# Treatment Effect Analysis

## Questions

**Assumptions/Preconditions:**

- Conditional Independence Assumptions (CIA)
- 
$$(BTUgas_i(0), BTUgas_i(1)) \perp Post_i | X_i$$

### 1. Treatment Effect

1. Describe in words the potential outcomes

   - The energy unit BTU change between treatment ($Post_i$) of year discrimination at a random assignment ($\perp$) in the independence of the other confounding variables $X_i$

2. Describe in words and mathematically the unit-level treatment Effect

   - Unit level treatment/ITE is defined by the difference in the treated and untreated (potential outcomes)
   - Mathematically: $ITE_i = Y_i(1) - Y_i(0)$

### 2. CIA Assumptions

1. Explain why $E[BTUgas_i | Post_i = 1] - E[BTUgas_i | Post_i = 0]$ is unlikely to identify ATE

   - During mean difference one could say omitted variable bias exists within the regression proven by:
   $$\frac{Cov(BTUgas_i, Post_i)}{Var(Post_i)} = \frac{Cov(\alpha + \beta Post_i + \gamma X_i + \epsilon)}{Var(D_i)} = \beta + \frac{\gamma Cov(Post_i, X_i)}{Var(Post_i)}$$
   - For true regression will always include $X_i$ along with the $D_i$ treatment

2. Describe in words the CIA Assumptions ($\perp Post_i | X_i$)

   - Only treatments are influenced through random assignment/sampling
   - No confounding variables should influence the outcome

3. Explain why the conditional independence assumption is more plausible than the unconditional independence assumption

   - Some observable factors are controlled to not influence the outcome in CIA
   - IA does not protect against confounding variable bias
   - One possible $X_i$ could be: habitant eating habit (e.g., if eat late: cook more, use dishwasher more, use more lighting)

### 3. COC Assumptions

1. Express the COC assumption mathematically and in words
   - $0 < P(Post_i = 1 | X_i = x) < 1 \quad \forall x$
   - For all possible observations of x, there is covariation of treated and untreated, making comparison possible; no complete independence of treatment to control

2. Why COC may fail at high-dimensional $X_i$
   - COC may fail due to the treatment and control units may not bound to all the pretreatment variables
   - $Pr(D_i | X_i = x) = 1$ or $0$

```r
df$Post_i <- ifelse(df$constYRS == 7, 1, 0)

# ATE Calculated by Difference in Means
treated_mean <- mean(df$BTUgas[df$Post_i == 1], na.rm = TRUE)
control_mean <- mean(df$BTUgas[df$Post_i == 0], na.rm = TRUE)
cat("Pre-code reform energy usage", control_mean, "\n")
```

```
## Pre-code reform energy usage 49.80924
```

```r
ATE_dim <- treated_mean - control_mean
cat('ATE (Difference in Means):', ATE_dim, '\n')
```

```
## ATE (Difference in Means): -2.651527
```

```r
# ATE Calculated by Regression
X <- c("cddzip30_00", "hddzip30_00", "lnsqftK", "numroom", "elecboth",
       "remodeled", "lnyrs_res", "lnrescnt", "lnhhinc", "nr0_5",
       "nr65_99", "college", "anydisabled", "hohblk1", "hohlat1",
       "own", "year09", "cecfast")
X <- na.omit(X)
linear_reg <- paste("BTUgas ~ Post_i +", paste(X, collapse = " + "))
linear_reg <- lm(as.formula(linear_reg), data = df)
ATE_reg <- coef(linear_reg)['Post_i']
cat('ATE (Regression):', ATE_reg, '\n')
```

```
## ATE (Regression): -2.532691
```

```r
# ATE Calculated by Matching on Covariates
X <- df[, c("cddzip30_00", "hddzip30_00", "lnsqftK", "numroom", "elecboth",
            "remodeled", "lnyrs_res", "lnrescnt", "lnhhinc", "nr0_5",
            "nr65_99", "college", "anydisabled", "hohblk1", "hohlat1",
            "own", "year09", "cecfast")]
X_std <- scale(X)
D <- df$Post_i
Y <- df$BTUgas
control_mean <- mean(Y[D == 0], na.rm = TRUE)

K_values <- c(1, 5, 10)
ATE_results <- numeric(length(K_values))
```

```r
for (k_index in seq_along(K_values)) {
    K <- K_values[k_index]
    Y_match <- numeric(nrow(X_std))

    for (i in 1:nrow(X_std)) {
        dist_i <- rowSums((X_std - X_std[i, ])^2)
        dist_i[D == D[i]] <- Inf
        neighbors <- order(dist_i)[1:K]
        Y_match[i] <- mean(Y[neighbors])
    }

    ATE_results[k_index] <- mean(Y_match[D == 1]) - control_mean
}

print(data.frame(K = K_values, ATE = ATE_results))
```

```
##    K       ATE
## 1  1 -7.755631
## 2  5 -7.101051
## 3 10 -5.320479
```

```r
# Reload to fix df being empty
df0 <- read_dta("table3-cols12.dta")
df <- df0[df0$constYRS==6| df0$constYRS ==7, ]
df <- df[!is.na(df$lnBTUgas), ]
df$BTUgas <- exp(df$lnBTUgas)
df$Post_i <- ifelse(df$constYRS == 7, 1, 0)

X <- c("cddzip30_00", "hddzip30_00", "lnsqftK", "numroom", "elecboth",
       "remodeled", "lnyrs_res", "lnrescnt", "lnhhinc", "nr0_5",
       "nr65_99", "college", "anydisabled", "hohblk1", "hohlat1",
       "own", "year09", "cecfast")

# Prepare Propensity Score
prop_model <- glm(as.formula(paste("Post_i", "~", paste(X, collapse = " + "))),
                  data = df, family = binomial)
df$prop_score <- plogis(fitted.values(prop_model))
print(range(df$prop_score))
```

```
## [1] 0.5741077 0.7023633
```

```r
# Perform KNN with distance based on abs(p(x_i) - p(x_j))
K_values <- c(1, 5, 10)
ATE_results <- data.frame(K = K_values, ATE = NA)
ATE_results <- numeric(length(K_values))

for (k_index in seq_along(K_values)) {
    K <- K_values[k_index]
    Y_match <- numeric(nrow(X_std))

    for (i in 1:nrow(X_std)) {
        dist_i <- abs(df$prop_score - df$prop_score[i])
```

```r
        dist_i[D == D[i]] <- Inf
        neighbors <- order(dist_i)[1:K]
        Y_match[i] <- mean(Y[neighbors])
    }

    ATE_results[k_index] <- mean(Y_match[D == 1]) - control_mean
}

print(data.frame(K = K_values, ATE = ATE_results))
```

```
##    K        ATE
## 1  1 -0.1695891
## 2  5 -0.9497518
## 3 10 -0.4213290
```

```r
# ATE by Inverse Probability Weighting
df$weights_treated <- ifelse(df$Post_i == 1, 1 / df$prop_score, 0)
df$weights_control <- ifelse(df$Post_i == 0, 1 / (1 - df$prop_score), 0)

weighted_outcome_treated <- sum(df$Post_i * df$BTUgas * df$weights_treated) /
                            sum(df$weights_treated)
weighted_outcome_control <- sum((1 - df$Post_i) * df$BTUgas * df$weights_control) /
                            sum(df$weights_control)

ATE_IPW <- weighted_outcome_treated - weighted_outcome_control
cat("The estimated ATE using IPW is:", ATE_IPW)
```

```
## The estimated ATE using IPW is: -2.628553
```

**Discussion**

1. Outside of KNN and propensity scored KNN, the ATE estimates are fairly consistent at 5 percent in negative direction (e.g., -2.5 to -2.6 out of 49.08). We see a sharp decrease from KNN result at -7/49, with some decrease as we increased the K for the classification. Using the propensity score for the distance calculated KNN, we have a marginal difference at fraction of decreased units.

2. The KNN implementation is harder than expected. I believe the simplest may be Regression, IPW and the Mean Difference. Writing algorithms gets self-doubt. One might need to try different methods of ATE estimation to see their consistent outcome. I am not sure why the propensity-KNN is returning a marginal difference.