

Screening and diagnostic tests

Gustaf Edgren, MD PhD
Associate professor of epidemiology
Karolinska Institutet

Outline

- Repetition of screening
 - Measures
 - Errors associated with screening
 - How to evaluate screening programs
- Advanced measures of diagnostic test performance
- Evaluation of screening programs
- Summary

Screening

- Screening is the strategy with which tests are performed in selected populations in order to find disease in otherwise asymptomatic individuals
- Importantly, unlike most other areas of medicine, screening is performed in **HEALTHY** people with **NO** sign of disease
 - This poses unique challenges!

Rationale for screening

- Intuitively, screening programs are implemented to:
 - Detect disease earlier, before spread has occurred,
 - and before symptoms have begun
 - Ultimately, the goal is to improve prognosis,
 - and to prevent disease spread between people (infectious disease)

When is screening suitable?

According to WHO guidelines from 1968:

1. The condition should be an important health problem
2. There should be a treatment for the condition
3. Facilities for diagnosis and treatment should be available
4. There should be a latent stage of the disease
5. There should be a test or examination for the condition
6. The test should be acceptable to the population
7. The natural history of the disease should be adequately understood
8. There should be an agreed policy on who to treat
9. The total cost of finding a case should be economically balanced in relation to medical expenditure as a whole
10. Case-finding should be a continuous process, not just a "once and for all" project

When is screening **really** suitable?

1. The disease should have serious consequences
2. The disease should be treatable and early treatment should improve prognosis
3. There should exist a simple, harmless and valid screening test
4. The prevalence of preclinical, asymptomatic disease should be sufficiently high in the screened population

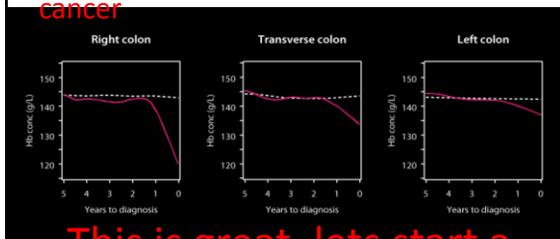
Screening tests

- A screening test is a test by which we want to be able to sort asymptomatic individuals into two groups:
 - Likely to have disease
 - Unlikely to have disease
- This division is usually quite unnatural, whereby most screening tests represent a dichotomy of a continuous measure which is usually:
 - A compromise between *too many false positives* and *too many false negatives*

Screening tests, cont.

- The ideal screening test would classify all persons with the disease as having the disease and all persons without the disease as not having the disease
- The ability of a test to correctly classify diseased and non-diseased is called the **test validity**
- The test's ability to perform exactly the same repeatedly is called the **test reliability**
- Which is more important?

Example – Hemoglobin and colon cancer

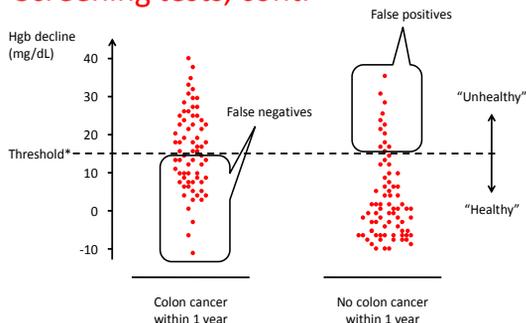


This is great, lets start a screening program!

Hemoglobin screening program

- Lets say we start measuring the hemoglobin concentration of everyone once per year to screen for colon cancer
- For every individual, we then compare this years result with the previous average to detect declining hemoglobin concentrations
- The question is, where do we draw the line between healthy and unhealthy?

Screening tests, cont.



*Criterion of positivity

The screening 2-by-2

		Gold standard / "truth"		
		+	-	
Screening test	+	True positives (a)	False positives (b)	a+b
	-	False negatives (c)	True negatives (d)	c+d
				a+b+c+d = N

Sensitivity

- The test sensitivity is a measure of the test's ability to correctly classify those **with** the disease:

$$\text{Sensitivity} = \frac{a}{a + c}$$

i.e. the proportion of those with the disease that are correctly classified as having the disease

Specificity

- The test specificity is a measure of the test's ability to correctly classify those **without** the disease:

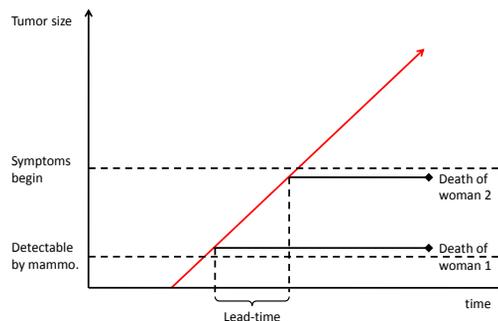
$$\text{Specificity} = \frac{d}{b + d}$$

i.e. the proportion of those without the disease that are correctly classified as not having the disease

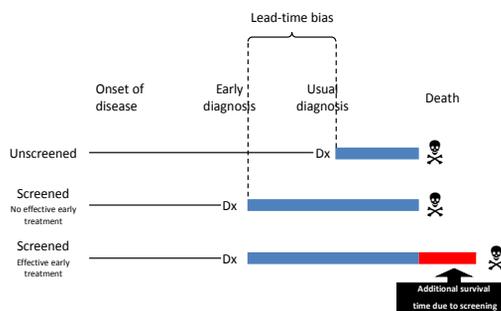
Lead-time bias

- In screening, the goal (and almost always the result) is to detect disease earlier in the disease progression
- Therefore, in an observational study assessing screening, cases detected through screening will appear to have a superior survival than cases detected clinically

Lead-time bias, cont.



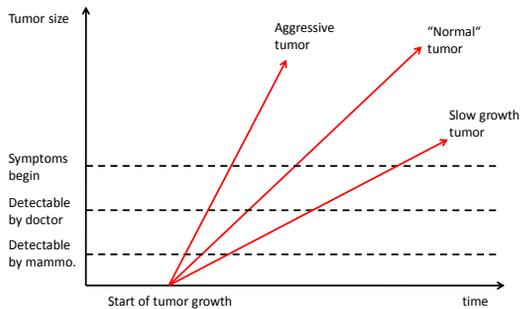
Lead-time bias, alternative explanation



Length-biased sampling

- Since all screening programs only screen participants at certain intervals (often several years), the probability of picking up pre-clinical disease depends on the aggressiveness of the disease (e.g. interval cases)
- Cases detected through screening are therefore often less aggressive and will have a better prognosis

Length-biased sampling



Over-diagnosis bias

- In addition to the classic biases (lead-time, length-bias and volunteer bias), observational studies are also susceptible to over-diagnosis bias:
 - For some conditions, the natural course of illness is often difficult to predict
 - For prostate cancer, as an example, there is considerable clinical heterogeneity and a paucity of methods for prognostication

Summary screening

- Screening is unique in medicine in that tests are performed in asymptomatic persons
- The ultimate goal of screening programs is to lessen the burden of disease by:
 - Earlier diagnosis
 - Diagnosis before start of symptoms
 - Prevent spread (both locally, as in disease progression, and between subjects)
- Due to unique types of bias, screening programs are best evaluated with RCTs

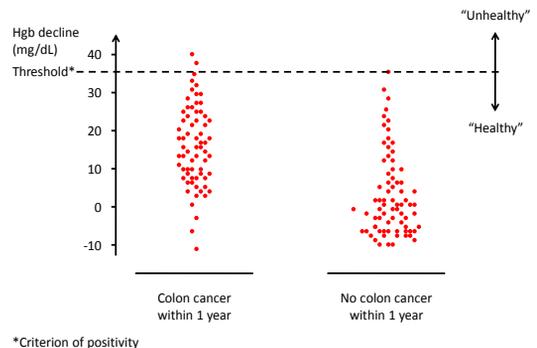
Summary screening, cont.

- Not all diseases are suitable for screening:
 1. The disease should have serious consequences
 2. The disease should be treatable and early treatment should improve prognosis
 3. There should exist a simple, harmless and valid screening test
 4. The prevalence of preclinical, asymptomatic disease should be sufficiently high in the screened population

Advanced test performance measures

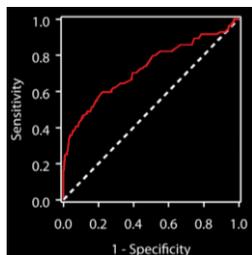
- In most screening tests, a continuous variable is measured and dichotomized into “sick” / “not sick” using a “criterion of positivity”
- The choice of threshold may seem arbitrary, but it can be optimized depending on the application of the test
- A common tool is the receiver operating characteristics (ROC) curve

ROC curves



ROC curves (2)

- So, what can we do with all these sensitivity/specificity values?
- Plot them of course!

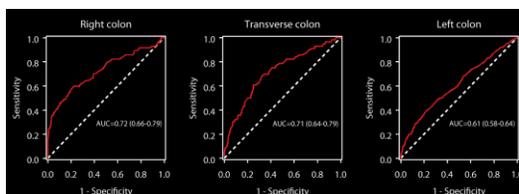


Comparing ROC curves

- The ROC curve is a graphical representation of what possible sensitivity/specificity values can be achieved with a certain test
- For each ROC curve, the area under the curve (AUC) can be estimated
- The AUC value can take any value between 0.5 and 1 (higher is better) and gives a summary measure of the test performance

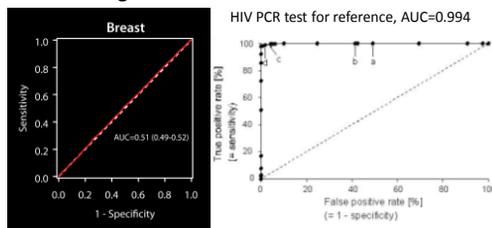
Comparing ROC curves

- The ROC curve is a graphical representation of what possible sensitivity/specificity values can be achieved with a certain test



Comparing ROC curves!

- Is hemoglobin changes a good test for breast cancer diagnostics?



Numbers needed to screen

- A classic measure of the successfulness of a screening program is numbers needed to screen (NNS)
- NNS is an analogue to NNT – Numbers needed to treat
- It is calculated as $1/ARR$ (=absolute risk reduction)
- NNS tells us how many individuals we need to screen to prevent one death (or whatever the outcome is)

Cost efficiency

- Cost efficiency calculations of screening notoriously difficult as they require factoring in of “costs” on many levels:
 - Cost of the screening program itself
 - Cost of “unnecessary” investigations and treatments*
 - Burden of false alarms
 - Burden of false reassurance
 - Pain and suffering

Screening RCT:s

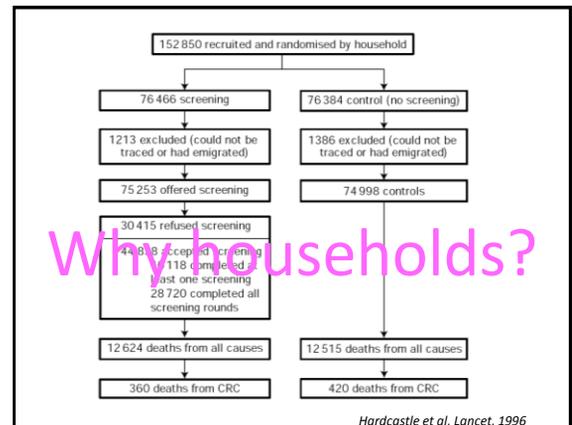
- We've already concluded that RCT:s are typically the study design of choice for screening program evaluation
- But, how would one design such a trial?

Pointers for screening RCT:s

- Depending of the expected gain of the screening, randomization can be on an individual level or in natural clusters
- In fact, in some cases, natural clusters can be necessary to promote acceptance
- Design the study for a HARD outcome: i.e. death*, or death from a certain cause
- Recall the ethical dilemmas of offering something only to one group – equipoise!

Screening RCT:s – limitations

- The screening RCT is typically an enormous undertaking, but ensures a high validity
- There are still some caveats, however:
 - If the prevalence of the screened disease is low, the trials have to be VERY large
 - If the disease of interest attracts public attention, your comparison group may be heavily polluted* and dilute the effects of the screening
 - While no threat to the internal validity, volunteer bias may limit the external validity**



Summary screening

- Screening is unique in medicine in that tests are performed in asymptomatic persons
- The ultimate goal of screening programs is to lessen the burden of disease by:
 - Earlier diagnosis
 - Diagnosis before start of symptoms
 - Prevent spread (both locally, as in disease progression, and between subjects)

Summary screening, cont.

- Not all diseases are suitable for screening:
 1. The disease should have serious consequences
 2. The disease should be treatable and early treatment should improve prognosis
 3. There should exist a simple, harmless and valid screening test
 4. The prevalence of preclinical, asymptomatic disease should be sufficiently high in the screened population

Summary screening, cont.

- The principal measures of the performance of a screening test are:
 - Sensitivity = the proportion of those with the disease that are correctly classified as having the disease
 - Specificity = the proportion of those without the disease that are correctly classified as not having the disease
- Due to unique types of bias, screening programs are best evaluated with RCTs