

L8: Midterm 1 review

What have we learned so far?

What have we learned so
far?

What have we learned so far?

Key topics covered on midterm :

- ▶ Course overview, PPDAC
- ▶ Describing data structure, visualization(bar chart and histogram)
- ▶ Describing data with numbers
- ▶ Scatterplots and correlation
- ▶ R basics working with data and visualizing data
- ▶ Intro to regression
- ▶ Two categorical variables
- ▶ Samples and observational studies
- ▶ Designing Experiments and Ethics

What have we learned so far?

Problem, Plan, Data, Analysis, Conclusion

You should know these terms.

You should also be familiar with the three types of problems we talked about and be able to identify them.

Types of problems

Q1: We are interested in projecting the number of immunization doses that will be needed in a clinic during the month of November based on the previous year's data.

Q2: We are creating a visualization of the mean exam scores in PH142 by program and year of student.

Q3: We are planning an intervention to reduce e-cigarette use and assessing the role of exposure to advertisements in e-cigarette use among young people.

What have we learned so far?

R basics working with data and visualizing data

You should be able to recognize what a code chunk is doing, and be able to fill in blanks in a code chunk for syntax we have seen multiple times in lecture, lab and problem sets.

What have we learned so far?

What function would i use to do the following:

restrict my dataset to a smaller number of variables?

restrict my dataset to a smaller set of observations?

check to see how many observations are in the dataset?

Create a new variable?

types of variables

What have we learned so far?

Identifying the unit of analysis

Differentiating between the types of variables

Creating visual summaries of variables

- Make bar charts using `ggplot()`'s `geom_bar()`
- Make histograms using `ggplot()`'s `geom_histogram()`
- Make scatterplots using `ggplot()`'s `geom_point()`

from Iverson et al. Abstract Exposure to industrial solvents has been associated with encephalopathy. Styrene is a neurotoxic industrial solvent, and we investigated the long-term risk of encephalopathy and unspecified dementia following styrene exposure. We followed 72,465 workers in the reinforced plastics industry in Denmark (1977–2011) and identified incident cases of encephalopathy ($n = 228$) and unspecified dementia ($n = 565$) in national registers. Individual styrene exposure levels were modeled from information on occupation, measurements of work place styrene levels, product, process, and years of employment.

check your knowledge

What have we learned so far?

What type of problem is being addressed here (descriptive, causative, predictive?) What type of variable is the outcome? What kind of a study design is this?

1. Investigate measures of centrality
 - ▶ mean and median, and when they're the same vs. different
2. Investigate measures of spread
 - ▶ IQR, standard deviation, and variance
3. Create a visualization of the “five number summary”
 - ▶ boxplots using `ggplot`
4. Calculate the variance and standard deviation

What have we learned so far?

- ▶ Determine which variable is explanatory and which is response, or when it doesn't matter
- ▶ Visually describe the relationship between two variables (form, direction, strength, and outliers)
- ▶ Numerically describe the relationship with the correlation coefficient r

What have we learned so far?

Be able to pull relevant pieces of information from r output and interpret them including - intercept - slope - r squared

Know how to use the equation to find predicted values of Y at a given X

Equation of the line of best fit

The line of best fit can be represented by the equation for a line:

$$y = a + bx$$

where a is the **intercept** and b is the **slope**.

Linear regression

Below are images of the output from code that will run a linear regression model on the relationship between age and charges for the subset population and produce the following outputs. Here the population is subset to smokers who are of normal BMI.

We will assign the regression model to the name `insure_model`. Write the missing commands to fill in blanks in the code

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  10609.    1325.     8.01 2.14e-10
## 2 age          244.     33.0     7.41 1.74e- 9
```

```
insure_model <- ____(formula = charges ~ age, data = insure_subset)
```

What have we learned so far?

Write the equation for the line of best fit for the subset data. Interpret the slope parameter, in one sentence what does the slope parameter tell you.

Using the model, predict the medical charges for a smoker of normal BMI who is 30 years old.

A relationship between your variable of interest (exposure, treatment) and your outcome of interest (disease status, health condition etc) is confounded when there is a variable that is associated with both the exposure and outcome, and is not on the causal pathway between the two.

Variables that are on the causal pathway are those that represent a way in which the exposure acts on the outcome.

What have we learned so far?

Two categorical variables

- ▶ Two way tables
 - ▶ marginal vs conditional distributions
- ▶ Bar graphs
 - ▶ side by side
 - ▶ stacked
- ▶ Simpson's paradox

Suppose I am a teacher in elementary school and there is a head lice outbreak at my school. I know that some of the children went on a field trip recently to a fire station. While they were at the fire station they were allowed to pass around a fire helmet and try it on. I suspect that this caused the outbreak. I collect the following data:

Group	head lice	no head lice
Field trip	49	62
No Field trip	15	89

Calculate the conditional probability of head lice among students who attended the field trip, and among those who did not.

What have we learned so far?

Why are we interested in the conditional rather than the marginal probabilities here?

What do these data suggest?

- ▶ Whether the treatment or exposure is controlled by an investigator
 - ▶ Experimental vs observational designs
- ▶ The population of interest
 - ▶ Target population
 - ▶ Study population
- ▶ How the sample was drawn from the population
 - ▶ Complete sample (census)
 - ▶ Random sampling
 - ▶ Convenience sampling
 - ▶ Volunteer sampling
- ▶ Was selection conditional on exposure or outcome

What have we learned so far?

- ▶ Understanding different types of experimental designs
- ▶ Thinking about sources of bias
 - ▶ bias from the design or conduct of sampling
 - ▶ bias from lack of adherence to protocol
 - ▶ bias in assessment
 - ▶ bias in analysis
- ▶ Ethics in randomization