# biostat/epi:
# for the USMLE Step 1 +
# a light primer for topics in EBM

Matthew Kraushar

MD/PhD student
Rutgers/RWJMS/Princeton
Physician Scientist Program
matthew.kraushar@rutgers.edu
matthew.kraushar@gmail.com

# how well do we measure something

precision
consistent, reproducible, reliable



precise and accurate



accuracy
validity



not precise and not accurate

images from First Aid for the USMLE Step1, 2011

# how well do we measure something

## precision
random errors = low precision



## precise and accurate
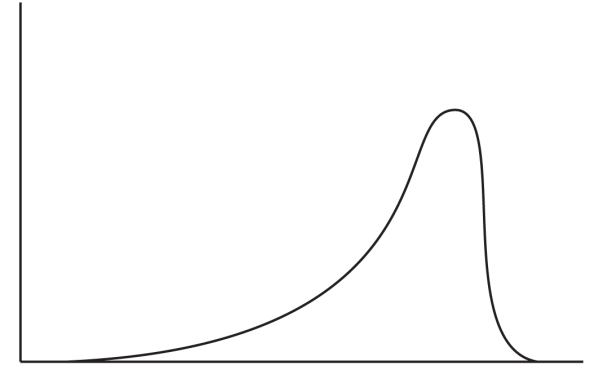


## accuracy
systematic errors = low accuracy



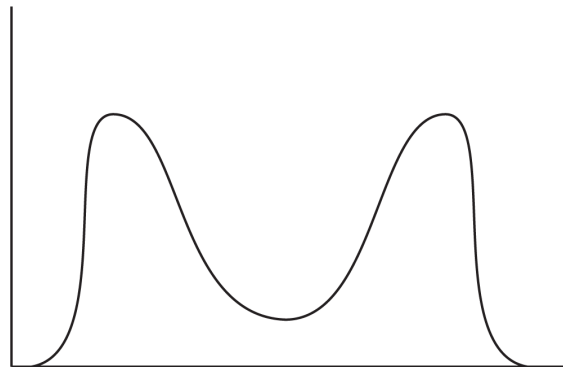## not precise and not accurate

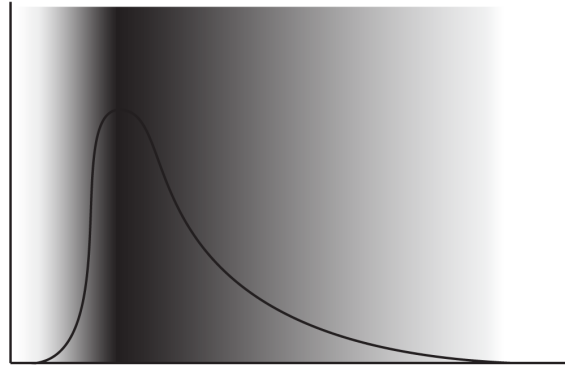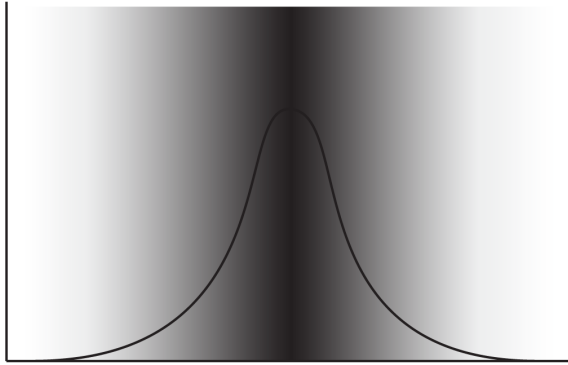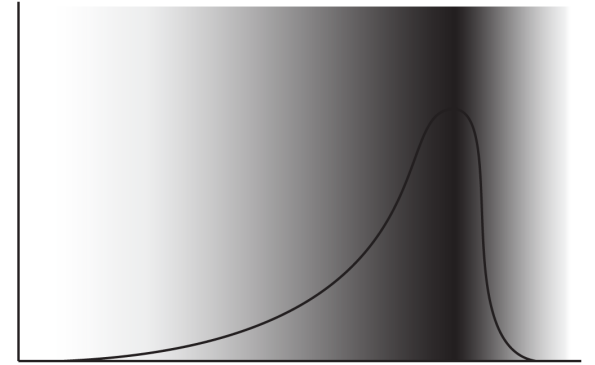images from First Aid for the USMLE Step1, 2011

# distribution

right skew

left skew

bimodal

# distribution

right skew

left skew

bimodal

sample density

# distribution



right skew

left skew

central tendency

# distribution

whatever

Mean Median Mode

0        whatever

Gaussian aka Normal

# distribution



Standard Normal

# distribution



Poisson:
average rate over time or distance

# distribution characterization

σ = standard deviation

σ/√n = standard error of the mean

# are they different?

2 means

# hypothesis testing

Study Results

|  | Ho | H1 |
|---|---|---|
| **Ho** | | false hit |
| **H1** | false miss | power |

Reality

# hypothesis testing

Study Results

|  | Ho | H1 |
|---|---|---|
| **Ho** |  | type 1 error |
| **H1** | type 2 error | power |

Reality

# hypothesis testing

Study Results

|  | Ho | H1 |
|---|---|---|
| **Ho** (Reality) | | $\alpha$ |
| **H1** (Reality) | $\beta$ | $(1-\beta)$ |

$\rightarrow$ p-value 0.05 (5%)

$\downarrow$ power

depends on:
    **sample size "n"** (more = better)
    size of effect
    subject compliance

# are they different?

confidence interval = range for repeated/multiple samples

mean +/- #

# are they different?

confidence interval = range for repeated/multiple samples
mean +/- #

mean

95%

# are they different?

2 means

# are they different?

2 means

# are they different?

3 means

# are they different?

percentage or proportions, compare to a known standard

(think Punnett square)

Chi square

$\chi 2$

is the observed frequency consistent with the expected frequency?

# are they different?

confidence interval = range for repeated/multiple samples
mean +/- 1.96(SEM)

z-score for
95% CI

-1σ  1σ

-2σ  2σ

-3σ  3σ

95%

# are they different?

confidence interval = range for repeated/multiple samples
   mean +/- 1.96(SEM)



beware: CI overlapping = not significantly different
beware: CI including 0 when difference between means
beware: CI including 1 when "relative" or "ratio" statistic

# linear vs. non-linear comparisons

non-linear

linear

# linear comparisons

regression ➜ r-value (correlation coefficient)

　　　-1 to 1

　　　closer to -1 or 1 = high correlation

　　　0 = no correlation

http://en.wikipedia.org/wiki/File:Linear_regression.svg

# linear comparisons

regression ➜ r-value (correlation coefficient)

-1 to 1

closer to -1 or 1 = high correlation

0 = no correlation



note: non-linear to linear transformations (think logs)

http://en.wikipedia.org/wiki/File:Linear_regression.svg

# contingency tables

|       |   +   |   -   |
|-------|-------|-------|
| **+** |       |       |
| **-** |       |       |

# contingency tables

diagnostic tests

|         | disease + | disease - |
|---------|-----------|-----------|
| test +  | TP        | FP        |
| test -  | FN        | TN        |

# contingency tables

diagnostic tests

disease

|  | + | - |
|---|---|---|
| test + | TP | FP |
| test - | FN | TN |

↓

sensitivity

# contingency tables

diagnostic tests



disease

|      | + | - |
|------|-----|-----|
| test + | TP | FP |
| test - | FN | TN |

TP/(TP+FN)

# contingency tables

diagnostic tests

disease

|  | + | - |
|---|---|---|
| test + | TP | FP |
| test - | FN | TN |

rules disease out

# contingency tables

diagnostic tests

|  | disease + | disease - |
|---|---|---|
| test + | TP | FP |
| test - | FN | TN |

rules disease out

# contingency tables

diagnostic tests

disease

|  | + | - |
|---|---|---|
| **+** | TP | FP |
| **-** | FN | TN |

test

example: ELISA

# contingency tables

diagnostic tests

|  |  | disease | |
|---|---|---|---|
|  |  | **+** | **-** |
| **test** | **+** | TP | FP |
|  | **-** | FN | TN |

specificity

# contingency tables

diagnostic tests

disease

|  |  | + | − |
|---|---|---|---|
| **test** | **+** | TP | FP |
|  | **−** | FN | TN |

TN/(TN+FP)

# contingency tables

diagnostic tests

disease

|       | + | − |
|-------|-----|-----|
| test + | TP | FP |
| test − | FN | TN |

rules disease in

# contingency tables

diagnostic tests



disease

|  | + | - |
|---|---|---|
| test + | TP | FP |
| test - | FN | TN |

rules disease in

# contingency tables

diagnostic tests

|       | disease + | disease - |
|-------|-----------|-----------|
| test + | TP | FP |
| test - | FN | TN |

example: Western

# contingency tables

diagnostic tests

# contingency tables

diagnostic tests

# contingency tables

diagnostic tests

disease

|  | + | - |  |
|---|---|---|---|
| test + | TP | FP | → positive predictive value |
| test - | FN | TN |  |

# contingency tables

diagnostic tests



disease

|  | + | - |  |
|---|---|---|---|
| **test +** | TP | FP | → TP/(TP+FP) |
| **test -** | FN | TN |  |

# contingency tables

diagnostic tests

# contingency tables

diagnostic tests

disease

|  | + | - |
|---|---|---|
| **+** | TP | FP |
| **-** | FN | TN |

test

→ TN/(TN+FN)

# contingency tables

risk factors

disease

|  | + | - |
|---|---|---|
| **+** | a | b |
| **-** | c | d |

risk factor

# contingency tables

risk factors: odds

disease

|        |   | +         | -         | odds ratio |
|--------|---|-----------|-----------|------------|
| risk factor | + | a | b | → a/b |
|        |   |           |           | ÷ |
|        | - | c | d | → c/d |

# contingency tables

risk factors: risk

disease

|  | + | − | relative risk |
|---|---|---|---|
| + | a | b | → a/(a+b) |
| − | c | d | → c/(c+d) |

risk factor

÷

# contingency tables

risk factors: risk

disease

# contingency tables

treatments

disease

|  | + | - |
|---|---|---|
| treatment | a | b |
| placebo | c | d |

# contingency tables

treatments: risk reduction

disease

|  | + | - |  |
|---|---|---|---|
| treatment | a | b | **a + b** |
| placebo | c | d | **c + d** |

↓

absolute risk reduction

everyone

# contingency tables

treatments: risk reduction

disease

|  | + | - | everyone |
|---|---|---|---|
| treatment | a | b | **a + b** |
| placebo | c | d | **c + d** |

(c/everyone) − (a/everyone)

**c + d**          **a + b**

# contingency tables

treatments: number needed to treat

# contingency tables

treatments: number needed to treat

|  | dead | alive | |
|---|---|---|---|
| treatment | a | b | **a + b** |
| placebo | c | d | **c + d** |

↓

1/(absolute risk reduction)

everyone

# contingency tables

treatments: number needed to treat



|  | dead | alive | |
|---|---|---|---|
| treatment | a | b | **a + b** |
| placebo | c | d | **c + d** |
| | **c + d** | **a + b** | everyone |

$$1/[(c/\text{everyone}) - (a/\text{everyone})]$$

# contingency tables

treatments: number needed to treat

10%

% dead

placebo  treatment

for more, see: http://www.thennt.com/

# contingency tables

treatments: number needed to treat

# contingency tables

treatments: number needed to treat



100%

% dead

} 2% avoid death

placebo    treatment

}8% die regardless of treatment

} 80% survive regardless of treatment

everyone!!!

for more, see: http://www.thennt.com/

# contingency tables

treatments: number needed to harm

# contingency tables

treatments: number needed to harm

|  | dead | alive | attributable risk |
|---|---|---|---|
| treatment | a | b | → $a/(a+b)$ |
|  |  |  | - |
| placebo | c | d | → $c/(c+d)$ |

everyone

# contingency tables

treatments: number needed to harm

# contingency tables

treatments: number needed to harm

# contingency tables

treatments: NNT vs. NNH

## What happens when NNT > NNH?

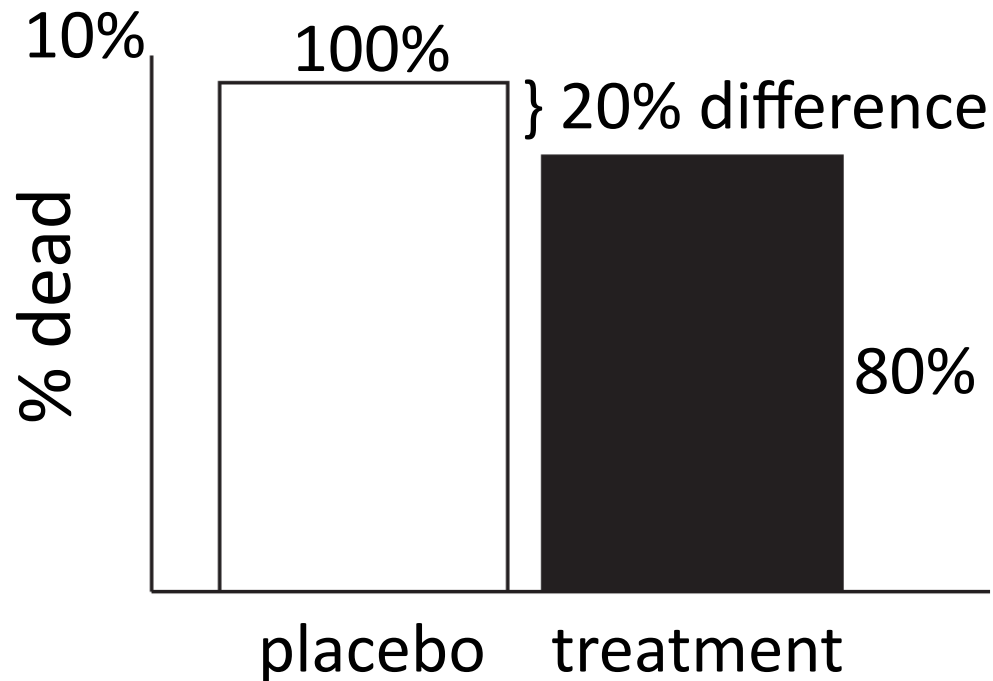# how many sick people are there?

prevalence = how many people are sick
    single time point
    point prevalence = (total # sick)/(total population)

incidence = how many people are getting sick, new cases
    time interval
    incidence = (new cases over time)/(total population **at risk***)

**at risk* = excludes currently have disease or previously positive**

# how many sick people are there?

prevalence = how many people are sick
    single time point
    point prevalence = (total # sick)/(total population)

incidence = how many people are getting sick, new cases
    time interval
    incidence = (new cases over time)/(total population **at risk***)

prevalence ≅ incidence x disease duration
chronic disease: prevalence > incidence
acute disease: prevalence ≅ incidence

# study types

- case control: retrospective, observational
  - disease vs. no disease → look for risk factor
  - statistic: **odds ratio**

disease

$+$       $-$

odds ratio

|  | | |
|---|---|---|
| a | b | → a/b |
| c | d | → c/d |

risk factor    $+$     $\div$

$-$

# study types

- cohort: prospective, observational
    risk factor present vs. absent → look for disease
    statistic: **relative risk**

disease

+        -

relative risk



risk factor

+ → a/(a+b)

÷

- → c/(c+d)

# study types

- cross sectional: single time point, observational
       disease or risk factor presence
       statistic: **prevalence**


       prevalence = how many people are sick
          single time point
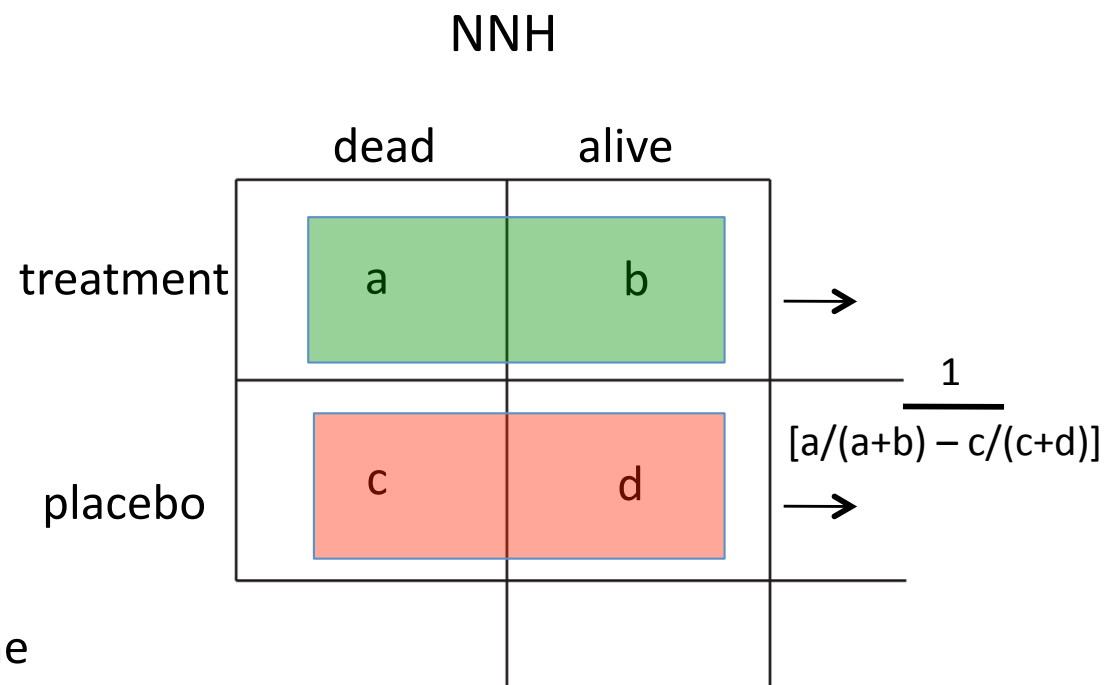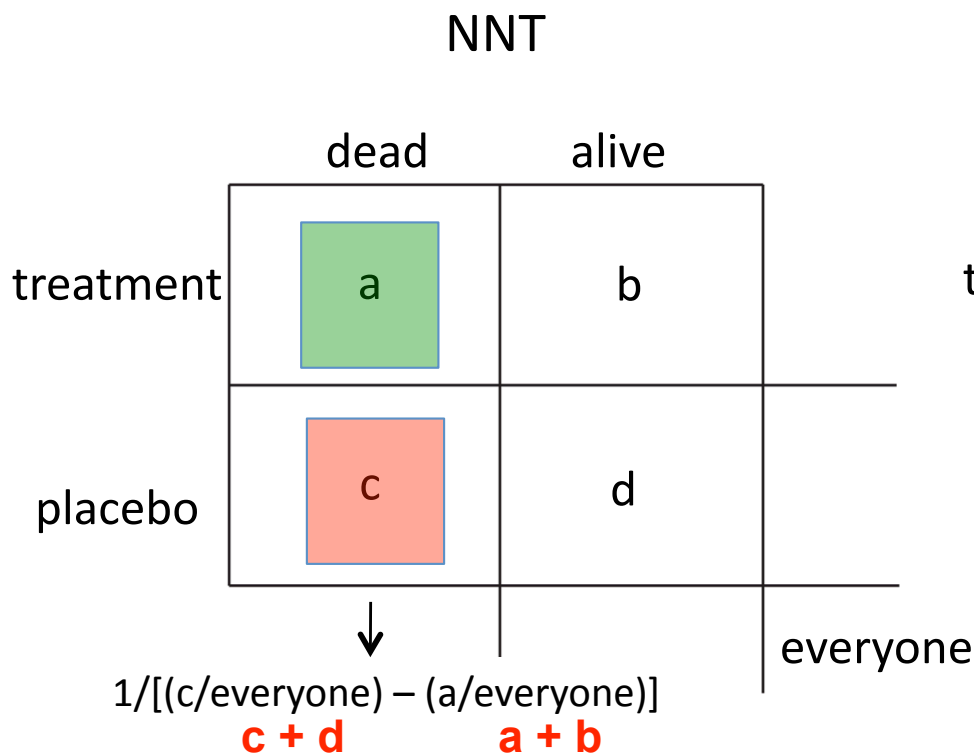          point prevalence = (total # sick)/(total population)

# study types

- randomized controlled trial: prospective, interventional, gold
  treatment vs. placebo → outcome
  statistic: NNT/NNH, multiple comparisons

# study types

- randomized controlled trial: prospective, interventional, gold
        treatment vs. placebo → outcome
        statistic: NNT/NNH, multiple comparisons

phase I:
small # pts, healthy volunteers = **safety**, toxicity, pharmacokinetics

phase II:
small # pts, have dz, placebo control = tx **efficacy**, dosing, adverse rxns

phase III:
large #pts, randomly assigned tx vs control = new vs old SOP tx or placebo

APPROVAL

phase IV:
post-marketing **surveillance** = for rare, long term effects

# study types

- meta-analysis: retrospective, chart review, platinum
         pools data from multiple studies, usually RCTs

# study types

- twin concordance study

      genes vs. environment questions

      genetic heritability: monozygotic vs. dizygotic

      think early life, autism etiology

# study types

- adoption study
> genes vs. environment questions
> think later life, schizophrenia etiology

# bias

bias = systematic error

- selection bias: nonrandom assignment into groups, loss to follow up

- recall bias: knowledge of presence of dz changes subject's response

- sampling bias: subjects are not representative of population, non-generalizable

- late-look bias: information gathered at inappropriate time, eg: giving survey on fatal dz to
  alive pts

- procedure bias: subjects in different groups not treated the same

- confounding bias: one factor distorts/confuses the effect of another closely related one

- lead-time bias: early detection confused with ↑ survival...seen with improved screening

- pygmalion effect: researchers belief changes outcome of tx...researcher's behavior

- hawthorne effect: subjects change behavior when the they know they're being
  studied...subject's behavior

- observer bias: researcher's decision affected by prior knowledge of subject's exposure status

~~matthew.kraushar@rutgers.edu~~

matthew.kraushar@gmail.com

www.matthewkraushar.com