

Discussion #2 Student Solutions

Name:

This discussion consists of a quick recap of sampling methods and biases, some tips and examples of how to calculate probabilities, and some SQL.

As stated in Lecture 1, some time in every discussion will be spent on selected homework problems.

Sampling and Bias

1. A campus organization wants to take a sample of Berkeley students who are registered for classes this semester. To do this, the organization takes a simple random sample of 20 classes from among all classes offered this semester, and then takes all students in those classes. You can assume that the organization has access to complete enrollment information all classes.

(a) Is this a simple random sample of students? Explain.

Solution: No. For example, a student who is taking 4 classes is more likely to get in than a student who is taking 3.

(b) Is this a probability sample of students? Explain.

Solution: Yes. The organization knows the total number of classes, and for each student it knows how many classes the student is taking. So it's possible to find the chance that the student gets in.

This is called a *cluster* sample. Each class is a cluster of students. The sample is an SRS of clusters, not an SRS of students.

2. The Current Population Survey is a national survey run by the Census Bureau. It is thorough and reliable, and thus is sometimes used as a benchmark to assess the accuracy of other surveys. As part of an assessment of its own phone surveys, the Pew Research Center found that the response rates have been dropping over the years. Still, on most measures, its estimates were comparable to those of the Current Population Survey. But for example 55% of respondents in the most recent Pew Survey said they did some type of volunteer work for or through an organization in the past year, compared with 27% in the Current Population Survey.

How do you think this difference might have arisen?

Solution: Non-respondents are different from respondents. For example, the qualities that motivate people to do volunteer work may be positively associated with their willingness to respond to phone surveys.

Finding Chances

Golden rules for finding the chance of an event:

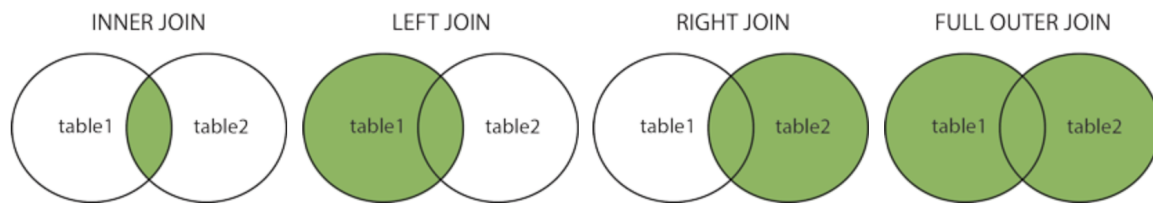
- List the ways: list all the distinct ways the event can happen, and add the chances of all the ways.
 - If the list above looks long and complicated, make the list of ways in which the event *doesn't* happen; it might be simpler.
 - If an event involves multiple trials, like a number of random draws, imagine yourself conducting the experiment one trial at a time.
3. Let n be a positive integer. Consider a sample of size n drawn at random with replacement from a population in which a proportion p of the individuals are called successes.
- (a) For an integer k such that $0 \leq k \leq n$, which of the following are equal to the chance of getting exactly k successes in the sample?
- (i) $p^k(1-p)^{n-k}$
 - (ii) $\binom{n}{k}p^k(1-p)^{n-k}$
 - (iii) $\binom{n}{n-k}p^k(1-p)^{n-k}$
 - (iv) $\frac{n!}{k!(n-k)!}p^k(1-p)^{n-k}$

Solution: All except (i).

- (b) Which of the following are equal to the chance of getting at least one success in the sample?
- (i) $np(1-p)^{n-1}$
 - (ii) $\sum_{k=2}^n \binom{n}{k}p^k(1-p)^{n-k}$
 - (iii) $\sum_{k=1}^n \binom{n}{k}p^k(1-p)^{n-k}$
 - (iv) $1-p^n$
 - (v) $1-(1-p)^n$

Solution: (iii) and (v)

SQL



Note: You do not always have to use the JOIN keyword to join sql tables. The following are equivalent:

```
SELECT column1, column2
FROM table1, table2
WHERE table1.id = table2.id;
```

```
SELECT column1, column2
FROM table1 JOIN table2
ON table1.id = table2.id;
```

- Describe which records are returned from each type of join in the figure above. How does a cross join relate to these types of joins?

Solution:

(INNER) JOIN: Returns records that have matching values in both tables

LEFT (OUTER) JOIN: Return all records from the left table, and the matched records from the right table

RIGHT (OUTER) JOIN: Return all records from the right table, and the matched records from the left table

FULL (OUTER) JOIN: Return all records when there is a match in either left or right table
CROSS JOIN: Return all pairs of records between the left and right tables.

- Consider the following real estate schema:

```
Homes(home_id int, city text, bedrooms int, bathrooms int,
area int)
Transactions(home_id int, buyer_id int, seller_id int,
transaction_date date, sale_price int)
```

Buyers(buyer_id int, name text)
Sellers(seller_id int, name text)

Fill in the blanks in the SQL query to find the id and selling price for each home in Berkeley. If the home has not been sold yet, **the price should be NULL**.

```
SELECT  _____ H.home_id, T.sale_price _____  
FROM    _____ Homes H _____  
_____ LEFT OUTER _____ JOIN _____ Transactions T _____  
ON      _____ H.home_id = T.home_id _____  
WHERE   _____ H.city = 'Berkeley' _____;
```

Solution: An alternate solution was to use Transactions in the FROM clause and perform a RIGHT OUTER JOIN with Homes.

```
SELECT H.home_id, T.sale_price  
FROM Transactions T  
RIGHT OUTER JOIN Homes H  
ON H.home_id=T.home_id  
WHERE H.city = 'Berkeley'
```