# Another "Doc Squirrel" and "Kid Cat" Adventure!!

# Cohorts

They sure did smoke alot in the 1950s and '60s!



By Stefan Tigges MD MSCR and David Schulman MD MPH

To show that cholera was most common in low lying parts of Oxford.

Correct. That distribution was used to support the miasma theory; poisonous vapors were thought to concentrate at the lowest elevations.

So the miasmists had persuasive maps too!

And graphs. This 1852 graph shows higher cholera mortality rates (per 10,000) at lower elevations. I have highlighted some death rates on the right and below.
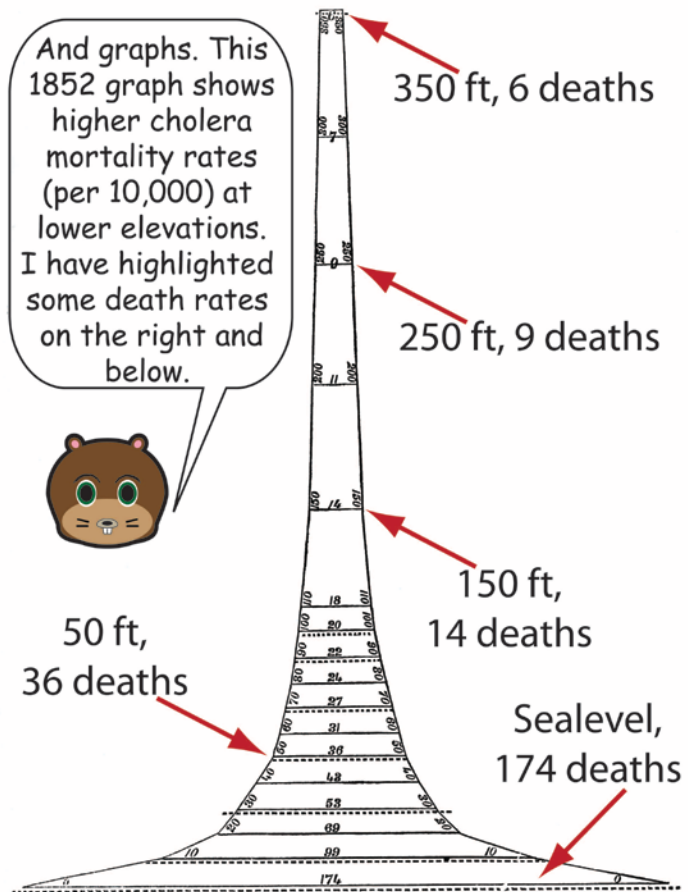
350 ft, 6 deaths

250 ft, 9 deaths

150 ft, 14 deaths

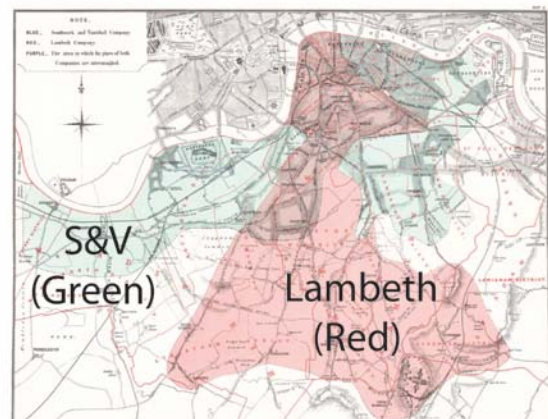50 ft, 36 deaths

Sealevel, 174 deaths

2

Why are you rehabilitating the miasma theory?

I'm not. I am trying to show why many of Snow's contemporaries were unconvinced by his map. In fact, Edmund Parkes thought Snow's map was evidence of an "atmospheric cause" and observed that "[t]here are, indeed, so many pumps in this district, that wherever the outbreak had taken place, it would most probably have had one pump or another in its vicinity."

Did Dr. Snow have other evidence for his water borne theory of cholera spread besides the Broad Street map?

He did. Here is another map showing parts of south London supplied by 2 rival drinking water companies. The Lambeth Company drew its' water upstream from London, far away from sewage dumped into the Thames. The Southwark and Vauxhall (S&V) Company piped its' water from further downstream, where the river was contaminated with sewage.

S&V
(Green)

Lambeth
(Red)

FUN.—August 18, 1866.

DEATH'S DISPENSARY.

OPEN TO THE POOR, GRATIS, BY PERMISSION OF THE PARISH.

|  | Lambeth | S&V |
|---|---|---|
| Cholera Deaths | 461 | 4,093 |
| Population | 173,748 | 266,516 |
| Mortality Rate per 1000 | 2.7 | 15.4 |

This table compares cholera deaths, total population, and cholera mortality rates of Lambeth and S&V water drinkers for a 14 week period in 1854. What do you think?

S&V customers are dying of cholera at much higher rates than Lambeth customers.

Correct. This type of study is called a cohort study. We compare outcomes in 2 (or more) groups or populations at risk of developing a disease; an exposed group and an unexposed group.

Exposed healthy people → Time → Number of diseased people

Unexposed healthy people → Time → Number of diseased people

Exposed to what?

In this case, S&V customers were exposed to water contaminated by sewage. Since Lambeth customers drank uncontaminated water, they are considered unexposed. The exposed group (S&V customers) had an almost 6-fold higher cholera mortality rate than the unexposed (Lambeth customers) group.

SOURCE OF THE SOUTHWARK WATER WORKS

4

Good. Now let's introduce the idea of disease prevalence. What percentage of people in Group 1 are zombies?

There are none, so zero percent.

Disease prevalence is defined as the proportion of people in a population who have the disease of interest at a given time. Prevalence=# diseased/total population. In Group 1, the prevalence of zombism is zero.

So in Group 2, the prevalence of zombism is 20%.

## Group 1

## Group 2

Let's try another example. You have a group of 10 people, 2 of whom are zombies. Ravenous, fast zombies, lurking behind shower curtains. What will happen to the number of zombies?

The number of zombies will increase over time.

Incidence is the metric we use to measure new cases of disease in a population. It is calculated by dividing new cases by the size of the population at risk. After one week, the number of zombies has increased from 2 to 6 out of 10. What is the incidence of zombism?

At the beginning of the week, there were 8 people in the original group who were normal and therefore at risk of turning into zombies. At the end of the week, there were 4 zombies out of a population at risk of 8 people, so the incidence is 4/8 or 50%.

↓ One Week Later ↓

8

9

10

11

You can calculate a relative risk using incidence rates in the same way that you did using cumulative incidence.

We divide the incidence rate in the exposed (IE) by the incidence rate in the unexposed (IU): IE/IU=0.50 cases per person week/0.25 cases per person week=2.

You can also calculate an absolute risk difference using incidence rates by subtracting the incidence in the unexposed from the exposed.

OK: IE-IU=0.50 cases/person week-0.25 cases/person week=0.25 cases/person week. That means that in 4 person weeks, we will get one new case in the exposed compared to the unexposed.

What would you expect the relative risk and the risk difference to be if the exposure did not increase the probability of becoming a zombie?

If closed shower curtains don't increase the chance of becoming a zombie, then the incidence in the exposed and unexposed groups is equal. In that case, the relative risk would be equal to one and the risk difference would be zero.

Correct. As we will see later when we discuss clinical trials, we explicitly assume that our exposure or intervention does not effect the outcome. That is called the null hypothesis.

Weird. If we assume that an exposure has no effect, why do a study?

Medical research is a little backwards. We assume that outcomes are the same in the exposed and unexposed groups so that we can do some statistical calculations that (we hope!) will provide evidence against the null hypothesis.

That is backwards.

So far we have only described spuriously harmful exposures, like shower curtains. Actual harmful exposures include things like tobacco, a high cholesterol diet and alcohol abuse. Of course, some exposures are healthy. Can you think of some?

Sure. Regular exercise, wearing seat belts and certain dietary habits are beneficial.

This is a comic-style page presenting a dialogue about cohort studies. The following is the text from each panel.

**Panel 1 (top left):**

Good. I think we know enough about cohort studies to evaluate a real one. We will look at a classic; Richard Doll and Austin Bradford Hill's 1956 paper "Lung Cancer and Other Causes of Death in Relation to Smoking". What is the first thing you would do if you wanted to perform a cohort study to determine if there was an association between smoking and lung cancer death?

Find a population at risk?



More Doctors smoke Camels than any other cigarette

**Panel 2 (top right):**

Correct. What population characteristics would make the trial easier to complete?

Since determining the presence or absence of the exposure and the outcome is crucial to the success of a cohort study, we would like to use a group where we can easily establish these things.

Right. You wouldn't choose, say, the crew of a tramp steamer traveling all over the world.



14

**Panel 3 (bottom left):**

Doll and Hill chose British doctors as their population. All 59,600 doctors in the UK in 1951 were listed in the Medical Register. How do you think Doll and Hill established whether or not the doctors smoked?
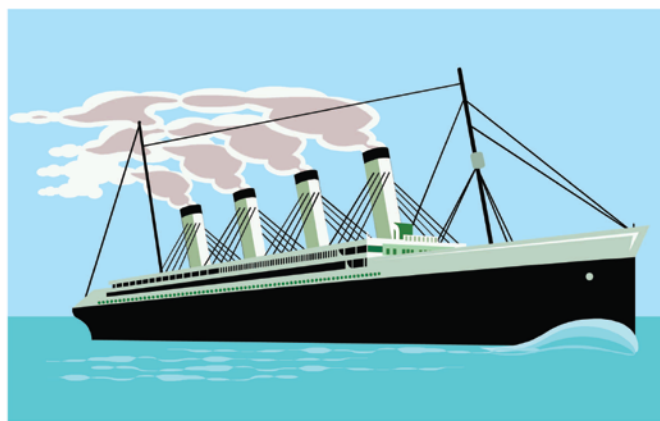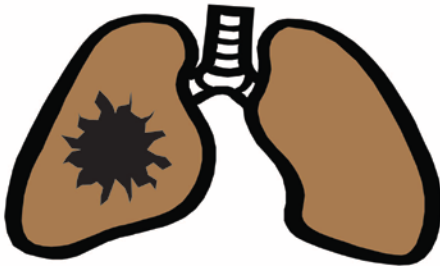
By asking them.



20,679* Physicians say "LUCKIES are less irritating"

**Panel 4 (bottom right):**

Correct. They sent out a simple "questionary" to all British doctors, asking docs to quantify the amount that they smoked. Respondents were classified as smokers and non-smokers. Smokers were further classified as former, light, moderate or heavy smokers. In a cohort study, it is crucial to correctly establish who is exposed. If an exposure is not socially accepted, it may be underreported; for example drug abusers may deny their habit. If exposure is misclassified, study results are unreliable.



So Doll and Hill compared outcomes between the smokers and non-smokers?
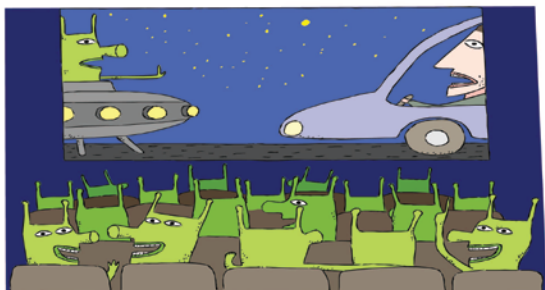
15

So Doll and Hill compared how often lung cancer deaths occurred in smoking versus non-smoking male British doctors over age 35. But how did they know that smokers and non-smokers were otherwise the same?

Good question. They could not have known for sure, but it seems a reasonable assumption since docs in the UK in 1956 would have been a fairly homogenous ethnic and socioeconomic group. We'll come back to this issue and see how the researchers dealt with one objection to this assumption later on. Unfortunately, a perfect cohort study is impossible; to do one, you would need to be able to travel between parallel universes.

What do you mean?

Let's say I take up smoking on my 40th birthday and you follow me for 30 years to determine the outcome. Who would be the ideal unexposed person to compare with my smoking self?

A version of you that does not smoke.

Exactly. The best unexposed individual to compare with an exposed individual is an unexposed version of the exposed person. Of course, that is impossible; we cannot simultaneously be smokers and non-smokers. This impossible situation has been called "counterfactual". Because we must do our study in one universe, our counterfactuals will be imperfect and we can never guarantee that our exposed and unexposed groups are otherwise identical.

So what happened with our British Doctors?

Doll and Hill followed the docs for 53 months. There were a total of 84 lung cancer deaths, one among the non-smokers and 83 among the smokers.
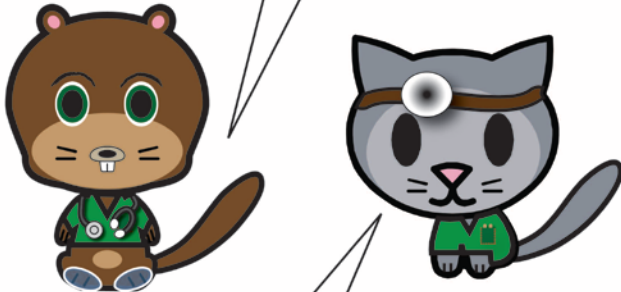
Wow, seems like a slam dunk.

Not so fast, let's look at the data in table form and see. What do you think?

That's a lot of person years of smoking!

|  | Non-Smokers | Smokers |
|---|---|---|
| Lung Cancer Deaths | 1 | 83 |
| Person Years | 15,107 | 98,090 |

**Panel 1:**

What about the cases where the cause of death was not certified as lung cancer?

They obtained no additional follow up in these cases.

Isn't that a problem? Perhaps 12 non-smokers who died of lung cancer are hidden among the docs who were classified as dying of another disease.

**Panel 2:**

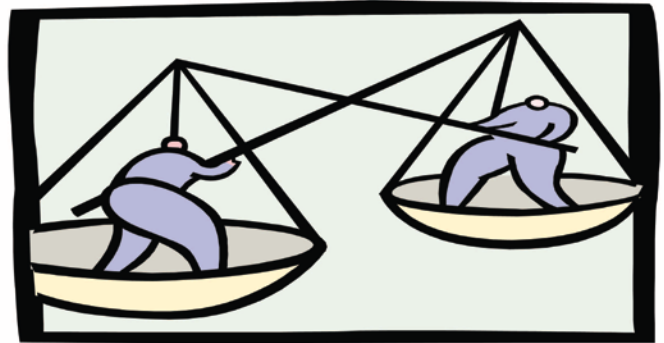True, they should have verified the cause of death in all cases, not just among people dying of lung cancer. Doll and Hill treated all non-lung cancer deaths the same way; they accepted the certified cause of death recorded in the relevant database. All lung cancer deaths underwent additional investigation to confirm the diagnosis. Any errors with respect to cause of death would be equally distributed among both the dead smokers and non-smokers. This type of misclassification is called non-differential.
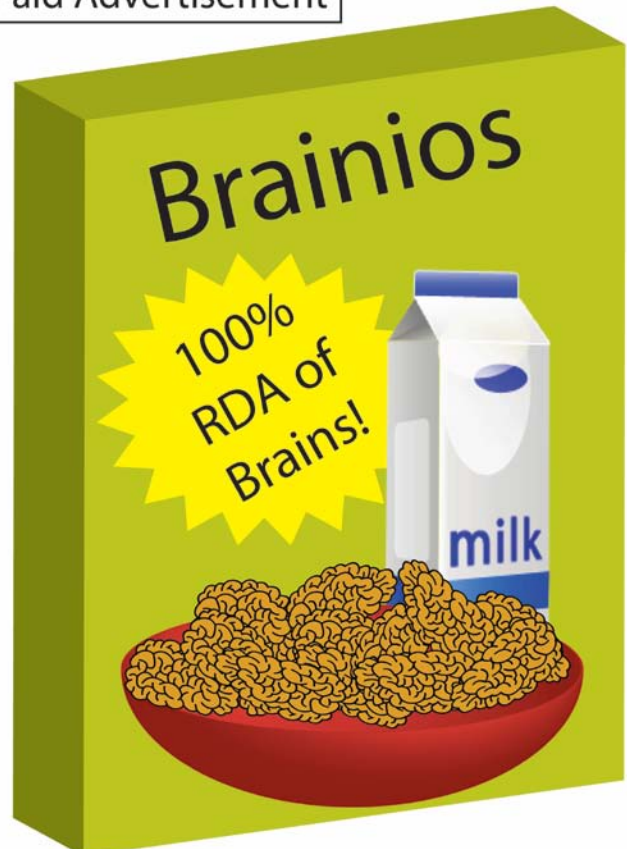
18

**Panel 3:**

How would non-differential misclassification effect your study?

Non-differential misclassification results in underestimation of the effect of the exposure; this type of error makes it harder to show that the null hypothesis (no effect of the exposure) is false.
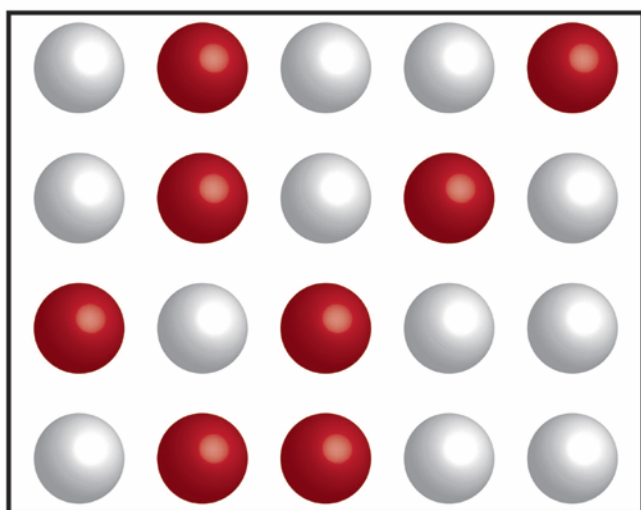
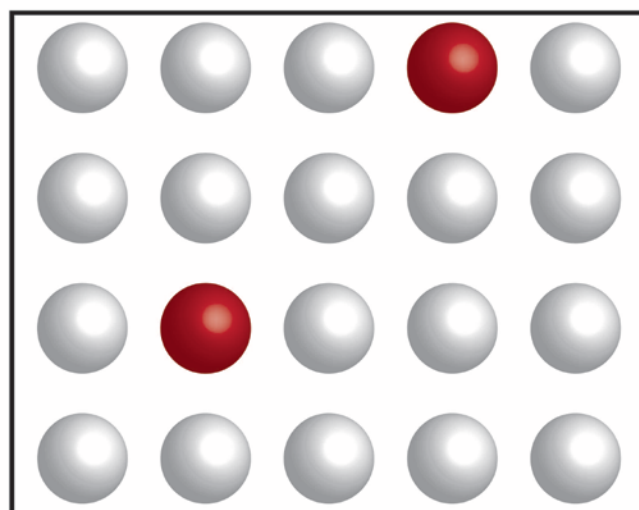I think I need an example.

**Panel 4:**

Let's say you are doing a cohort study and have an exposed and an unexposed group. Healthy people are represented by the white balls and diseased people are represented by the red balls. The true health status of each group is shown in the top two boxes: 8 out of 20 exposed are diseased while 2 out of 20 unexposed are diseased. Clearly, there is an association between the exposure and the disease. If there is non-differential misclasification, the association is harder to recognize. The two bottom boxes shows what the groups look like if we mistakenly classify 33% of the healthy people in each group as diseased. With non-differential misclassification, the groups look more alike; now we classify 12 of 20 exposed and 8 of 20 non-exposed as diseased.
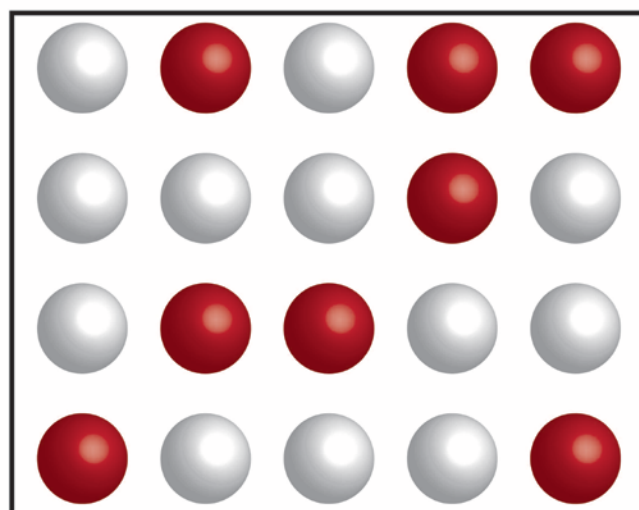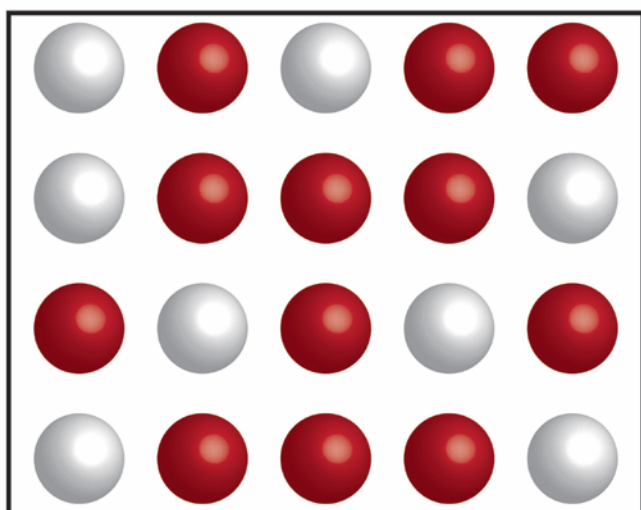
Exposed: True Status
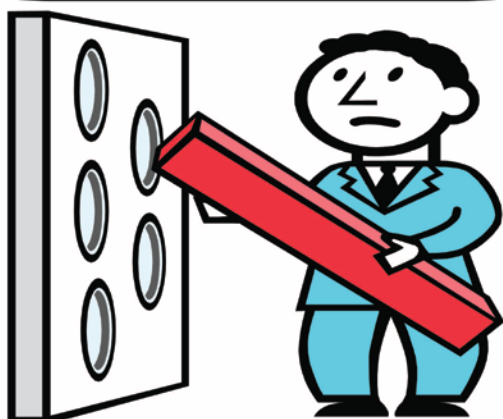
Unexposed: True Status

19

Apparent Status
after Non-differential
Misclassification

Doll and Hill presented evidence that showed that there was no differential misclassification that led to overdiagnosis of lung cancer among smokers. Although imperfect, this cohort study provided powerful evidence for the association of smoking with lung cancer. Doll and Hill also found that there was a dose response, with much higher lung cancer mortality among heavier smokers than light or moderate smokers. In addition, ex-smokers had decreased lung cancer deaths compared to present smokers.

That reminds me of Hill's criteria for establishing a causal relationship between an exposure and an outcome. This cohort study demonstrated at least five of his criteria: 1) Strength of Association, 2) Specificity, 3) Temporal Relationship, 4) Biological Gradient and 5) Experiment.

Doll and Hill were incredibly thorough and addressed multiple potential objections to their conclusions. For example, they investigated whether the difference in lung cancer mortality was due to atmospheric pollution and not smoking. Their analysis showed that more non-smokers and fewer heavy smokers lived in big cities, evidence against an atmospheric cause.

Those miasmists are as unkillable as zombies! By the way, how did Doll and Hill ever decide to undertake such a huge, expensive and time-consuming study to examine an exposure that many people thought was harmless? After all, most of the doctors that they studied smoked!

They had evidence from a different type of study published in 1950 that showed that smoking was harmful. This other study is called a case-control study and we will discuss it in the next chapter.

21