

# Social Data Science

## **Course Recap**

# Goals of the course

- ▶ connect theories about inequality to quantitative empirical evidence
- ▶ evaluate the effects of hypothetical interventions to reduce inequality
- ▶ conduct data analysis using the R programming language

# Structure of the Course

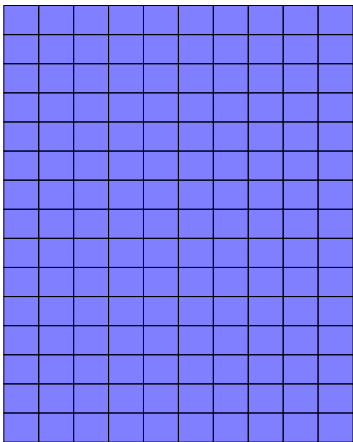
We learned key concepts in social data science.

1. Population Sampling
2. Models for Subgroup Summaries
3. Causal Inference

# Population Sampling

**Full Count  
Enumeration**

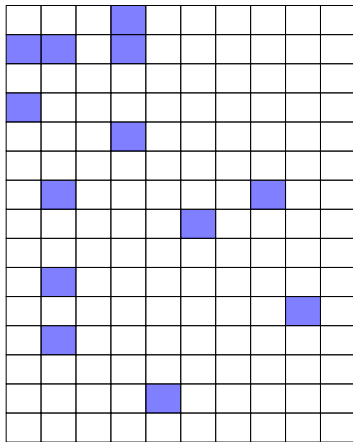
Back of Room



Front of Room

**Probability  
Sample**

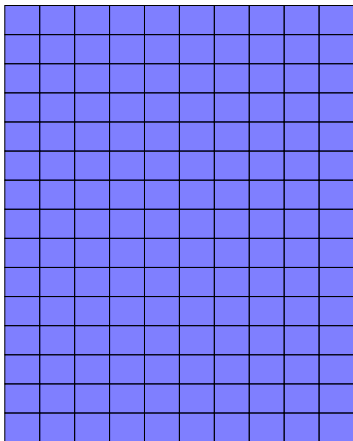
Back of Room



Front of Room

# Full Count Enumeration

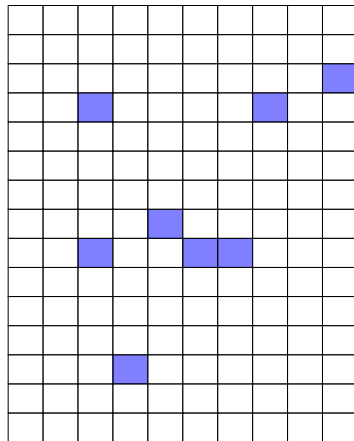
Back of Room



Front of Room

# Probability Sample

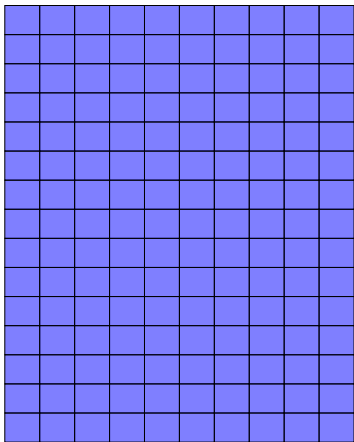
Back of Room



Front of Room

**Full Count  
Enumeration**

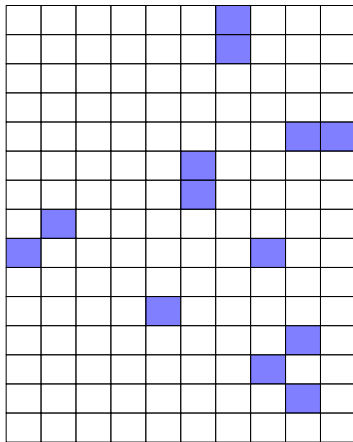
Back of Room



Front of Room

**Probability  
Sample**

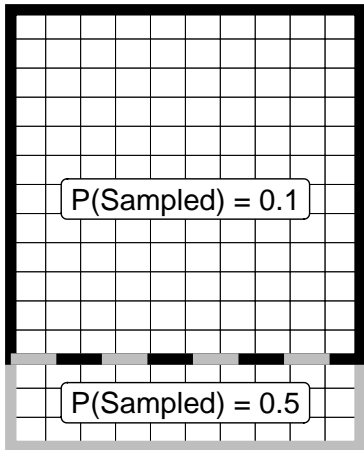
Back of Room



Front of Room

## Sample Design

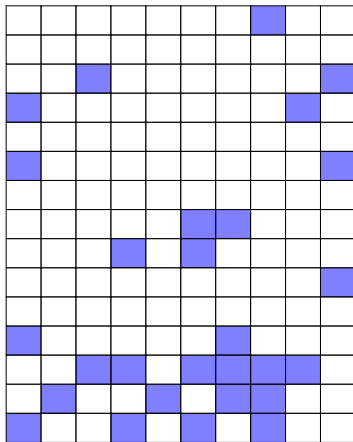
Back of Room



Front of Room

## Sample

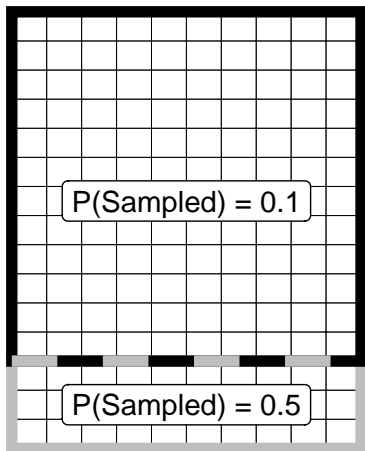
Back of Room



Front of Room

## Sample Design

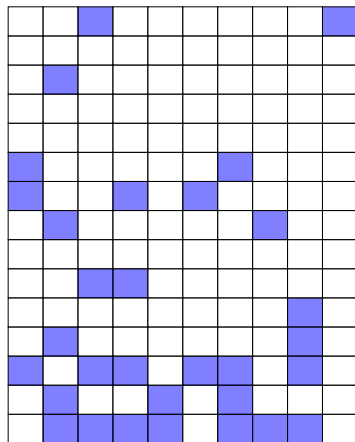
Back of Room



Front of Room

## Sample

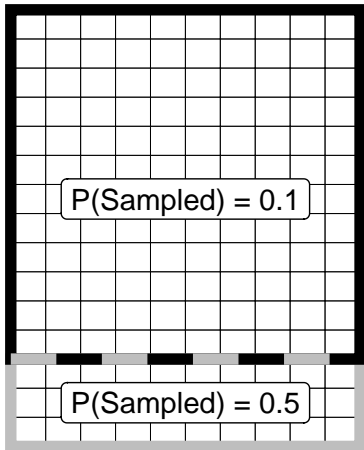
Back of Room



Front of Room

## Sample Design

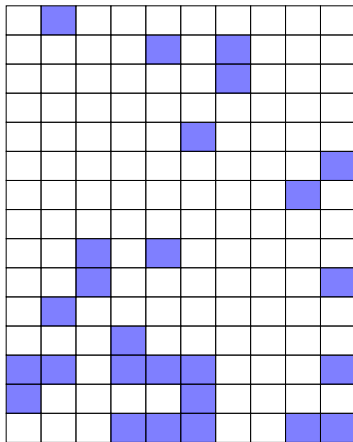
Back of Room



Front of Room

## Sample

Back of Room



Front of Room

# Working with data

- ▶ Creating objects in R
- ▶ Piping data into functions
- ▶ Summarizing a sample
- ▶ Visualizing with ggplot

## Working with data

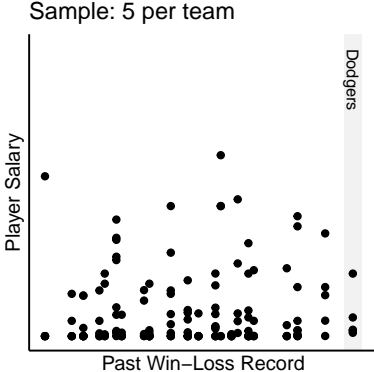
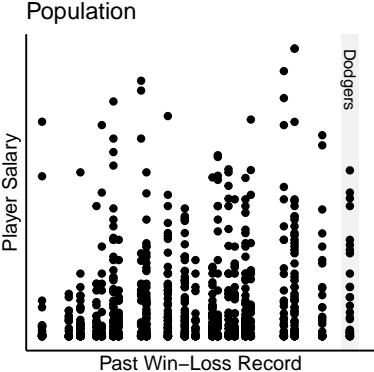
Questions of the form:

Among -----, what is the mean of -----?

(population inference from a sample)

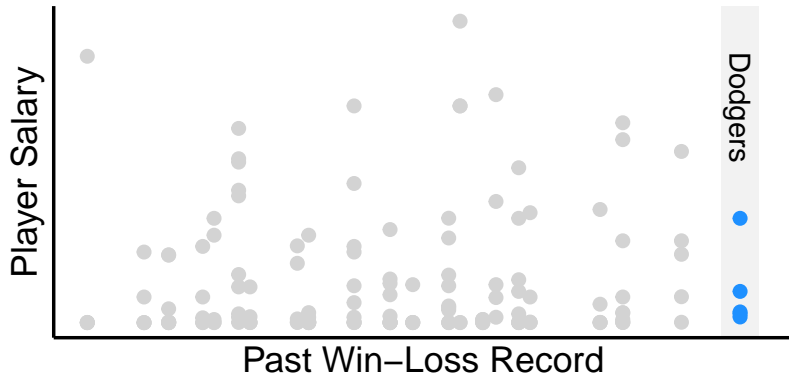
# Models for Subgroup Summaries

With only the sample, estimate the mean salary of the Dodgers



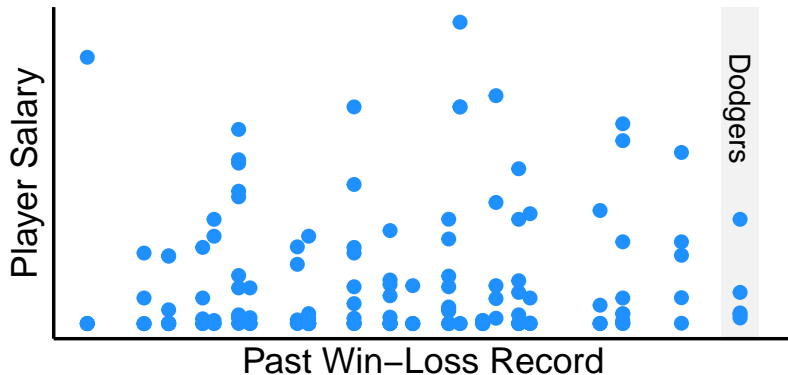
# Working with Data: Statistical Learning

## Estimator 1: Subgroup sample mean



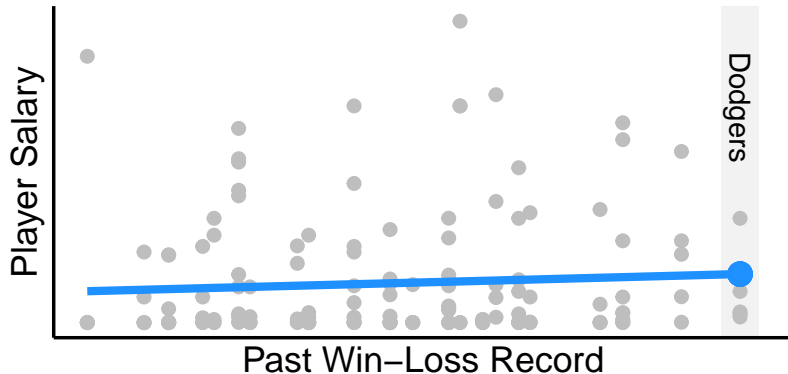
# Working with Data: Statistical Learning

## Estimator 2: Full sample mean

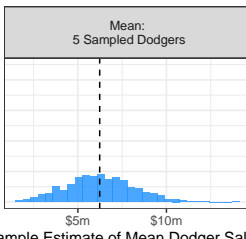
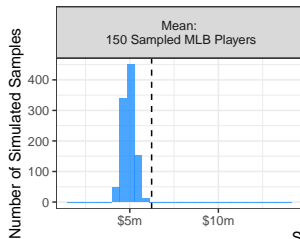
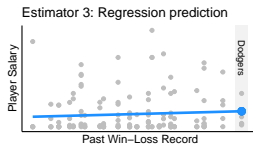
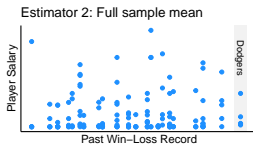
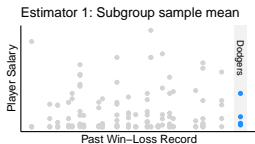


# Working with Data: Statistical Learning

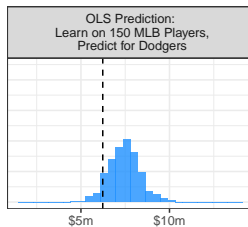
## Estimator 3: Regression prediction



# Working with Data: Statistical Learning



Sample Estimate of Mean Dodger Salary




# Working with Data: Statistical Learning

Learning Set

	Respondent Income	Respondent Education	Parent Education	Grandparent Education	Sex	Race	Grandparent Income	Parent Income
Case 1								
Case 2								
Case 3								
Case 4								
Case 5								

Learn a prediction function

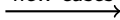


Respondent Income

Holdout Set

Case 6								
Case 7								
Case 8								

Predict for new cases



?
?
?

Causal inference



# Fundamental problem of causal inference

Holland 1986

Descriptive evidence



Causal claim



Causal inference is a **missing data** problem

Person 1	lifespan	missing	lifespan	lifespan
Person 2	missing	lifespan	lifespan	lifespan
Person 3	lifespan	missing	lifespan	lifespan
Person 4	missing	lifespan	lifespan	lifespan
Person 5	lifespan	missing	lifespan	lifespan
Person 6	lifespan	missing	lifespan	lifespan
Person 7	missing	lifespan	lifespan	lifespan
Person 8	lifespan	missing	lifespan	lifespan
	Outcome under Mediterranean diet	Outcome under standard diet	Outcome under Mediterranean diet	Outcome under standard diet

## Potential outcomes

$$Y_i^a$$

the outcome  $Y$   
of person  $i$   
if exposed to treatment  $A = a$

## Causal identification with DAGs



**Non-causal path** starts with an edge pointing in to  $A$  and ends at  $Y$

Identified causal effects by an adjustment set that blocks all non-causal paths.

# Estimation

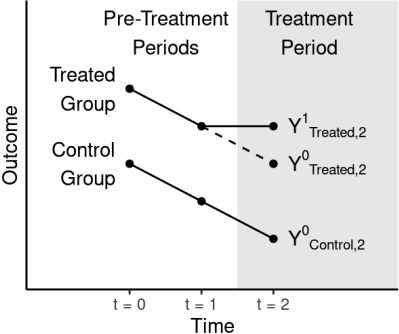
- ▶ regression
- ▶ inverse probability weighting
- ▶ matching

# Causal inference with unmeasured confounding

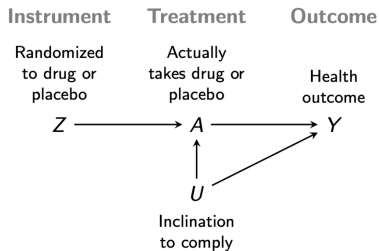
# Causal inference with unmeasured confounding



# Causal inference with unmeasured confounding



# Causal inference with unmeasured confounding



# Course structure

- ▶ concepts introduced in lecture
- ▶ conceptual practice with quizzes
- ▶ coding practice with problem sets
- ▶ support in office hours and on Piazza

# Goals of the course

- ▶ connect theories about inequality to quantitative empirical evidence
- ▶ evaluate the effects of hypothetical interventions to reduce inequality
- ▶ conduct data analysis using the R programming language

# Your thoughts

- ▶ What could we do to make this course better?
- ▶ What is your favorite thing you learned?
- ▶ What parts do you anticipate being most useful for your future work?

# Course evaluations

- ▶ Specific examples are especially helpful
  - ▶ Comments on course content
  - ▶ Comments on course delivery
- ▶ Specific positive experiences with your TA!