



**CS224C: NLP for CSS**

# Casual Inference

Diyi Yang  
Stanford CS

# Lecture Overview

- ◆ Prediction vs. Understanding
- ◆ Randomized controlled trial (RCT)
- ◆ Observation data and studies
- ◆ Propensity score methods

# Prediction vs. Understanding

*Statistical Science*  
2010, Vol. 25, No. 3, 289–310  
DOI: 10.1214/10-STS330  
© Institute of Mathematical Statistics, 2010

## To Explain or to Predict?

Galit Shmueli

Two main uses of statistical models:

**Prediction:** inferring the most likely values (+ prediction intervals) for data where you don't know the answer

**Understanding:** estimating the relationship between a predictor variable and some outcome (+ quantifying uncertainty about that relationship)

# Starting with Regression

Logistic regression

$$P(y = 1 | x, \beta) = \frac{\exp(\sum_{i=1}^F x_i \beta_i)}{1 + \exp(\sum_{i=1}^F x_i \beta_i)}$$

Linear regression

$$y = \sum_{i=1}^F x_i \beta_i + \epsilon$$

# Features and Coefficients

$x_i$  refers to each feature

such as "speaking English", "mentioning Clinton on Twitter"

$\beta_i$  refers to the coefficient associated with  $x_i$

# A Simple Example

$$P(y = 1 | x, \beta) = \frac{\exp(x_0\beta_0 + x_1\beta_1)}{1 + \exp(x_0\beta_0 + x_1\beta_1)}$$

$x_0$ : whether the user speaks English

$x_1$ : how many times the user mentions Clinton on Twitter

$y$ : 1 if the user votes for Clinton, otherwise 0

# A Simple Example

$$P(y = 1 | x, \beta) = \frac{\exp(x_0\beta_0 + x_1\beta_1)}{1 + \exp(x_0\beta_0 + x_1\beta_1)}$$

$$\frac{P(y = 1 | x, \beta)}{1 - P(y = 1 | x, \beta)} = \exp(x_0\beta_0 + x_1\beta_1) = \exp(x_0\beta_0) \cdot \exp(x_1\beta_1)$$

If  $x_1$  increases by 1,

$$\exp(x_0\beta_0) \cdot \exp((x_1 + 1)\beta_1) = \exp(x_0\beta_0) \cdot \exp(x_1\beta_1 + \beta_1) = \exp(x_0\beta_0) \cdot \exp(x_1\beta_1) \cdot \exp(\beta_1)$$

$\exp(\beta)$  Refers to the factor by which the odds change with a 1-unit increase in  $x$

# Interpreting the coefficient for “explanation”

We can assess how significant is the relationship between a predictor and its outcome (aka correlations) with a hypothesis test

But are these reliable?

Can we add control variables?

Refined correlations!

# Correlation vs. Causation

Understand the causal relationship of a treatment Z on some outcome Y

Treatment	Outcome
take a drug	cured of disease
graduate high school	earnings
cast John Goodman	box office
living in Berkeley	political preference

Slides Credit to David Bamman

# Terminology

Treatment:  $Z(0), Z(1)$

The predictor variable whose causal relationship we're interested in

Potential outcomes:  $Y=0, Y=1$

The dependent variable

We're interested in the causal relationship between the treatment  $Z$  and  $Y$

# Counterfactual

John doesn't brush his teeth ( $Z=0$ ) and developed heart disease ( $Y=1$ )

What would have happened if he did brush his teeth ( $Z=1$ )?

For any data point, we only ever get to observe one outcome. We never observe the counterfactual.

# Observational Data

Hypothesis tests for observational data assess the relationship between variables but don't establish causality

**Examples:** if we intervened and relocated someone to Palo Alto, would they become liberal?

# Experimental Data

Data that allows you to perform an intervention and determine the value of some variable

**Clinical data:** treatment vs. placebo

**Web design:** one or two homepage designs

**Political email campaigns:** one of two (differently worded) solicitations

A potential confound exists if any other variable is correlated with your intervention decision:

E.g., users volunteering to receive a drug (and not the placebo)

# Randomization Experiments

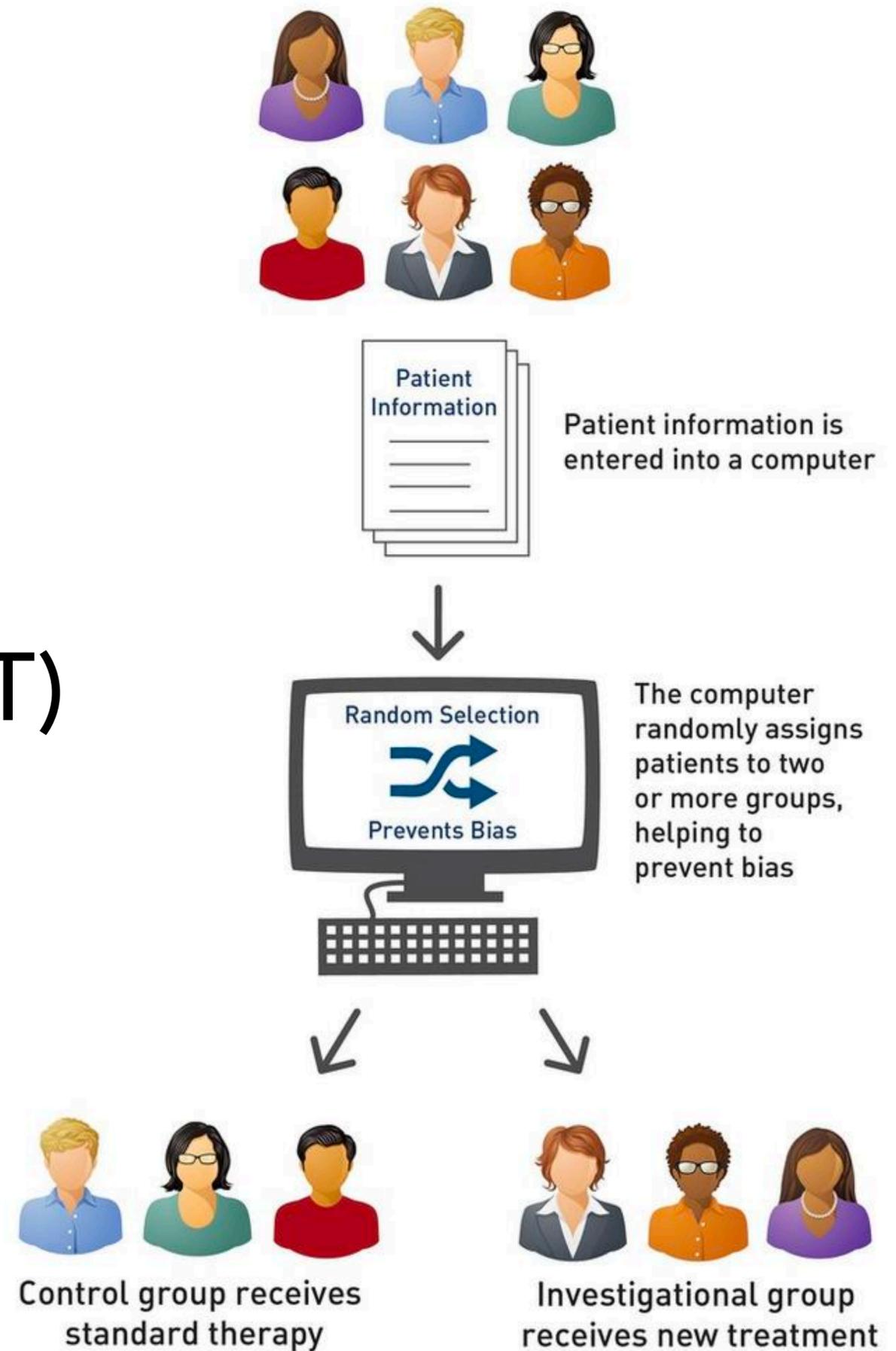
Users are randomly assigned an outcome (which web page), which allows us to better establish causality

**A/B testing** = significance test in randomized experiment with two outcomes

We can run a standard regression, but now if the  $\beta_{\text{design}_A}$  is significant, we can interpret it causally. By randomly assigning the treatment, we are ensuring that its value is uncorrelated with any other variable

# Randomized Control Trail (RCT)

<https://www.cancer.gov/about-cancer/treatment/clinical-trials/what-are-trials/randomization/clinical-trial-randomization-infographic>



# RCT Estimation

$$E[Y(1) - Y(0)] = E[Y|Z = 1] - E[Y|Z = 0]$$

RCT gives an unbiased estimate of the average effect of the treatment

# Randomization May Not Be Feasible

- Ethical Issues
- Controlled or the treatment conditions may be harmful

# Observational Data

Observational data can't be intervened to establish an causal relationship

Instead, we could:

Accounting for confounding variables

**Assume** there is a randomization experiment **lurking** in the data

# Propensity Score

**Propensity score:** the probability of treatment assignment conditional on observed baseline covariates - also called a ***balancing score***

$$e_i = Pr(Z_i = 1 | \mathbf{X}_i)$$

**In RCTs**, propensity score is known and defined by the study design.

**In observational studies**, the true propensity score is not known, but can be estimated using the study data

# Four Propensity Score Methods

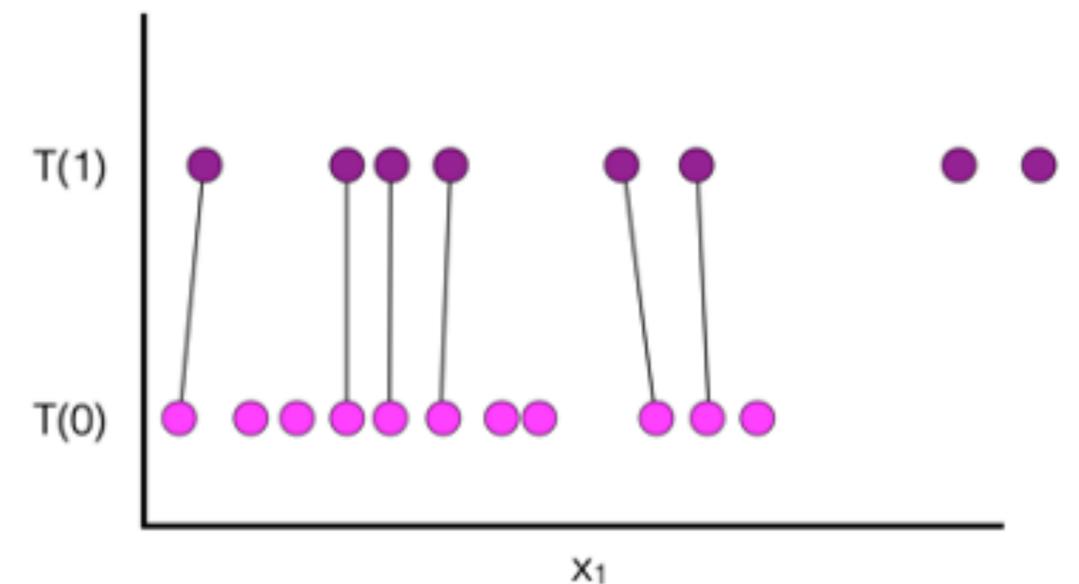
1. Propensity score matching
2. Stratification on the propensity score
3. Inverse probability of treatment weighting
4. Covariate adjustment using the propensity score

# 1 Propensity Score Matching

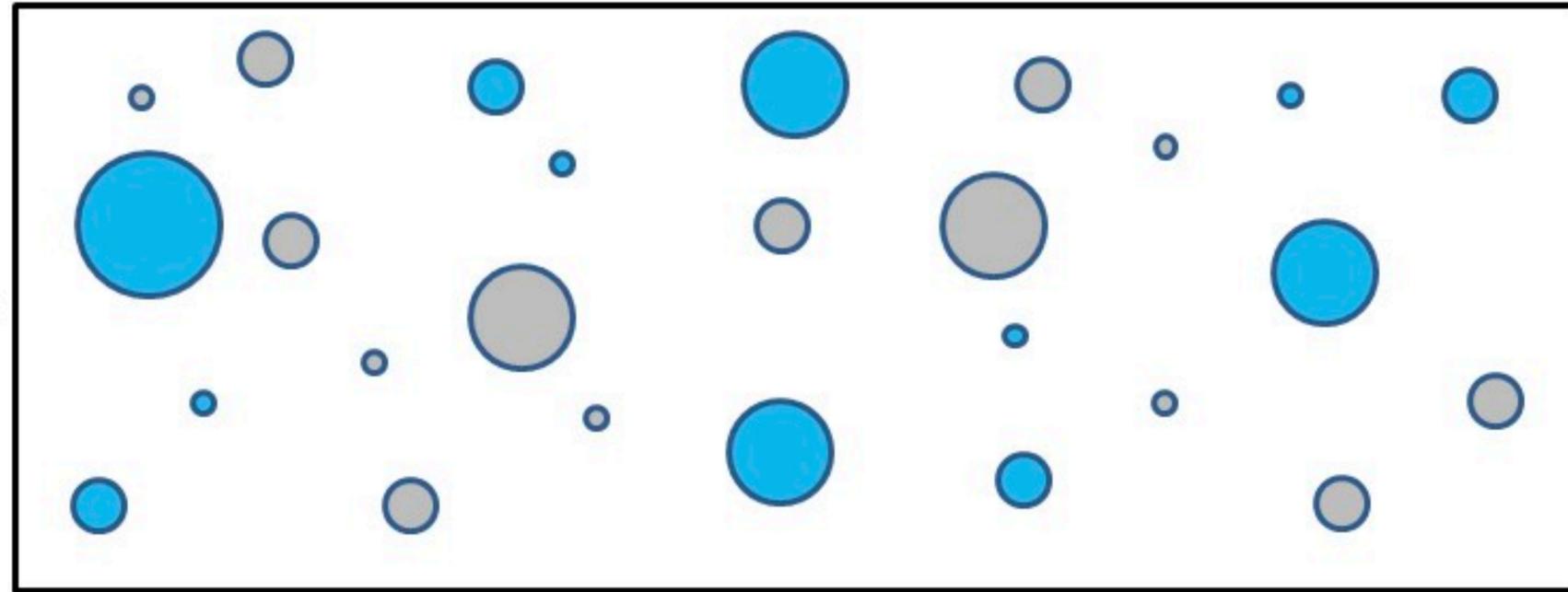
*"Form matched sets of treated and untreated subjects who share a similar value of propensity score"*

Common approach: **one-to-one or pair matching**

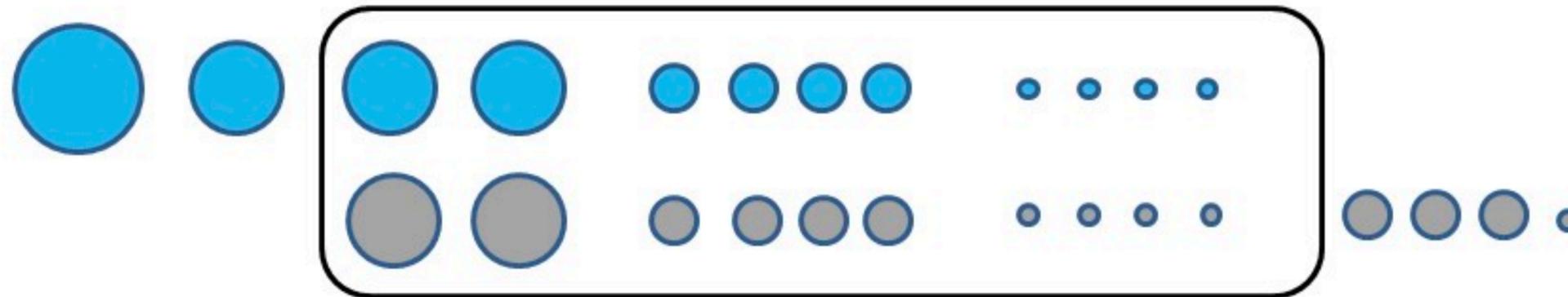
One can directly compare outcomes between treated and untreated subjects within the propensity score matched sample



Population  
with varying  
characteristics



Study Group with Matching



 Treatment  Control

# Decisions on how to form matched pairs

1. Choose between matching ***without replacement*** and ***with replacement***
2. Go with ***greedy*** or ***optimal matching***
  - ▶ **Greedy**: a treated subject is *first selected* at random, and the untreated subject whose propensity score is closest to that is chosen for matching
  - ▶ **Optimal**: matches are formed so as to minimize the total within-pair difference of the propensity score

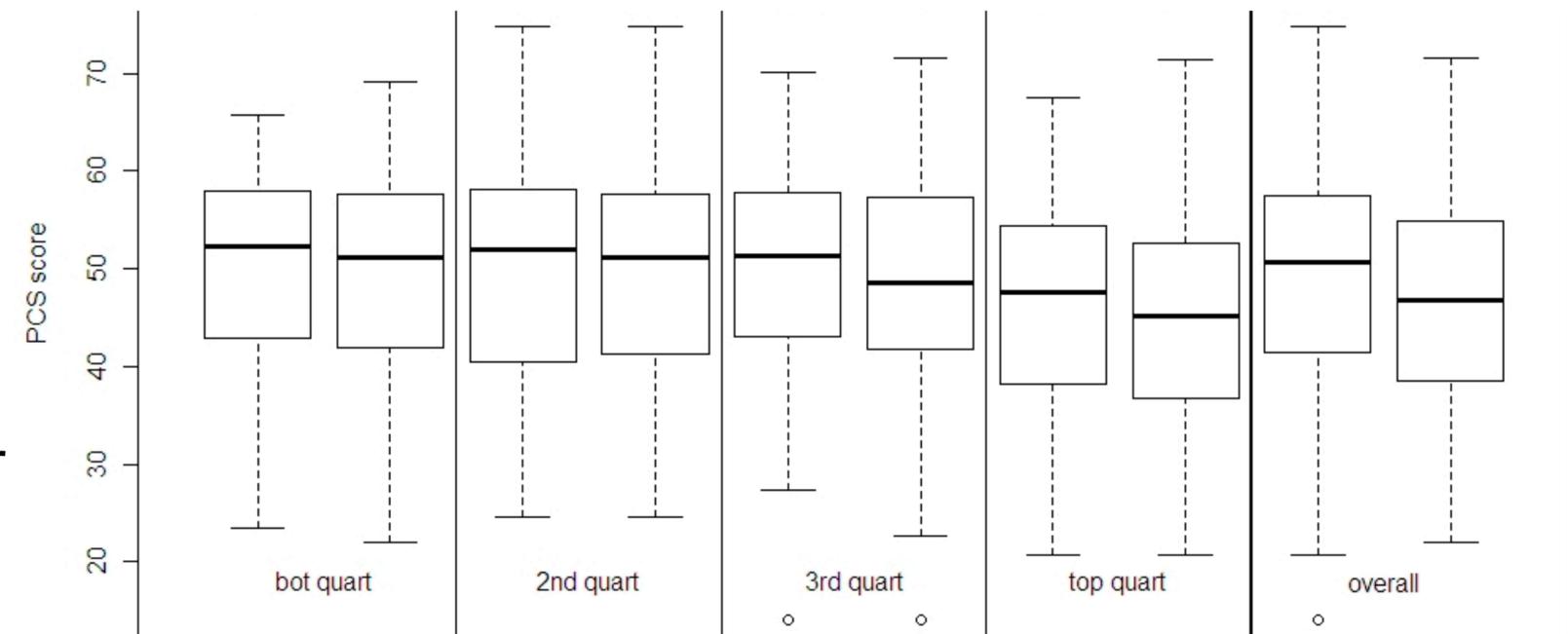
# Decisions on how to determine the “close”

Two primary methods for selecting untreated subjects whose propensity score is “close” to that of a treated subject

- ▶ Nearest neighbor matching
- ▶ Nearest neighbor matching within a specified caliper distance

## 2. Stratification on the Propensity Score

- ▶ Stratify subjects into mutually exclusive subsets based on the rank-ordered propensity score.
- ▶ Overall treatment effect is pooled over stratum-specific treatment effects – a *meta-analysis of a set of quasi-RCTs*



<http://sas-and-r.blogspot.com/2010/05/example-736-propensity-score.html>

# 3 Inverse Probability of Treat Weighting (IPTW)

IPTW using the propensity score creates weights based on the probability score to create a synthetic dataset

$$w_i = \frac{Z_i}{e_i} + \frac{(1 - Z_i)}{1 - e_i}$$

Aka, the *inverse probability* of the treatment received.

This is very similar to the use of survey sampling weights to weigh survey samples so that they are representative of specific populations.

# 4 Regression Adjustment using Propensity Score

Outcome is regressed on an indicator of the treatment status and the estimated propensity scores.

*Continuous outcome: linear models*

*Dichotomous outcome: logistic regression*

The effect of treatment is determined using the estimated regression coefficient from the fitted regression model.

# Comparing Different Propensity Score Methods

## ***The shared goal***

To remove confounding so that the treatment condition is independent of baseline characteristics between treated and untreated subjects

## ***Differences:***

Matching, stratification and weighting separate the design of the study from the analysis of the study, while regression requires both the propensity score and the outcome to be in the same model

Different tolerance to sensitivity

Primary study analysis method	 Pros	 Cons
Traditional covariate adjustment	<ul style="list-style-type: none"> <li>• Performed well</li> <li>• Provides prognostic model for outcome of interest</li> </ul>	<ul style="list-style-type: none"> <li>• May not be suitable with many covariates in smaller studies</li> </ul>
Propensity score (PS) stratification	<ul style="list-style-type: none"> <li>• Retains data from all study participants</li> <li>• Opportunity to explore interactions between treatment and PS on outcome risk</li> <li>• Provides effect estimates for every stratum</li> </ul>	<ul style="list-style-type: none"> <li>• Performs less well in datasets with few outcomes, particularly when the number of strata is large</li> <li>• May not account for strong confounding</li> </ul>
PS matching	<ul style="list-style-type: none"> <li>• Reliable; provides excellent covariate balance in most circumstances</li> <li>• Simple to analyze, present and interpret</li> </ul>	<ul style="list-style-type: none"> <li>• Some patients are unmatched leading to information excluded from the analysis</li> <li>• Less precise</li> </ul>
PS inverse probability weighting	<ul style="list-style-type: none"> <li>• Retains data from all study participants</li> <li>• Easy to implement</li> <li>• Creates a pseudo population with perfect covariate balance</li> </ul>	<ul style="list-style-type: none"> <li>• Can be unstable when extreme weights occur</li> </ul>
PS covariate adjustment (use of PS as a covariate)	<ul style="list-style-type: none"> <li>• Performed well</li> </ul>	<ul style="list-style-type: none"> <li>• Adding the PS as an additional covariate produced results very similar (and not necessarily superior) to traditional covariate adjustment</li> </ul>

# Balance Diagnostics

*"The true propensity score is a balancing score"*

**Standardized differences** to compare the similarity of treated and untreated subjects in the matched samples

For continuous variables: 
$$d = \frac{(\bar{x}_{\text{treatment}} - \bar{x}_{\text{control}})}{\sqrt{\frac{s^2_{\text{treatment}} + s^2_{\text{control}}}{2}}}$$