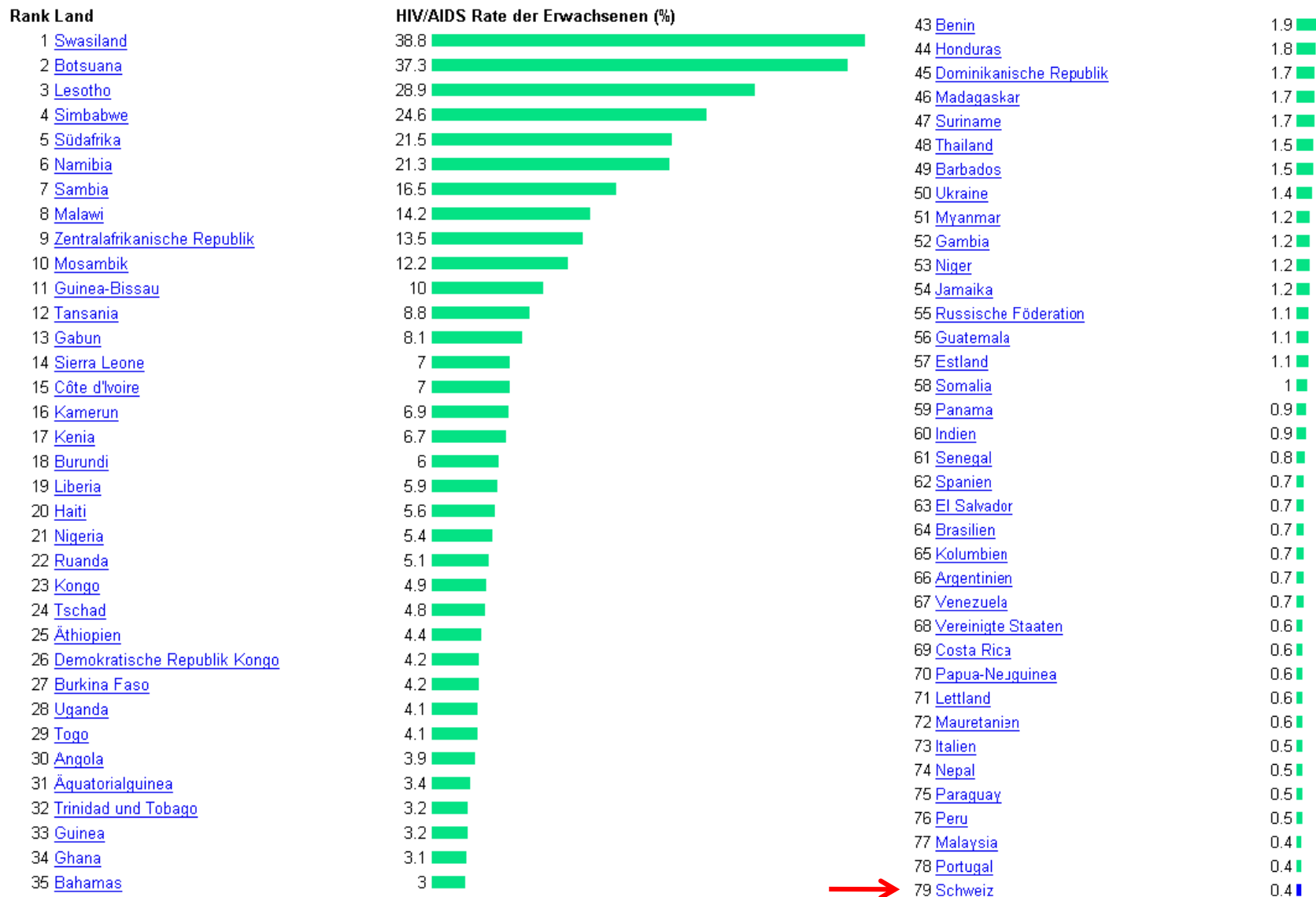
This is the **prevalence** of AIDS in Switzerland

Sensitivity of the ELISA-Test to detect a HIV+ blood sample:
0.999

Specificity of the ELISA-Test to identify a HIV- blood sample correctly: :
0.997

-> in-class exercise with topic screening with the Aids test:

HIV+/AIDS proportions in different countries



Confusion Matrix

From the tree diagram given in the in-class exercise we can read of the content of the corresponding confusion matrix.

	T +	T -	Summe
HIV +	30'769	31	30'800
HIV -	23'008	7'646'192	7'669'200
sum	53'777	7'646'223	7'700'000

Prevalence

$$P(HIV^+) = \frac{30800}{7700000} = 0.004$$

Sensitivity

$$P(T+ | HIV^+) = \frac{30769}{30800} = 0.999$$

Specificity

$$P(T- | HIV^-) = \frac{7646192}{7669200} = 0.997$$

Definition of the conditional probability

The conditional probability of an event (e.g. A or D+) given that some other event (e.g. B or T+) has already occurred is written as $P(A|B)$ and defined as the quotient of the probability of the joint of events A and B, and the probability of B. Der vertical dash means „given that“ or „under the condition“ B has already occurred.

A and B are two events and $P(B) \neq 0$. The conditional probability of A given B is defined as:

$$P(A | B) := \frac{P(A \cap B)}{P(B)}$$

Remark: If A and B are **independent**, we get:

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A) \cdot P(B)}{P(B)} = P(A)$$

Bayes's theorem

Inversion of a conditional probability

Bayes's theorem gives the rule how to invert a conditional probability, and how **to update the probability** by using some additional information:

Bayes' s theorem:

$$P(B | A) = \frac{P(A | B) \cdot P(B)}{P(A)} = \frac{P(A | B)}{P(A)} \cdot P(B)$$

posteriori probability for B
or updated probability
or predictive value

info

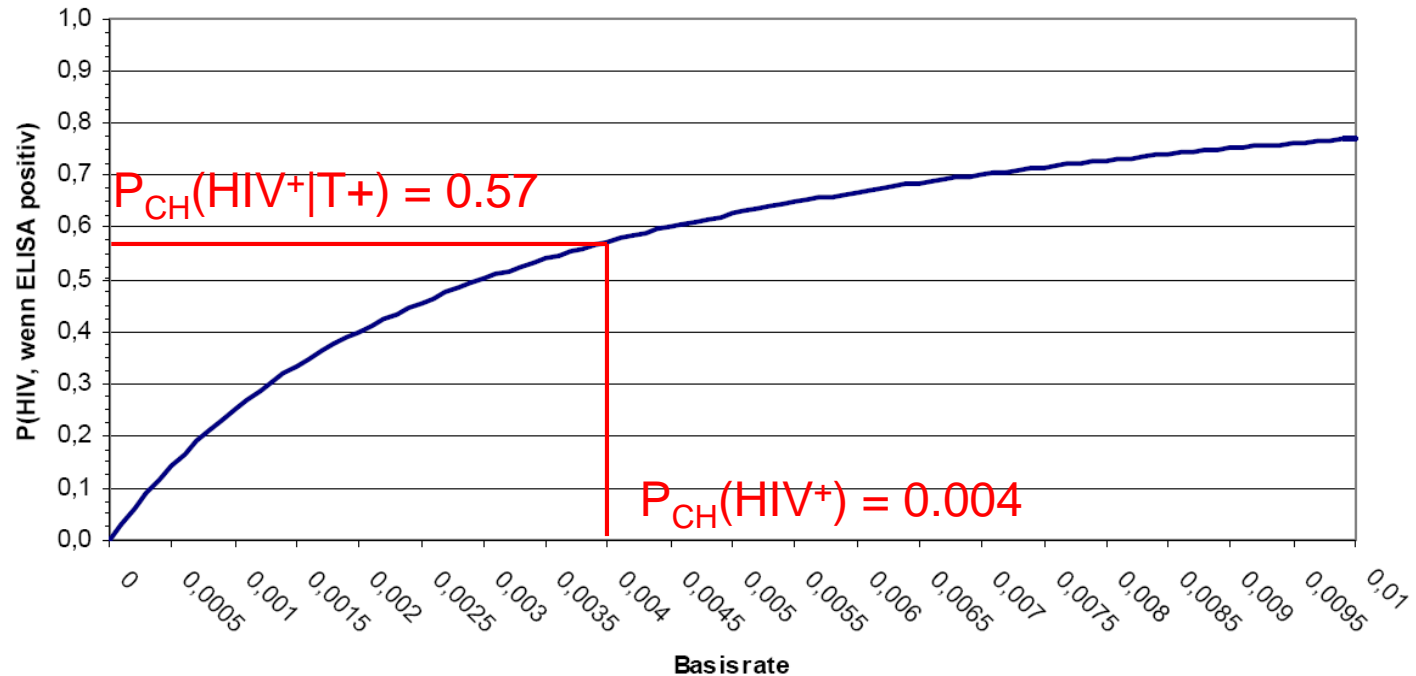
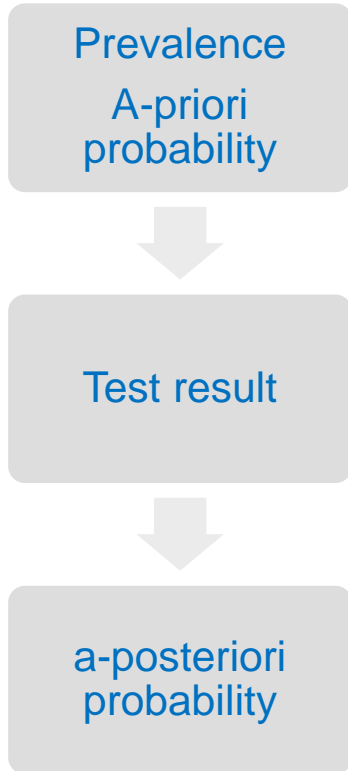
a-priori probability for B
or prevalence of B

proof:

$$P(B | A) := \frac{P(A \cap B)}{P(A)} = \frac{\frac{P(A \cap B)}{P(B)} \cdot P(B)}{P(A)} = \frac{P(A | B) \cdot P(B)}{P(A)}$$

$$P(A | B) := \frac{P(A \cap B)}{P(B)}$$

Wahrscheinlichkeiten neu bewerten: A-priori- und a-posteriori-Wahrscheinlichkeit



Inversion of a conditional probability

In general: $P(T+ | HIV^+) \neq P(HIV^+ | T+)$

Often we know a conditional probability as e.g.:

$$\text{Sensitivity: } P(T+ | HIV^+) = \frac{P(T+ \cap HIV^+)}{P(HIV^+)}$$

$$\text{Specificity: } P(T- | HIV^-) = \frac{P(T- \cap HIV^-)}{P(HIV^-)}$$

But we are interested in the predictive value of the diagnostic test which are the inversed conditional probabilities:

$$\text{positive predictive Value} \quad PPV = P(HIV+ | T+) = \frac{P(T_p | HIV^+) \cdot P(HIV^+)}{P(T_p)} = \frac{TP}{TP + FP}$$

$$\text{negative predictive Value} \quad NPV = P(HIV- | T-) = \frac{P(T- | HIV^-) \cdot P(HIV^-)}{P(T-)} = \frac{TN}{TN + FN}$$

Review: Power and level of significance, sensitivity and specificity of a test

A worked example

The **fecal occult blood** (FOB) screen test was used in 2030 people to look for bowel cancer:

		Patients with bowel cancer (as confirmed on endoscopy)		
		Condition Positive	Condition Negative	
Fecal Occult Blood Screen Test Outcome	Test Outcome Positive	True Positive (TP) = 20	False Positive (FP) = 180	Positive predictive value = $TP / (TP + FP)$ = $20 / (20 + 180)$ = 10%
	Test Outcome Negative	False Negative (FN) = 10	True Negative (TN) = 1820	Negative predictive value = $TN / (FN + TN)$ = $1820 / (10 + 1820)$ ≈ 99.5%
		Sensitivity = $TP / (TP + FN)$ = $20 / (20 + 10)$ ≈ 67%	Specificity = $TN / (FP + TN)$ = $1820 / (180 + 1820)$ = 91%	

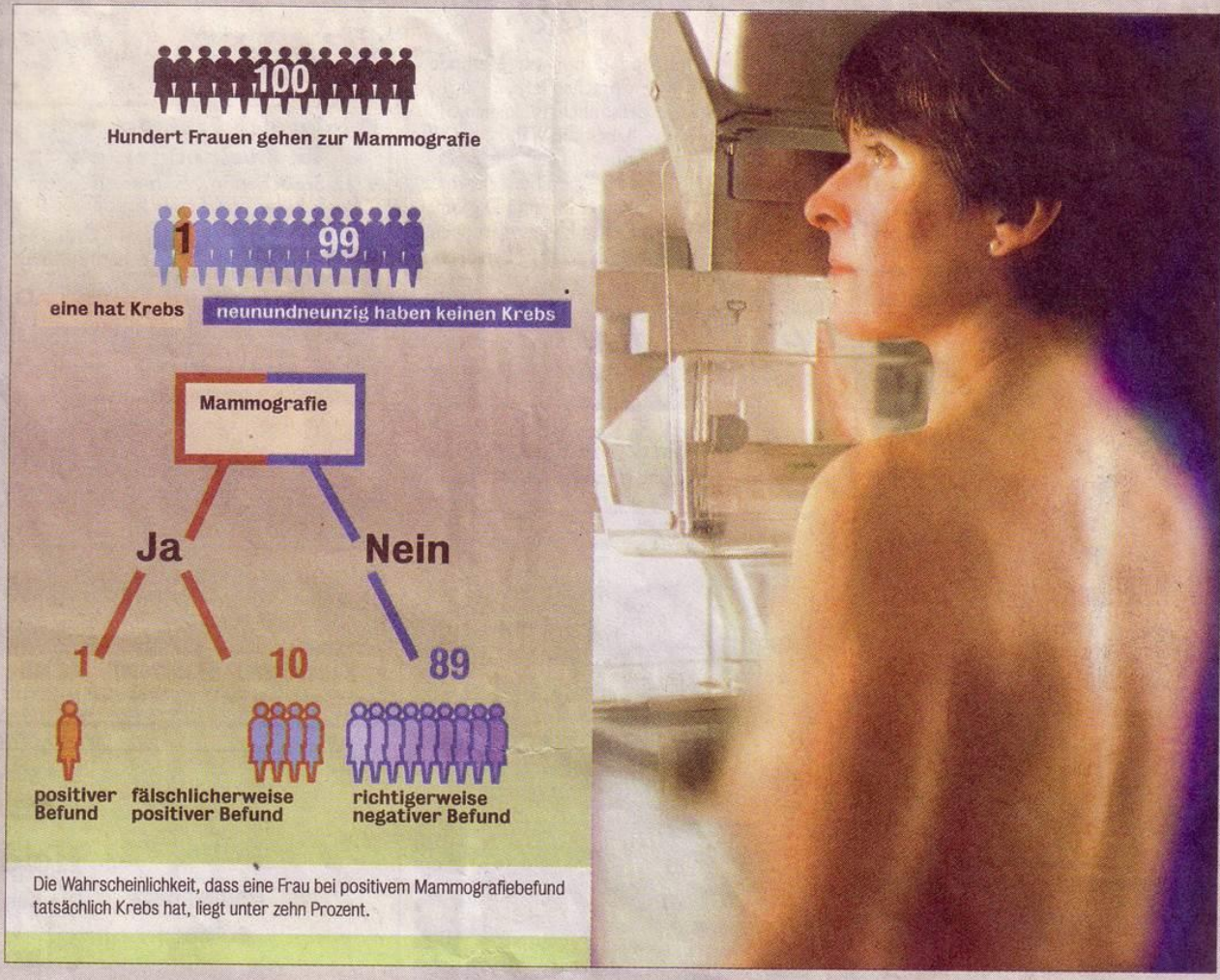
Characterizes
quality of test

Posteriori probability
Depends on
population e.g. the
prevalence

fecal occult blood test (FOBT) checks for hidden (occult) blood in the stool (feces, excrements)

«Tagesanzeiger» explains a-priori und a-posteriori probabilities

Brustkrebs: Was bedeutet ein positiver Mammografiebefund?



How to interpret a Mammography result

We can use the Bayes's theorem to determine the PPV and NPV of a Mammography result dependent on the prevalence.

prevalence	sensitivity	specificity	PPV	NPV
1.0%	86.6%	96.8%	21.5%	99.9%
4.5%	86.6%	96.8%	56.4%	99.3%
10.0%	86.6%	96.8%	75.1%	98.5%
50.0%	86.6%	96.8%	96.4%	87.9%

The breast cancer prevalence among British women aged 59 is 4.5%.
(<http://www.cancerresearchuk.org/cancer-info/>)

The negative predictive value (NPV) is with 99.3% much higher than the PPV of 56%

In the "One Million Women Study" (Banks et al. 2004) 122'355 50- 64 year old women who had a Mammography were followed for one year and the histological confirmed breast cancer incidences were determined.

The Fagan-Nomogram allows to graphically determine the posteriori probability

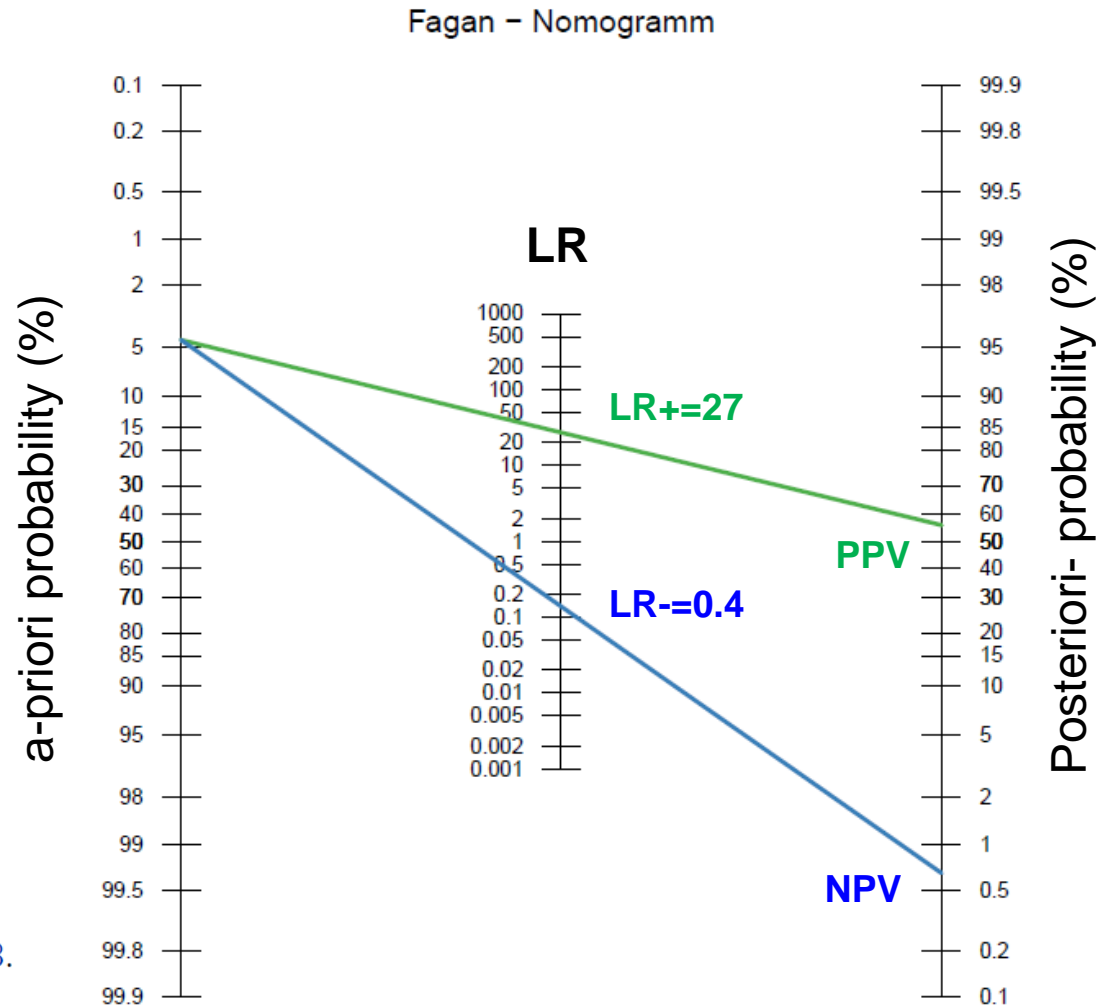
We can use **likelihood ratios** **LR+** and **LR-** to get graphically from the prevalence to the predictive value of a test:

$$LR+ = \frac{\text{sensitivity}}{1 - \text{specificity}}$$

$$LR- = \frac{1 - \text{sensitivity}}{\text{specificity}}$$

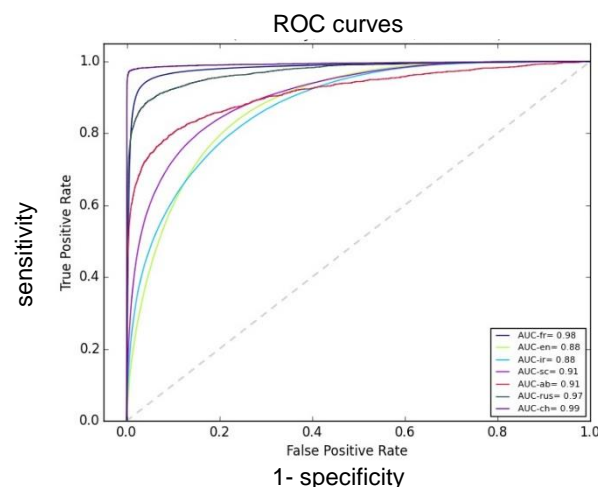
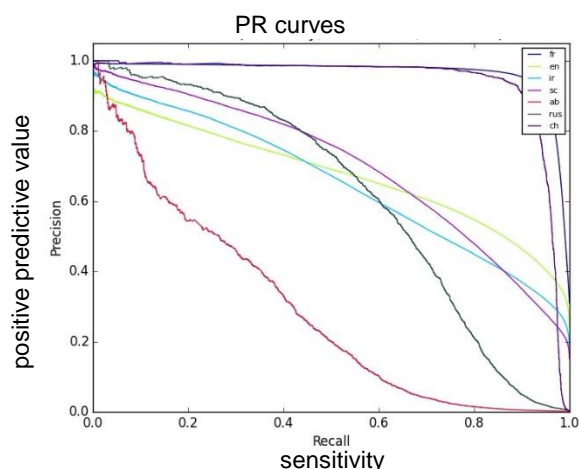
Mammography:
sensitivity=86.6%
specificity=96.8%

$$LR^+ = \frac{86.6\%}{3.2\%} \approx 27.1 \quad LR^- = \frac{13.4\%}{96.8\%} \approx 0.138.$$



Summary as extended confusion table & ROC and PR curves

		predicted condition			
total population		prediction positive	prediction negative		
true condition	condition positive	True Positive (TP)	False Negative (FN) (type II error)	True Positive Rate (TPR), Sensitivity, Recall, Probability of Detection $= \frac{\sum TP}{\sum \text{condition positive}}$	False Negative Rate (FNR), Miss Rate $= \frac{\sum FN}{\sum \text{condition positive}}$
	condition negative	False Positive (FP) (Type I error)	True Negative (TN)	False Positive Rate (FPR), Fall-out, Probability of False Alarm $= \frac{\sum FP}{\sum \text{condition negative}}$	True Negative Rate (TNR), Specificity (SPC) $= \frac{\sum TN}{\sum \text{condition negative}}$
$\text{Accuracy} = \frac{\sum TP + \sum TN}{\sum \text{total population}}$		Positive Predictive Value (PPV), Precision $= \frac{\sum TP}{\sum \text{prediction positive}}$	False Omission Rate (FOR) $= \frac{\sum FN}{\sum \text{prediction negative}}$	Positive Likelihood Ratio (LR+) $= \frac{TPR}{FPR}$	Diagnostic Odds Ratio (DOR) $= \frac{LR+}{LR-}$
		False Discovery Rate (FDR) $= \frac{\sum FP}{\sum \text{prediction positive}}$	Negative Predictive Value (NPV) $= \frac{\sum TN}{\sum \text{prediction negative}}$	Negative Likelihood Ratio (LR-) $= \frac{FNR}{TNR}$	



Remark: Unlike the ROC curve, PR curves are very sensitive to imbalance. A classifier that is optimized for good AUC, might yield poor precision-recall results on an unbalanced data.

Summary

- We need a (new) **test set with known true binary outcome to evaluate the performance** of a diagnostic test (or classifier)
- A binary diagnostic test (classifier) can be evaluated based on the
 - **confusion matrix** (determined in real world conditions) that allows to compute
 - test specific performance measures that do not depend on the disease prevalence
 - **sensitivity**: Probability that the test classifies a positive case as positive
 - **specificity**: Probability that the test classifies a negative case as negative
 - **accuracy**: overall classification rate
 - predictive performance measures that depend on the disease prevalence
 - **positive predictive value**: probability that a positive tested subject is sick
 - **negative predictive value**: probability that a negative tested subject is healthy
- A **diagnostic scoring test with continuous score** as outcome can be evaluated by using different **score-cutoffs to define positive and negative predictions**
 - by moving the cutoff we can determine a
 - **ROC curve** (sensitivity vs 1-specificity) and use the **AUC** (area under the curve) as performance measure
 - **PR curve** (Precision=positive-predictive-value vs Recall=sensitivity) and its AUC