# DATA37711_AS3_Adrian Cao

November 23, 2024

# 1 Question 1: Causal DAGs

## 1.1 (a) Negative association between COVID-19 severity and smoking

- **Confounder**: A confounder (e.g., age, health status, or other underlying conditions) influences both smoking and the severity of COVID-19. For example, smoking might correlate with age, and younger individuals are more likely to smoke also be at lower risk for severe COVID-19 infections. Thus, **age acts as a confounder**, leading to the observed negative association.

  Causal DAG: `Smoking <-- Age --> COVID-19 Severity`

  Without controlling for age, the observed association between smoking and COVID-19 severity may not reflect a direct causal relationship.

- **Collider Bias**: If severity of COVID-19 and smoking both influence hospitalization or testing rates (e.g., only severe cases or smokers are hospitalized/tested more often), selection bias (conditioning on hospitalization) could create a spurious association.

  Causal DAG with collider: `Smoking --> Hospitalization <-- COVID-19 Severity`

  Conditioning on hospitalization (a collider) might create an artificial negative association.

## 1.2 (b) Simpson's Paradox in Berkeley Admissions

This is a **classic example of Simpson's Paradox**. The overall acceptance rates can differ between men and women due to **confounding by department**. For example, Male applicants may disproportionately apply to departments with higher acceptance rates (e.g., STEM fields). Female applicants may disproportionately apply to departments with lower acceptance rates (e.g., humanities).

Thus, while female applicants may have higher acceptance rates **within each department**, the overall acceptance rate for women is lower due to the distribution of applications across departments.

Causal DAG: `Applicant Gender --> Department Applied To --> Acceptance Rate`

In this DAG: - Gender influences which department is applied to. - The department determines the acceptance rate. - If not properly stratified by department, the observed gender disparity in acceptance rates can appear misleading.

Simpson's Paradox illustrates that aggregating across groups without accounting for such confounders can lead to paradoxical conclusions.

# 2 Question 2: Average Treatment Effect

## 2.1 (a): Why $\tau_{ATT}$ is Preferred in Certain Cases

The Average Treatment Effect on the Treated ($\tau_{ATT}$) is preferred when some units are very unlikely to be treated because:

1. **Rare Treatment Groups**: For certain subpopulations (e.g., older unmarried men), there might be insufficient data to accurately estimate the effect of the treatment due to sparse overlap between treated and untreated units.
2. **Focus on Treated Population**: $\tau_{ATT}$ focuses exclusively on the population that actually receives the treatment, providing more reliable insights into the treatment's effect on this specific group.
3. **Captures Qualitative Goal**: While $\tau_{ATT}$ does not represent the overall causal effect for the entire population, it still provides meaningful information about the treatment's impact on the treated group.

## 2.2 (b): Plug-In Estimator for $\tau_{ATT}$

To estimate $\tau_{ATT}$ using a plug-in estimator: 1. **Estimate Conditional Expectations**: - Use a model (e.g., Random Forest) to estimate $\mathbb{E}[Y|A=1,X]$ and $\mathbb{E}[Y|A=0,X]$ for the treated group $A=1$. 2. **Compute Differences**: - For each treated unit ($A=1$), compute the difference $\hat{\tau}(X) = \mathbb{E}[Y|A=1,X] - \mathbb{E}[Y|A=0,X]$. 3. **Average Over Treated Units**: - Average the computed differences for all treated units to estimate $\tau_{ATT}$.

Formula:

$$\hat{\tau}_{ATT} = \frac{1}{n_1} \sum_{i:A_i=1} (\mathbb{E}[Y|A=1,X_i] - \mathbb{E}[Y|A=0,X_i])$$

where $n_1$ is the number of treated units ($A=1$).

## 2.3 (c) and (d) Code Implementation

```
[10]: import numpy as np
      import pandas as pd
      from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
      from sklearn.model_selection import StratifiedKFold, KFold
      from sklearn.utils import resample

      def make_data_lalonde(df):
          df_new = df.drop(['nodegree'], axis=1)
          df_new['pos74'] = (df_new['RE74'] > 0).astype(int)
          df_new['pos75'] = (df_new['RE75'] > 0).astype(int)
          df_new['treatment'] = df_new['treatment'].astype(int)
          return df_new


      col_names = ['treatment', 'age', 'education', 'black',
                   'hispanic', 'married', 'nodegree', 'RE74', 'RE75', 'RE78']
```

```python
control = pd.read_csv('https://raw.githubusercontent.com/anishazaveri/
 ↪austen_plots/master/data/imbens-raw/psid_controls.txt', header=None,
 ↪sep=r"\s\s", names=col_names, engine='python')
treatment = pd.read_csv('https://raw.githubusercontent.com/anishazaveri/
 ↪austen_plots/master/data/imbens-raw/nswre74_treated.txt', header=None,
 ↪sep=r"\s\s", names=col_names, engine='python')

lalonde1 = pd.concat([control, treatment]).reset_index(drop=True)
lalonde1 = make_data_lalonde(lalonde1)

confounders = lalonde1.drop(columns=['RE78', 'treatment'])
outcome = lalonde1['RE78']
treatment = lalonde1['treatment']
```

```python
[12]: # Define model functions
def make_Q_model():
    return RandomForestRegressor(random_state=42, n_estimators=500)

def make_g_model():
    return RandomForestClassifier(random_state=42, n_estimators=100,
 ↪max_depth=5)

# Cross-fitting functions
def treatment_k_fold_fit_and_predict(make_model, X, A, n_splits):
    predictions = np.full_like(A, np.nan, dtype=float)
    kf = StratifiedKFold(n_splits=n_splits, shuffle=True, random_state=42)
    for train_index, test_index in kf.split(X, A):
        X_train, A_train = X.iloc[train_index], A.iloc[train_index]
        model = make_model()
        model.fit(X_train, A_train)
        predictions[test_index] = model.predict_proba(X.iloc[test_index])[:, 1]
    return predictions

def outcome_k_fold_fit_and_predict(make_model, X, y, A, n_splits, output_type):
    predictions0 = np.full_like(A, np.nan, dtype=float)
    predictions1 = np.full_like(A, np.nan, dtype=float)
    kf = KFold(n_splits=n_splits, shuffle=True, random_state=42)
    X_w_treatment = X.copy()
    X_w_treatment['A'] = A
    X0, X1 = X_w_treatment.copy(), X_w_treatment.copy()
    X0['A'], X1['A'] = 0, 1
    for train_index, test_index in kf.split(X, y):
        model = make_model()
        model.fit(X_w_treatment.iloc[train_index], y.iloc[train_index])
        predictions0[test_index] = model.predict(X0.iloc[test_index])
        predictions1[test_index] = model.predict(X1.iloc[test_index])
    return predictions0, predictions1
```

```python
# Predict nuisance parameters
g = treatment_k_fold_fit_and_predict(make_g_model, confounders, treatment,
 ↪n_splits=10)
Q0, Q1 = outcome_k_fold_fit_and_predict(make_Q_model, confounders, outcome,
 ↪treatment, n_splits=10, output_type="continuous")

# Data preparation for ATT
data_and_nuisance_estimates = pd.DataFrame({'g': g, 'Q0': Q0, 'Q1': Q1, 'A':
 ↪treatment, 'Y': outcome})

# Function for plug-in ATT estimation
def plugin_att_estimator(Q1, Q0, A):
    # Subset for treated group
    treated_idx = A == 1
    # Compute the ATT as the average difference for treated units
    tau_att = np.mean(Q1[treated_idx] - Q0[treated_idx])
    return tau_att

# Calculate the plug-in ATT estimate
tau_ATT_plugin = plugin_att_estimator(data_and_nuisance_estimates['Q1'],
                                      data_and_nuisance_estimates['Q0'],
                                      data_and_nuisance_estimates['A'])

# Bootstrap confidence interval for plug-in estimator
n_bootstrap = 1000
bootstrap_estimates_plugin = []

# Perform bootstrap
for _ in range(n_bootstrap):
    # Resample data with replacement
    bootstrap_sample = resample(data_and_nuisance_estimates, replace=True,
 ↪random_state=42 + _)

    # Extract Q1, Q0, A from the bootstrap sample
    Q1_bootstrap = bootstrap_sample['Q1']
    Q0_bootstrap = bootstrap_sample['Q0']
    A_bootstrap = bootstrap_sample['A']

    # Estimate ATT using the plug-in estimator
    tau_att_bootstrap = plugin_att_estimator(Q1_bootstrap, Q0_bootstrap,
 ↪A_bootstrap)
    bootstrap_estimates_plugin.append(tau_att_bootstrap)

# Compute 95% confidence interval
ci_lower_plugin = np.percentile(bootstrap_estimates_plugin, 2.5)
ci_upper_plugin = np.percentile(bootstrap_estimates_plugin, 97.5)
```

```
# Display results
print(f"Plug-In ATT Estimate: {tau_ATT_plugin}")
print(f"95% Confidence Interval (Plug-In): ({ci_lower_plugin},␣
 ↪{ci_upper_plugin})")
```

```
Plug-In ATT Estimate: 1085.3404537748204
95% Confidence Interval (Plug-In): (654.0181089834856, 1532.9075526666634)
```

## 2.4 Conclusion

The plug-in estimator provides a simpler and computationally efficient method for estimating $\tau_{ATT}$, yielding a narrower confidence interval and a slightly lower point estimate (1085.34) compared to the AIPTW estimator (1300.98). The plug-in approach avoids the sensitivity to extreme propensity scores seen in AIPTW, making it preferable in cases with severe overlap issues. However, the AIPTW estimator is doubly robust and more reliable in datasets with well-balanced treatment assignment, as it adjusts for propensity score imbalance. Overall, the plug-in estimator is ideal for scenarios prioritizing simplicity and lower variance, while AIPTW offers robustness at the cost of higher variance and potential instability in small or imbalanced samples.

# 3 Question 3: Spurious features and Fairness

## 3.1 (a) Statistics and Fairness notion

### 3.1.1 1. Demographic Parity (DP):

$$\text{DP} := \text{average}(\sigma(\hat{f}(x_i)) \,|\, z_i = 1) - \text{average}(\sigma(\hat{f}(x_i)) \,|\, z_i = 0)$$

- **Criterion**: Demographic parity ensures that the predictions $\hat{f}(X)$ are independent of the sensitive attribute $Z$. This means that $\hat{f}(X) \perp\!\!\!\perp Z$. - **Equal to 0 When**: DP equals 0 if the model $\hat{f}(X)$ satisfies demographic parity, i.e., the average predicted probability of the outcome is the same across groups $Z = 1$ and $Z = 0$.

### 3.1.2 2. Equalized Odds (EO):

$$\text{EO} := \text{average}(\sigma(\hat{f}(x_i)) \,|\, z_i = 1, y_i = 1) - \text{average}(\sigma(\hat{f}(x_i)) \,|\, z_i = 0, y_i = 1)$$

- **Criterion**: Equalized odds ensures that the predictions $\hat{f}(X)$ are independent of $Z$, conditional on the true outcome $Y$. This means $\hat{f}(X) \perp\!\!\!\perp Z \,|\, Y$. - **Equal to 0 When**: EO equals 0 if the model $\hat{f}(X)$ satisfies equalized odds, i.e., the true positive rate and false positive rate are the same across groups $Z = 1$ and $Z = 0$.

### 3.1.3 3. Predictive Parity (PP):

$$\text{PP} := \text{average}(|y_i - \sigma(\hat{f}(x_i))| \,|\, z_i = 1) - \text{average}(|y_i - \sigma(\hat{f}(x_i))| \,|\, z_i = 0)$$

- **Criterion**: Predictive parity ensures that the sensitive attribute $Z$ is independent of the true outcome $Y$, conditional on the predicted probability $\hat{f}(X)$. This means $Z \perp\!\!\!\perp Y \,|\, \hat{f}(X)$. - **Equal to 0 When**: PP equals 0 if the model $\hat{f}(X)$ satisfies predictive parity, i.e., the residuals (absolute differences between true $Y$ and predicted $\hat{f}(X)$) are the same across groups $Z = 1$ and $Z = 0$.

## 3.2 (b) Fairness through unawareness

```python
# Import necessary libraries
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.metrics import accuracy_score
import numpy as np

# The dataset contains:
# - Features for predicting income (e.g., age, education, work hours, etc.)
# - "income" column: target variable (Y), binary (0/1)
# - "Z" column: sensitive attribute
df = pd.read_csv('fairness_data.csv')

# Separate features (X), target (Y), and sensitive attribute (Z)
X = df.drop(columns=["income", "Z"])  # Features: exclude "income" (Y) and "Z"
 ↪(sensitive attribute)
Y = df["income"]  # Target variable: income (binary)
Z = df["Z"]  # Sensitive attribute: Z (binary)

X_train, X_test, Y_train, Y_test, Z_train, Z_test = train_test_split(X, Y, Z,
 ↪test_size=0.2, random_state=42)

# Use a Gradient Boosting Classifier
clf_ftu = GradientBoostingClassifier(random_state=42)
clf_ftu.fit(X_train, Y_train)

# Predict both class labels and probabilities
Y_pred = clf_ftu.predict(X_test)
Y_pred_proba = clf_ftu.predict_proba(X_test)[:, 1]  # Predicted probabilities
 ↪for the positive class

# Calculate test-set accuracy
test_accuracy = accuracy_score(Y_test, Y_pred)

# Demographic Parity (DP): Difference in average predicted probabilities across
 ↪sensitive groups
dp = np.mean(Y_pred_proba[Z_test == 1]) - np.mean(Y_pred_proba[Z_test == 0])

# Equalized Odds (EO): Difference in average predicted probabilities
 ↪conditioned on Y=1 across sensitive groups
eo = np.mean(Y_pred_proba[(Z_test == 1) & (Y_test == 1)]) - np.
 ↪mean(Y_pred_proba[(Z_test == 0) & (Y_test == 1)])

# Predictive Parity (PP): Difference in average residuals (true Y - predicted
 ↪probabilities) across sensitive groups
```

```
pp = np.mean(np.abs(Y_test[Z_test == 1] - Y_pred_proba[Z_test == 1])) - \
     np.mean(np.abs(Y_test[Z_test == 0] - Y_pred_proba[Z_test == 0]))

print("Results:")
print(f"Test Set Accuracy: {test_accuracy:.2f}")
print(f"Demographic Parity (DP): {dp:.4f}")
print(f"Equalized Odds (EO): {eo:.4f}")
print(f"Predictive Parity (PP): {pp:.4f}")
```

```
Results:
Test Set Accuracy: 0.84
Demographic Parity (DP): -0.1697
Equalized Odds (EO): -0.1911
Predictive Parity (PP): -0.2439
```

### 3.3 (c) Marginalizing out the sensitive attribute

```
[ ]: clf_with_z = GradientBoostingClassifier(random_state=42)
     X_with_z_train = X_train.copy()
     X_with_z_train["Z"] = Z_train   # Include the sensitive attribute
     clf_with_z.fit(X_with_z_train, Y_train)

     # Create a new feature set without the sensitive attribute for prediction
     X_without_z_test = X_test.copy()

     # Predict for Z=0 and Z=1, then average the predictions
     X_with_z_test_0 = X_without_z_test.copy()
     X_with_z_test_0["Z"] = 0   # Set Z=0 for all test data
     Y_pred_proba_z0 = clf_with_z.predict_proba(X_with_z_test_0)[:, 1]   # Predicted
      ↪probabilities for Z=0

     X_with_z_test_1 = X_without_z_test.copy()
     X_with_z_test_1["Z"] = 1   # Set Z=1 for all test data
     Y_pred_proba_z1 = clf_with_z.predict_proba(X_with_z_test_1)[:, 1]   # Predicted
      ↪probabilities for Z=1

     # Average the predictions for Z=0 and Z=1 to get marginalized predictions
     Y_pred_proba_marginalized = (Y_pred_proba_z0 + Y_pred_proba_z1) / 2
     Y_pred_marginalized = (Y_pred_proba_marginalized > 0.5).astype(int)   # Convert
      ↪probabilities to binary predictions

     # Calculate test-set accuracy
     test_accuracy_marginalized = accuracy_score(Y_test, Y_pred_marginalized)

     # Calculate fairness statistics
     # Demographic Parity (DP)
```

```python
dp_marginalized = np.mean(Y_pred_proba_marginalized[Z_test == 1]) - np.
 ↪mean(Y_pred_proba_marginalized[Z_test == 0])

# Equalized Odds (EO)
eo_marginalized = np.mean(Y_pred_proba_marginalized[(Z_test == 1) & (Y_test ==␣
 ↪1)]) - \
                  np.mean(Y_pred_proba_marginalized[(Z_test == 0) & (Y_test ==␣
 ↪1)])

# Predictive Parity (PP)
pp_marginalized = np.mean(np.abs(Y_test[Z_test == 1] -␣
 ↪Y_pred_proba_marginalized[Z_test == 1])) - \
                  np.mean(np.abs(Y_test[Z_test == 0] -␣
 ↪Y_pred_proba_marginalized[Z_test == 0]))

# Step 4: Print results
print("Results for Marginalized Classifier:")
print(f"Test Set Accuracy: {test_accuracy_marginalized:.2f}")
print(f"Demographic Parity (DP): {dp_marginalized:.4f}")
print(f"Equalized Odds (EO): {eo_marginalized:.4f}")
print(f"Predictive Parity (PP): {pp_marginalized:.4f}")
```

```
Results for Marginalized Classifier:
Test Set Accuracy: 0.83
Demographic Parity (DP): -0.0038
Equalized Odds (EO): 0.0565
Predictive Parity (PP): -0.1551
```

### 3.4 (d) Proof for marginalization

1. **Causal Graph**:
   - $Z$ is a **common cause** for $W$ and $Y$.
   - $P$ influences $\tilde{Y}$, which is the observed, noisy version of $Y$.
   - $X = (P, W)$, meaning the features $X$ include $P$ (a direct cause of $\tilde{Y}$) and $W$ (caused by $Z$).
2. **Goal**:
   - Prove $f_{\text{inv}}(x)$, the marginalized classifier, depends only on $P$ and is independent of $Z$ and $W$.
3. **Classifier Definitions**:
   - $f(x, z) = P(Y = 1 \mid X = x, Z = z) = P(Y = 1 \mid P = p, W = w, Z = z)$, where $X = (P, W)$.
   - Marginalized classifier:

$$f_{\text{inv}}(x) = \frac{1}{|Z|} \sum_{z \in Z} P(Y = 1 \mid P = p, W = w, Z = z).$$

### 3.4.1 Step 1: Independence Structure From the Graph

From the causal graph: - $Z$ is a **common cause** for $W$ and $Y$, so:

$$P(Y \mid P, W, Z) = P(Y \mid P, Z).$$

This is because $P$ directly influences $Y$ (via $\tilde{Y}$), and $Z$ is a common cause.

Thus, $f(x, z)$ simplifies to:

$$f(x, z) = P(Y = 1 \mid P = p, Z = z).$$

### 3.4.2 Step 2: Marginalize Over $Z$

The marginalized classifier $f_{\text{inv}}(x)$ is defined as:

$$f_{\text{inv}}(x) = \frac{1}{|Z|} \sum_{z \in Z} P(Y = 1 \mid P = p, Z = z) = \mathbb{E}_Z[P(Y = 1 \mid P = p, Z)].$$

Since $P(Y = 1 \mid P = p)$ is independent of $Z$, the summation simplifies to:

$$f_{\text{inv}}(x) = P(Y = 1 \mid P = p).$$

### 3.4.3 Step 3: Dependence on $P$ Only

Because $f_{\text{inv}}(x)$ simplifies to $P(Y = 1 \mid P = p)$: 1. $Z$ is marginalized out. 2. $W$ has no remaining influence because it depends on $Z$, and $Z$'s influence has already been accounted for.

Thus, $f_{\text{inv}}(x)$ depends only on $P$.

## 3.5 (e) Real data analysis

```
[15]: # Step 1: Prepare the dataset with "sex" as the sensitive attribute
      # Assuming "sex" is a column in the dataset
      # Drop the synthetic sensitive attribute "Z" and use "sex" as the sensitive
       ↪attribute
      sensitive_attribute = "sex"
      X_real_sensitive = df.drop(columns=["income", sensitive_attribute])  # Features
       ↪without income (Y) and sex
      Y_real_sensitive = df["income"]  # Target variable
      Z_real_sensitive = df[sensitive_attribute]  # Sensitive attribute (sex)

      # Split into training and test sets
      X_train_real, X_test_real, Y_train_real, Y_test_real, Z_train_real, Z_test_real
       ↪= train_test_split(
          X_real_sensitive, Y_real_sensitive, Z_real_sensitive, test_size=0.2,
       ↪random_state=42
      )

      # Step 2: Train a Fairness Through Unawareness (FTU) classifier
      clf_ftu_real = GradientBoostingClassifier(random_state=42)
```

```python
clf_ftu_real.fit(X_train_real, Y_train_real)

# Predict on the test set
Y_pred_real = clf_ftu_real.predict(X_test_real)
Y_pred_proba_real = clf_ftu_real.predict_proba(X_test_real)[:, 1]

# Calculate test-set accuracy
test_accuracy_real_ftu = accuracy_score(Y_test_real, Y_pred_real)

# Calculate fairness statistics
# Demographic Parity (DP)
dp_real_ftu = np.mean(Y_pred_proba_real[Z_test_real == 1]) - np.
 ↪mean(Y_pred_proba_real[Z_test_real == 0])

# Equalized Odds (EO)
eo_real_ftu = np.mean(Y_pred_proba_real[(Z_test_real == 1) & (Y_test_real ==␣
 ↪1)]) - \
              np.mean(Y_pred_proba_real[(Z_test_real == 0) & (Y_test_real ==␣
 ↪1)])

# Predictive Parity (PP)
pp_real_ftu = np.mean(np.abs(Y_test_real[Z_test_real == 1] -␣
 ↪Y_pred_proba_real[Z_test_real == 1])) - \
              np.mean(np.abs(Y_test_real[Z_test_real == 0] -␣
 ↪Y_pred_proba_real[Z_test_real == 0]))

# Step 3: Train a classifier including the sensitive attribute "sex"
clf_with_sex = GradientBoostingClassifier(random_state=42)
X_with_sex_train = X_train_real.copy()
X_with_sex_train[sensitive_attribute] = Z_train_real  # Include the sensitive␣
 ↪attribute
clf_with_sex.fit(X_with_sex_train, Y_train_real)

# Marginalize the sensitive attribute for predictions
# Predict for sex=0 and sex=1, then average the predictions
X_with_sex_test_0 = X_test_real.copy()
X_with_sex_test_0[sensitive_attribute] = 0  # Set sex=0
Y_pred_proba_sex_0 = clf_with_sex.predict_proba(X_with_sex_test_0)[:, 1]

X_with_sex_test_1 = X_test_real.copy()
X_with_sex_test_1[sensitive_attribute] = 1  # Set sex=1
Y_pred_proba_sex_1 = clf_with_sex.predict_proba(X_with_sex_test_1)[:, 1]

# Average the predictions for sex=0 and sex=1
Y_pred_proba_marginalized_sex = (Y_pred_proba_sex_0 + Y_pred_proba_sex_1) / 2
Y_pred_marginalized_sex = (Y_pred_proba_marginalized_sex > 0.5).astype(int)
```

```python
# Step 4: Evaluate the marginalized classifier
# Test-set accuracy
test_accuracy_marginalized_sex = accuracy_score(Y_test_real,␣
 ↪Y_pred_marginalized_sex)

# Fairness statistics
# Demographic Parity (DP)
dp_marginalized_sex = np.mean(Y_pred_proba_marginalized_sex[Z_test_real == 1])␣
 ↪- \
                    np.mean(Y_pred_proba_marginalized_sex[Z_test_real == 0])

# Equalized Odds (EO)
eo_marginalized_sex = np.mean(Y_pred_proba_marginalized_sex[(Z_test_real == 1)␣
 ↪& (Y_test_real == 1)]) - \
                    np.mean(Y_pred_proba_marginalized_sex[(Z_test_real == 0)␣
 ↪& (Y_test_real == 1)])

# Predictive Parity (PP)
pp_marginalized_sex = np.mean(np.abs(Y_test_real[Z_test_real == 1] -␣
 ↪Y_pred_proba_marginalized_sex[Z_test_real == 1])) - \
                    np.mean(np.abs(Y_test_real[Z_test_real == 0] -␣
 ↪Y_pred_proba_marginalized_sex[Z_test_real == 0]))

# Step 5: Print results
print("Results for Fairness Through Unawareness (FTU) Classifier:")
print(f"Test Set Accuracy: {test_accuracy_real_ftu:.2f}")
print(f"Demographic Parity (DP): {dp_real_ftu:.4f}")
print(f"Equalized Odds (EO): {eo_real_ftu:.4f}")
print(f"Predictive Parity (PP): {pp_real_ftu:.4f}")

print("\nResults for Marginalized Classifier Including Sensitive Attribute␣
 ↪(Sex):")
print(f"Test Set Accuracy: {test_accuracy_marginalized_sex:.2f}")
print(f"Demographic Parity (DP): {dp_marginalized_sex:.4f}")
print(f"Equalized Odds (EO): {eo_marginalized_sex:.4f}")
print(f"Predictive Parity (PP): {pp_marginalized_sex:.4f}")
```

```
Results for Fairness Through Unawareness (FTU) Classifier:
Test Set Accuracy: 0.85
Demographic Parity (DP): -0.0612
Equalized Odds (EO): -0.0952
Predictive Parity (PP): -0.0562

Results for Marginalized Classifier Including Sensitive Attribute (Sex):
Test Set Accuracy: 0.85
Demographic Parity (DP): -0.0589
```

```
Equalized Odds (EO): -0.0909
Predictive Parity (PP): -0.0556
```

### 3.5.1 Conclusion:

The results using the real sensitive attribute ("sex") show higher fairness across all metrics compared to the synthetic sensitive attribute ("Z"). For the Fairness Through Unawareness (FTU) classifier, the test accuracy slightly improved (85% vs. 84%), while fairness metrics such as Demographic Parity (DP), Equalized Odds (EO), and Predictive Parity (PP) have significantly smaller disparities (-0.0612, -0.0952, and -0.0562, respectively, compared to -0.1697, -0.1911, and -0.2439 for the synthetic attribute). Similarly, for the marginalized classifier, fairness improved dramatically in DP (-0.0589 vs. -0.0038), EO (-0.0909 vs. 0.0565), and PP (-0.0556 vs. -0.1551), while test accuracy increased from 83% to 85%. These results suggest that the real sensitive attribute is less correlated with unfair model behavior compared to the synthetic attribute, likely due to its causal structure in the data.