

# EHR Working Group Workshop

## Practical: Propensity scores

10<sup>th</sup> July 2024

In this practical, we will explore the relationship between treatment (*exposure*) and outcome using propensity score approaches for the simulated cohort data.

**Exercise 1** Read in the analysis dataset (cohort.dta) and explore the data.

For today's session, since we wish to focus our attention on the estimation of causal effects we will remove patients with missing data. In real life, for variables that contain missing values we would instead use a range of techniques to handle the missing data.

```
library(tidyverse)
```

```
library(MatchIt)
```

```
library(survey)
```

```
library(cobalt)
```

```
library(tableone)
```

```
library(haven)
```

```
data <- read_dta(file = "cohort.dta")
```

**Exercise 2** Estimate the effect of treatment on the outcome

```
model <- glm(trt ~ outcome, family=binomial(), data=data)
```

Repeat, adjusting for covariates

```
adjusted_model <- glm(trt ~ outcome + age + female+ ses+ smoke + alc + bmicat + nsaid_rx  
+ cancer + hyper, family=binomial(), data=data)
```

**Exercise 3** Estimate the propensity score

As described in the earlier session, the propensity score is defined as the probability of exposure/treatment assignment conditional on measured (or observed) characteristics. Our aim is to identify an exposed group and a non-exposed group that are similar in terms of observed characteristics, on average. To do that, for each individual, we estimate their probability of being exposed (i.e. of initiating the study drug instead of the comparator drug). Here, we use a multiple logistic regression model to estimate those probabilities.

First, we will fit a multiple logistic regression model for being exposed to the study drug:

Fit model

```
PS_model <- glm(trt ~ age + female+ ses+ smoke + alc + bmicat + nsaid_rx + cancer +  
hyper, family=binomial(), data=data)
```

Display model coefficients

```
broom::tidy(PS_model, exponentiate=TRUE, conf.int=TRUE)
```

Using this model, we predict the probability of each individual being exposed given their observed confounders/characteristics. We will save these probabilities as a variable in the analysis dataset.

```
data$PS <- PS_model %>% predict(data, type = "response")
```

#### Exercise 4 Check distribution

After fitting the propensity score, it is important to investigate its distribution!

Ideally, we do not want the propensity score to completely explain the probability of being in the exposed or the unexposed group (i.e. we do not want scores at zero or 1). If this does occur, this would indicate that some individuals have almost no chance of being in the other exposure group, which would violate the positivity assumption (which states that each individual has a non-zero probability of being in either exposure group).

Draw a histogram of the estimated propensity scores. Are there very extreme propensity score values?

```
hist(data$PS, xlab="Propensity score", ylab= "Number of individuals", main="Propensity  
score distribution", col="darkgrey", breaks=50)
```

#### Exercise 5 Check distribution cont.

Second, we wish to check that the range of propensity score values in the two exposed groups overlaps. Such overlap is sometimes called *\*common support\**. A lack of common support (i.e. lack of overlap) would also indicate a violation of the positivity assumption. For example, if we have estimated propensity score values between 0.8 and 0.9 in the exposed group but the highest propensity score in the unexposed group is 0.78, this would indicate that we have a set of people in the exposed group who do not have anyone comparable to them in the unexposed group.

To explore the area of common support, you can use a histogram or plot a box plot stratifying by treatment.

```
boxplot(PS ~ trt, data= data, ylab="Propensity score", xlab="Treatment")
```

The R package ``ggplot2`` allows us to produce more sophisticated and publishable graphs. Also, it allows us to plot the distribution of propensity scores according to treatment in the form of a histogram. First, install the package (``install.packages("ggplot")``). Then load the package:

```
library(ggplot2)
```

And produce a histogram of the propensity score according to treatment:

```
ggplot(data, aes(x=PS, color=factor(trt), fill=factor(trt))) +  
  geom_histogram(position="identity", bins=50, alpha=0.5) +  
  labs(x = "Propensity score", y="Number of individuals")
```

Summarise the propensity score:

```
data %>%  
  group_by(trt) %>%  
  summarize(n=n(), min=min(PS), mean_PS=mean(PS), median_PS=median(PS),  
    max=max(PS))
```

Based on the results above, how concerned are you about violations of the positivity assumption?

The study drug group has estimated propensity score values higher than the comparator group (0.6545169 vs 0.6532809). However, these are so close that we wouldn't consider this concerning in practice. Similarly, the difference between the minimum values is very close.

We are not particularly concerned about common support here. However, it is worth bearing in mind that violations of positivity can occur that are not obvious in the plots and statistics we have looked at.

## Exercise 6 Assessing covariate balance

Assessing differences in the distribution of covariates between exposed and non-exposed groups before incorporating any type of adjustment for confounding is a key step to identify potential confounding (by measured covariates). Here we will calculate, for each covariate, the standardized difference (StD) between exposed and non-exposed groups. We use StD instead of hypothesis testing because the StD is less influenced by sample size. We will consider values of StD of over 0.1 to indicate meaningful imbalance of that covariate between groups.

To calculate the standardized differences between two groups for continuous variables we use the means and standard deviations of the mean:

$$StD = \frac{(mean_{exposed} - mean_{unexposed})}{sd_{pool}}$$

where  $mean_{exposed}$  and  $mean_{unexposed}$  are the sample mean of the covariate in the exposed and unexposed groups, respectively, and " $sd_{pool}$ " represents the standard deviation of the covariate in the entire sample. There are various ways of obtaining an estimate of the latter. We will calculate this from the sample standard deviations of the covariate in the two groups ( $sd_{exposed}$  and  $sd_{unexposed}$ ) as:

$$sd_{pool} = \sqrt{\frac{(sd_{exposed}^2 + sd_{unexposed}^2)}{2}}$$

# Use this formula to calculate the standardized difference of age between the study drug and the comparator groups.

```
data %>%
  group_by(trt) %>%
  summarize(mean=mean(age), sd=sd(age))
```

```
(46.1 - 40.5) / (sqrt((16.5^2 + 15.4^2) / 2))
```

To calculate the standardized difference between two groups for categorical variables we use proportions and standard deviations of proportions:

$$StD = \frac{(prop_{exposed} - prop_{unexposed})}{sd_{pool}}$$

where  $prop_{exposed}$  and  $prop_{unexposed}$  are the sample proportion with the characteristic of interest in the exposed and unexposed groups, respectively, and  $sd_{pool}$  can be calculated as:

$$sd_{pool} = \sqrt{\frac{((prop_{exposed} \times (1 - prop_{exposed})) + (prop_{unexposed} \times (1 - prop_{unexposed})))}{2}}$$

# Use these formulae to calculate the standardized difference of gender between the study drug and comparator groups.

```
data %>%
  group_by(trt, female) %>%
  summarize(n=n()) %>%
  mutate(perc=prop.table(n))
```

```
(0.427 - 0.372) / (sqrt((0.427 * (1 - 0.427) + 0.372 * (1 - 0.372)) / 2))
```

Alternatively, various R packages can calculate standardized differences for you. For example, you can use the package `TableOne` to do this using the `smd` option:

```
library(tableone)
```

```
vars <- c("age", "female", "ses ", "alc", "bmicat", "nsaid", "cancer", "hyper")
SD_crude <- CreateTableOne(vars, data=data, strata="trt", test=FALSE)
print(SD_crude, smd=TRUE)
```

```
library(cobalt)
```

```
covs <- subset(data, select=c(age, female, alc, bmicat, nsaid, cancer, hyper))
bal.tab(covs, treat=data$trt, binary="std", continuous="std", s.d.denom="pooled")
```

What are the differences between the different ways of obtaining the standardized differences?

How balanced are the covariates?

#### **Exercise 7** Adjust on propensity score

```
model <- glm(outcome ~ trt * ps, data = data, family = binomial)
summary(model)
```

#### **Exercise 8** Generate inverse probability treatment weights

Inverse probability weighting can be used to estimate the average treatment effect (ATE) using the formula below.

$$ATE = E(Y|X_{exposed}) - E(Y|X_{unexposed})$$

# Generate the weights for the treated and untreated group

$$weight_{exposed} = \frac{1}{PS}$$

$$weight_{unexposed} = \frac{1}{(1 - PS)}$$

```
data$wt <- ifelse(data$trt == 1, 1 / data$ps, 1 / (1 - data$ps))
```

# Check balance after weighting

#### **Exercise 10** Perform matching

Matching is one of the most common propensity score methods used in health research. It involves comparing each individual exposed to the study drug to a similar individual exposed to the comparator drug (i.e. they have similar propensity scores).

By selecting two groups with similar distributions of observed confounders, we attempt to replicate the situation of a trial. Matching involves matching exposed and unexposed individuals on their propensity score values. In most applications, we aim to achieve balance in the baseline covariates between the exposed and unexposed groups (i.e. matching on the propensity score should remove differences in observed confounders).

Here, we use the `MatchIt` package to perform 1:1 nearest neighbour matching with replacement and a caliper width of  $0.2 \times \text{SD}(\text{PS})$ .

```
match_data <- matchit(trt ~ outcome + age + female+ ses+ smoke + alc + bmicat + nsaid_rx  
+ cancer + hyper, method="nearest", ratio=1, replace=TRUE, caliper=0.2*sd(data$PS),  
data=data)
```

```
summary(match_data)
```

Summarise matching weights

```
hist(match_data$weights, xlab="Weights", ylab="Number of individuals", main="Distribution  
of matching weights", col="darkgrey", breaks=50)
```

Create matched dataset

```
matched <- match.data(match_data)
```

```
head(matched)
```

Check common support after matching

```
boxplot(PS ~ trt, data= matched, ylab="Propensity score", xlab="Treatment")
```

```
ggplot(matched, aes(x=PS, color=factor(trt), fill=factor(trt))) +  
  geom_histogram(position="identity", bins=50, alpha=0.5) +  
  labs(x = "Propensity score", y="Number of individuals")
```

In a real-life scenario, it is unlikely data would be balanced prior to matching. If imbalances still remain after matching, it may be that the functional form of continuous variables is not correctly specified or there may be some interactions between variables that need to be accounted for in the model.

Now calculate the standardized differences to assess whether matching on the propensity score has created balance on all measured confounders.

```
SD_matched <- CreateTableOne(vars, data=matched, strata="trt", test=FALSE)  
print(SD_matched, smd=TRUE)
```

Finally, let us re-examine the relationship between treatment and the outcome in our matched dataset. Fit a logistic regression model on the matched dataset (using the matching weights):

```
matched_logit <- glm(outcome ~ trt, weights=weights, family=binomial(), data=matched)
```

summary(matched\_logit)

**Exercise 11** Complete the table using the estimates generate from earlier questions. How do the methods compare?

Method	N	Estimate	SE
Adjustment			
Weighting			
Matching			

In conclusion, propensity scores are a powerful method to adjust for observed confounding when estimating causal effects in observational data. However, they rely on strong assumptions and careful implementation to yield valid results. The methods we have covered here are a starting point, and further refinement and sensitivity analysis may be necessary to ensure robust conclusions in your own research.