

Directory Location

```
In [22]: 1 #set root as a relative path as users may use different OS
2 ROOT = "./"
3 DAT = ROOT # QBCollegeStats.xsv is stored in the root
4 FIG = ROOT + "outputs/"
```

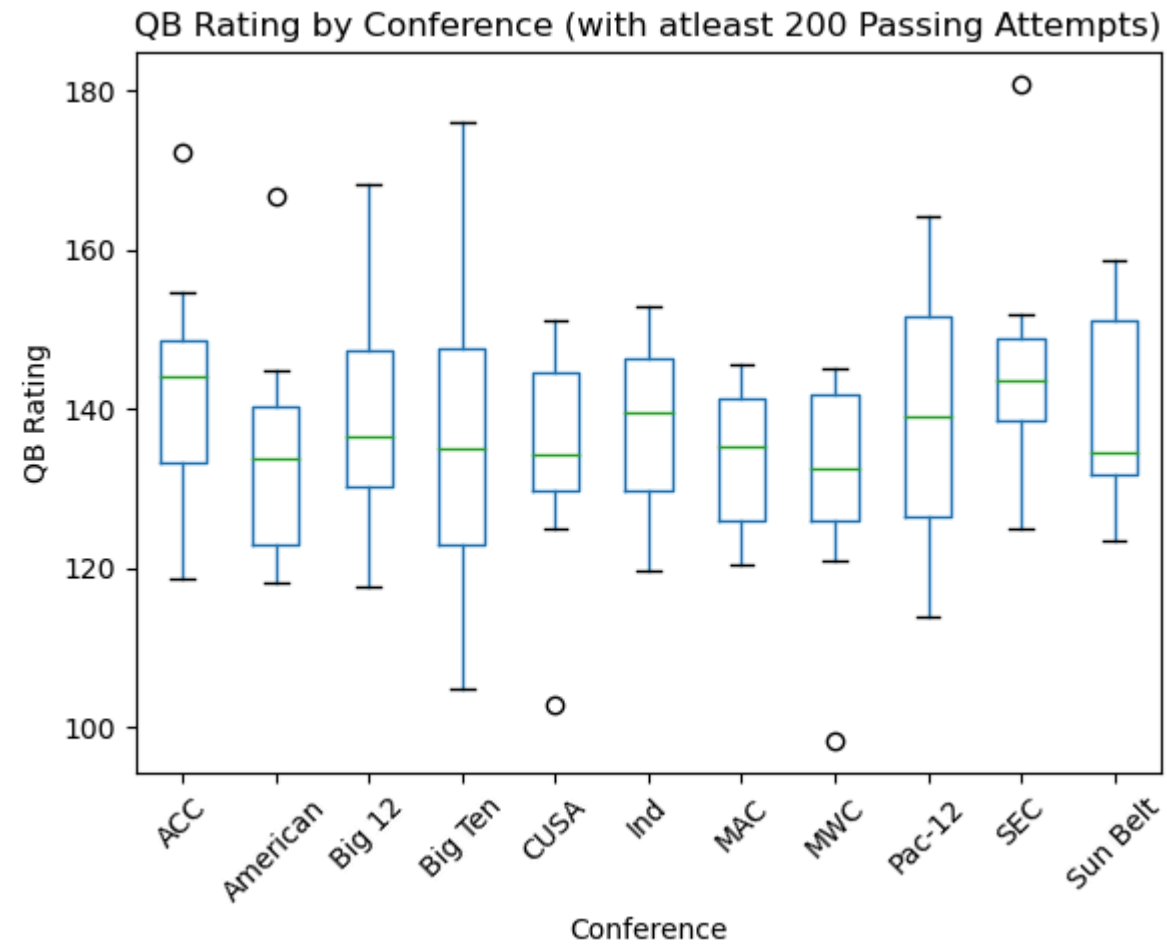
Configure the data base, taken from <https://www.sports-reference.com/cfb/years/2024-passing.html> (<https://www.sports-reference.com/cfb/years/2024-passing.html>).

```
In [23]: 1 import pandas as pd
2 import numpy as np
3
4 # load the data
5 qb_data = pd.read_csv("QBCollegeStats.csv")
6
7 # data is stored in one column, separate it into separate columns
8 if qb_data.shape[1] == 1:
9     qb_data = qb_data[qb_data.columns[0]].str.split(",", expand=True)
10
11 # Set column headers
12 qb_data.columns = ["Rk", "Player", "Team", "Conf", "G", "Cmp", "Att", "Cmp%", "Yds",
13 "TD", "TD%", "Int", "Int%", "Y/A", "AY/A", "Y/C", "Y/G", "Rate",
14 "Awards", "Player-additional"]
15
16 #rename conferences of these two teams as they are stored wrong in database
17 qb_data.loc[qb_data["Team"].isin(["Oregon State", "Washington State"]), "Conf"] = "Pac-12"
```

Box plot of QB rating by Conference

```
In [24]: 1 import matplotlib.pyplot as plt
2
3 #ensure QB rating and QB passing attempts is numeric values
4 qb_data["Rate"] = pd.to_numeric(qb_data["Rate"])
5 qb_data["Att"] = pd.to_numeric(qb_data["Att"])
6
7 # Filter for QBs with atleast 200 passing attempts
8 qb_data_filtered = qb_data[(qb_data["Att"] >= 200)]
9
10 # Plot figure
11 plt.figure(figsize=(10, 6))
12
13 # plot boxplot with data grouped conference
14 qb_data_filtered.boxplot(column="Rate", by="Conf", grid = False)
15
16 #plot axis
17 plt.title("QB Rating by Conference (with atleast 200 Passing Attempts)")
18 plt.suptitle("")
19 plt.xlabel("Conference")
20 plt.ylabel("QB Rating")
21
22 #rotate labels of x axis to make it more legible
23 plt.xticks(rotation=45)
24
25 #store it as pdf in a folder
26 FIG = "./outputs/"
27 plt.savefig(FIG+"QBRatingBoxPlot.pdf")
28
29 plt.show()
```

<Figure size 1000x600 with 0 Axes>



Import advance modeling tools

```

In [25]: 1 from econml.dml import CausalForestDML
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import GradientBoostingRegressor, GradientBoostingClassifier
4 from statsmodels.regression.linear_model import OLS
5 from statsmodels.tools import add_constant
6
7 #name the power 5 conferences
8 power5_conferences = ["SEC", "Big Ten", "Big 12", "ACC", "Pac-12"]
9
10 # make a copy to avoid SettingWithCopyWarning
11 qb_data_filtered = qb_data_filtered.copy()
12
13 # Create a binary Power5 column, uses loc to avoid SettingWithCopyWarning
14 qb_data_filtered.loc[:, "Power5"] = qb_data_filtered["Conf"].apply(lambda x: 1 if x in power5_conferences
15
16 # Select numeric columns for the model, input variable
17 X = qb_data_filtered[[
18     "Rk", "G", "Cmp", "Att", "Cmp%", "Yds", "TD", "TD%", "Int", "Int%",
19     "Y/A", "AY/A", "Y/C", "Y/G"
20 ]]
21
22 #dependent variable
23 D = qb_data_filtered["Power5"]
24
25 #output variable
26 Y = qb_data_filtered["Rate"]
27
28 #combine input variable with dependent variable
29 X_with_D = pd.concat([X, D], axis=1)
30

```

Linear Regression Model

In [26]:

```
1 # Combine all variables
2 combined = pd.concat([Y, X_with_D], axis=1)
3
4 # Drop rows that contain values such as "NaNs"
5 combined_clean = combined.dropna()
6 Y_clean = combined_clean["Rate"]
7 X_clean = combined_clean.drop(columns=["Rate"])
8
9 # Data stored were still objects, convert it all to numerical
10 X_clean = X_clean.apply(pd.to_numeric)
11 Y_clean = pd.to_numeric(Y_clean)
12
13 # merge again and drop rows that converted into "NaNs" value after conversioj
14 combined_final = pd.concat([Y_clean, X_clean], axis=1).dropna()
15
16 # Split it back to X and Y
17 Y_clean = combined_final.iloc[:, 0]
18 X_clean = combined_final.iloc[:, 1:]
19
20 #add a constant term for X
21 X_clean = add_constant(X_clean)
22
23 #create a OLS regression model
24 OLS_regression_model = OLS(Y_clean, X_clean).fit()
25
26 # Create a figure (to store it as a pdf)
27 fig, ax = plt.subplots(figsize=(12, 8))
28 #remove axis
29 ax.axis("off")
30
31 # Convert the summary to a string
32 OLS_summary = OLS_regression_model.summary().as_text()
33
34 # Plot the text which was converted as a string onto the figure
35 ax.text(0, 1, OLS_summary, fontsize=10, va="top", family="monospace")
36
37 # Store as PDF
38 plt.savefig(FIG + "OLS_Regression_Summary.pdf", bbox_inches="tight")
39
40 #output
```

```
41 plt.show()
```

OLS Regression Results

```

=====
Dep. Variable:          Rate    R-squared:          1.000
Model:                  OLS      Adj. R-squared:        1.000
Method:                 Least Squares    F-statistic:        2.823e+05
Date:                   Thu, 24 Apr 2025    Prob (F-statistic):    2.03e-234
Time:                   15:03:07      Log-Likelihood:        145.38
No. Observations:       121      AIC:                  -258.8
Df Residuals:           105      BIC:                  -214.0
Df Model:                15
Covariance Type:        nonrobust
=====

```

	coef	std err	t	P> t	[0.025	0.975]
const	-1.6657	1.094	-1.522	0.131	-3.836	0.504
Rk	0.0017	0.001	1.617	0.109	-0.000	0.004
G	-0.0056	0.030	-0.187	0.852	-0.065	0.054
Cmp	-0.0035	0.003	-1.373	0.173	-0.009	0.002
Att	-4.73e-05	0.002	-0.026	0.979	-0.004	0.004
Cmp%	1.0343	0.019	55.440	0.000	0.997	1.071
Yds	0.0005	0.000	2.498	0.014	0.000	0.001
TD	-0.0013	0.008	-0.159	0.874	-0.017	0.015
TD%	1.7179	0.069	24.955	0.000	1.581	1.854
Int	-0.0146	0.015	-1.003	0.318	-0.044	0.014
Int%	1.6200	0.132	12.310	0.000	1.359	1.881
Y/A	0.1327	0.276	0.481	0.632	-0.415	0.680
AY/A	7.9417	0.301	26.344	0.000	7.344	8.539
Y/C	0.1298	0.068	1.909	0.059	-0.005	0.265
Y/G	-0.0006	0.002	-0.366	0.715	-0.004	0.003
Power5	-0.0125	0.016	-0.779	0.438	-0.044	0.019

```

=====
Omnibus:                 0.816    Durbin-Watson:          1.999
Prob(Omnibus):           0.665    Jarque-Bera (JB):        0.848
Skew:                    -0.001    Prob(JB):                0.654
Kurtosis:                 2.590    Cond. No.                 4.20e+05
=====

```

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 4.2e+05. This might indicate that there are strong multicollinearity or other numerical problems.

Estimate casual forest

```
In [27]: 1 cf = CausalForestDML(
2         model_y=GradientBoostingRegressor(),
3         model_t=GradientBoostingClassifier(), #use a classifier for treatment
4         discrete_treatment=True,
5         random_state=42
6     )
7
8     #fitting the treatment effects into the model
9     cf.fit(Y, D, X=X)
10    tau_hat = cf.effect(X)
11    tau_hat_se = cf.effect_interval(X)
12
13    #check and print summary
14    cf.summary()
```

Population summary results are available only if `cache_values=True` at fit time!

Out [27]:

Doubly Robust ATE on Training Data Results

	point_estimate	stderr	zstat	pvalue	ci_lower	ci_upper
ATE	-8.53	6.951	-1.227	0.22	-22.154	5.094

Doubly Robust ATT(T=0) on Training Data Results

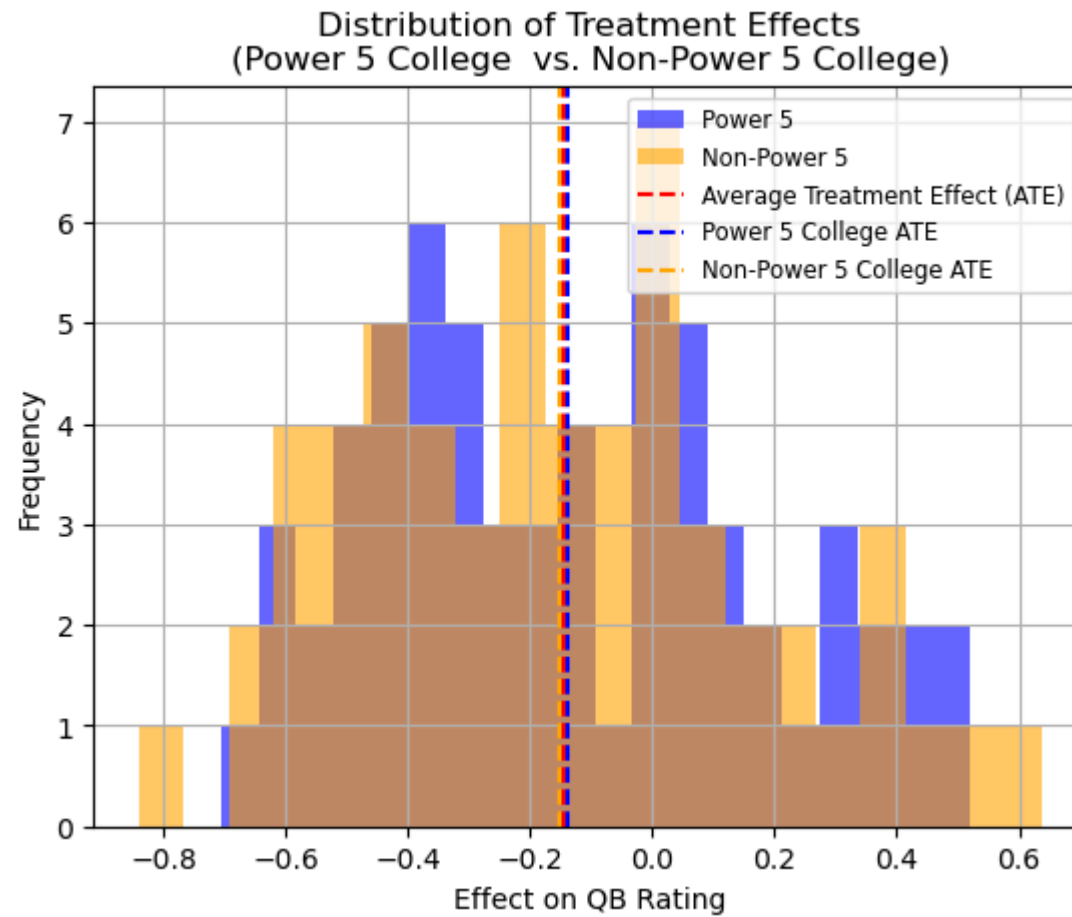
	point_estimate	stderr	zstat	pvalue	ci_lower	ci_upper
ATT	-9.661	11.014	-0.877	0.38	-31.249	11.926

Doubly Robust ATT(T=1) on Training Data Results

	point_estimate	stderr	zstat	pvalue	ci_lower	ci_upper
ATT	-7.488	8.682	-0.862	0.388	-24.504	9.528

Histogram of treatment effects

```
In [28]: 1 # Separate based on actual Power 5 status
2 tau_power5 = tau_hat[D == 1]
3 tau_nonpower5 = tau_hat[D == 0]
4
5 #plot histogram
6 plt.hist(tau_power5, bins=20, alpha=0.6, label="Power 5", color="blue")
7 plt.hist(tau_nonpower5, bins=20, alpha=0.6, label="Non-Power 5", color="orange")
8 #plot average treatment effect line
9 plt.axvline(np.mean(tau_hat), color="red", linestyle="--", label= "Average Treatment Effect (ATE)")
10 plt.axvline(np.mean(tau_power5), color="blue", linestyle="--", label= "Power 5 College ATE")
11 plt.axvline(np.mean(tau_nonpower5), color="orange", linestyle="--", label= "Non-Power 5 College ATE")
12
13 #plot axis, legend and grid
14 plt.title("Distribution of Treatment Effects\n(Power 5 College vs. Non-Power 5 College)")
15 plt.xlabel("Effect on QB Rating")
16 plt.ylabel("Frequency")
17 plt.legend(fontsize = "small")
18 plt.grid(True)
19
20 #store it as pdf in a folder
21 plt.savefig(FIG+"QBRatingATE.pdf")
22
23 #output diagram
24 plt.show()
25
26 #print the values of ATE
27 print("ATE (All):", np.mean(tau_hat))
28 print("ATE (Power 5):", np.mean(tau_power5))
29 print("ATE (Non-Power 5):", np.mean(tau_nonpower5))
```

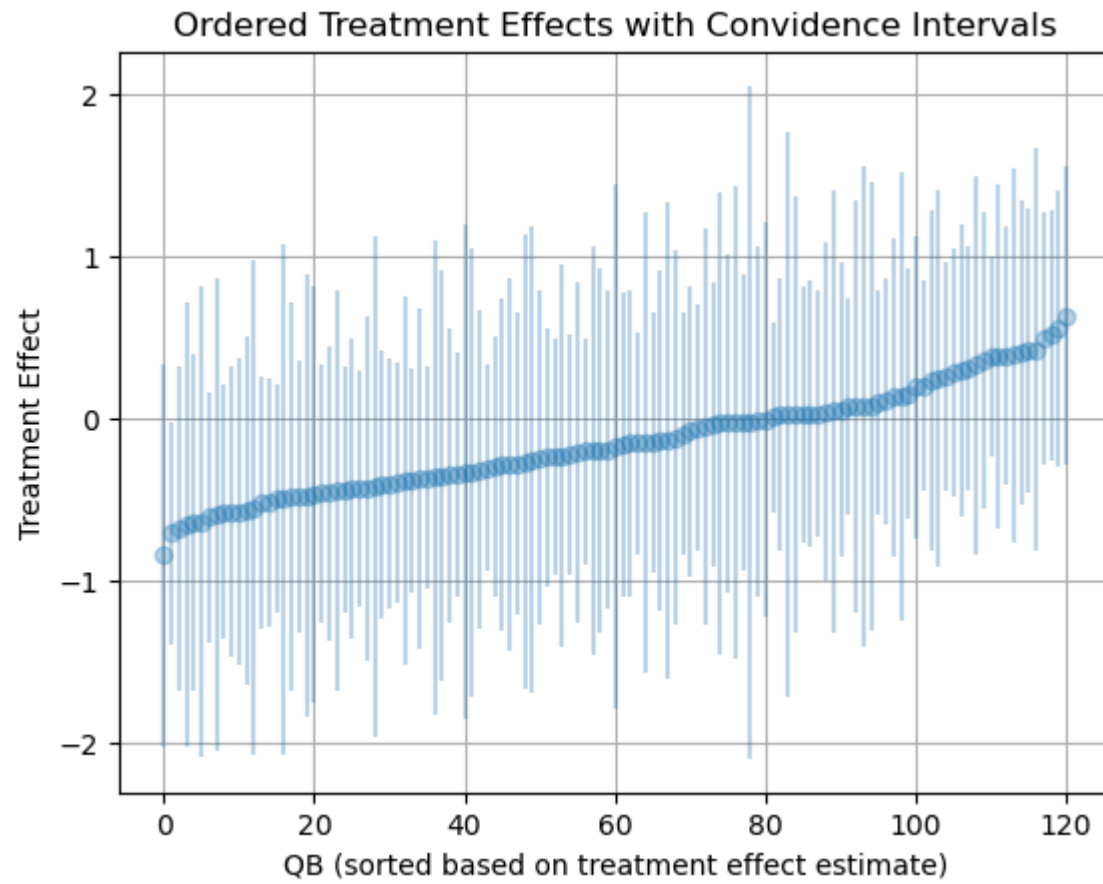
ATE (All): -0.14431797340206004

ATE (Power 5): -0.1381118511303841

ATE (Non-Power 5): -0.1510591062143977

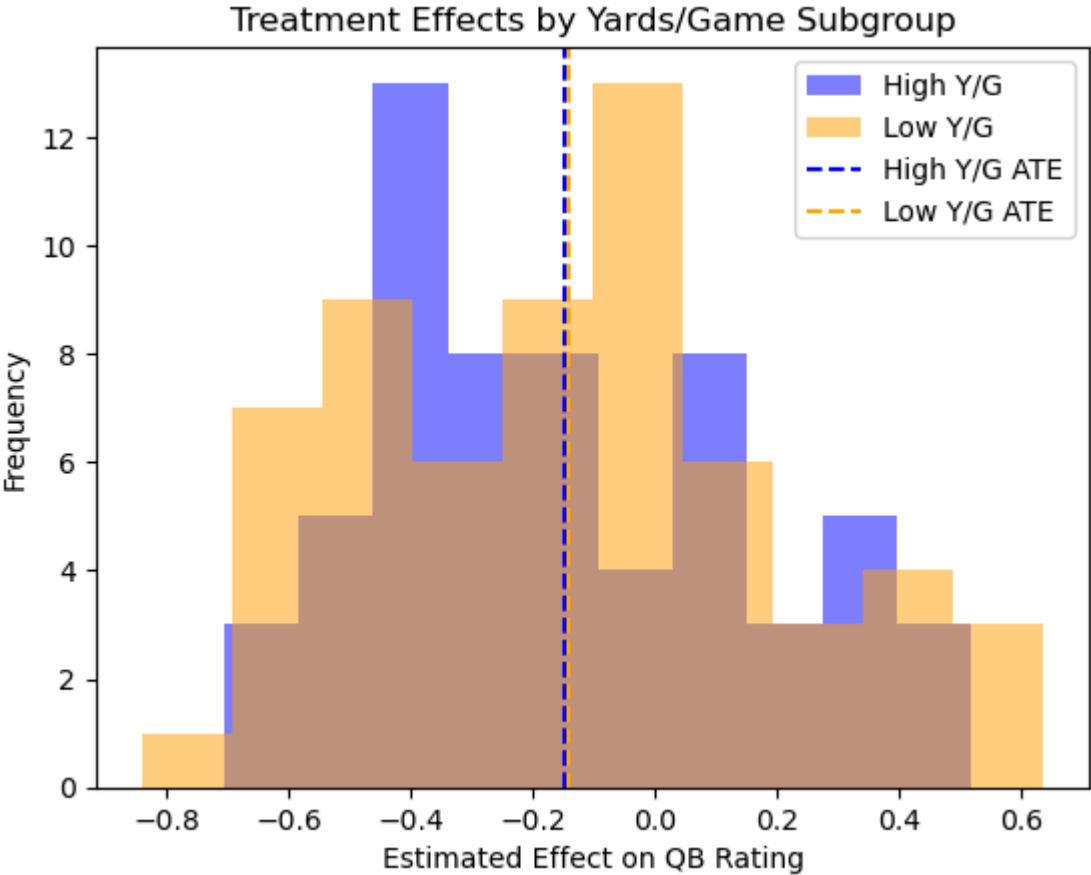
Ordered treatment effects

```
In [29]: 1 #flatten treatment effect to estimate effects of a 1D array
2 effects = tau_hat.flatten()
3
4 #calculate confidence intervals
5 ci = tau_hat_se[1] - effects
6
7 #sort treatment effects for a clearer visualisation
8 sorted_index= np.argsort(effects)
9
10 #plot treatment effects with error bars which are confidence intervals,
11 #each point is a QB and sorted by their estimated treatment effect
12 plt.errorbar(
13     np.arange(len(effects)),
14     effects[sorted_index],
15     yerr=ci[sorted_index],
16     fmt='o',
17     alpha=0.3)
18
19 #plot axis
20 plt.title("Ordered Treatment Effects with Convidence Intervals")
21 plt.xlabel("QB (sorted based on treatment effect estimate)")
22 plt.ylabel("Treatment Effect")
23 plt.grid(True)
24
25 #store as pdf
26 plt.savefig(FIG+"CIOTE.pdf")
27
28 #output diagram
29 plt.show()
```



Subgroup Treatment effect

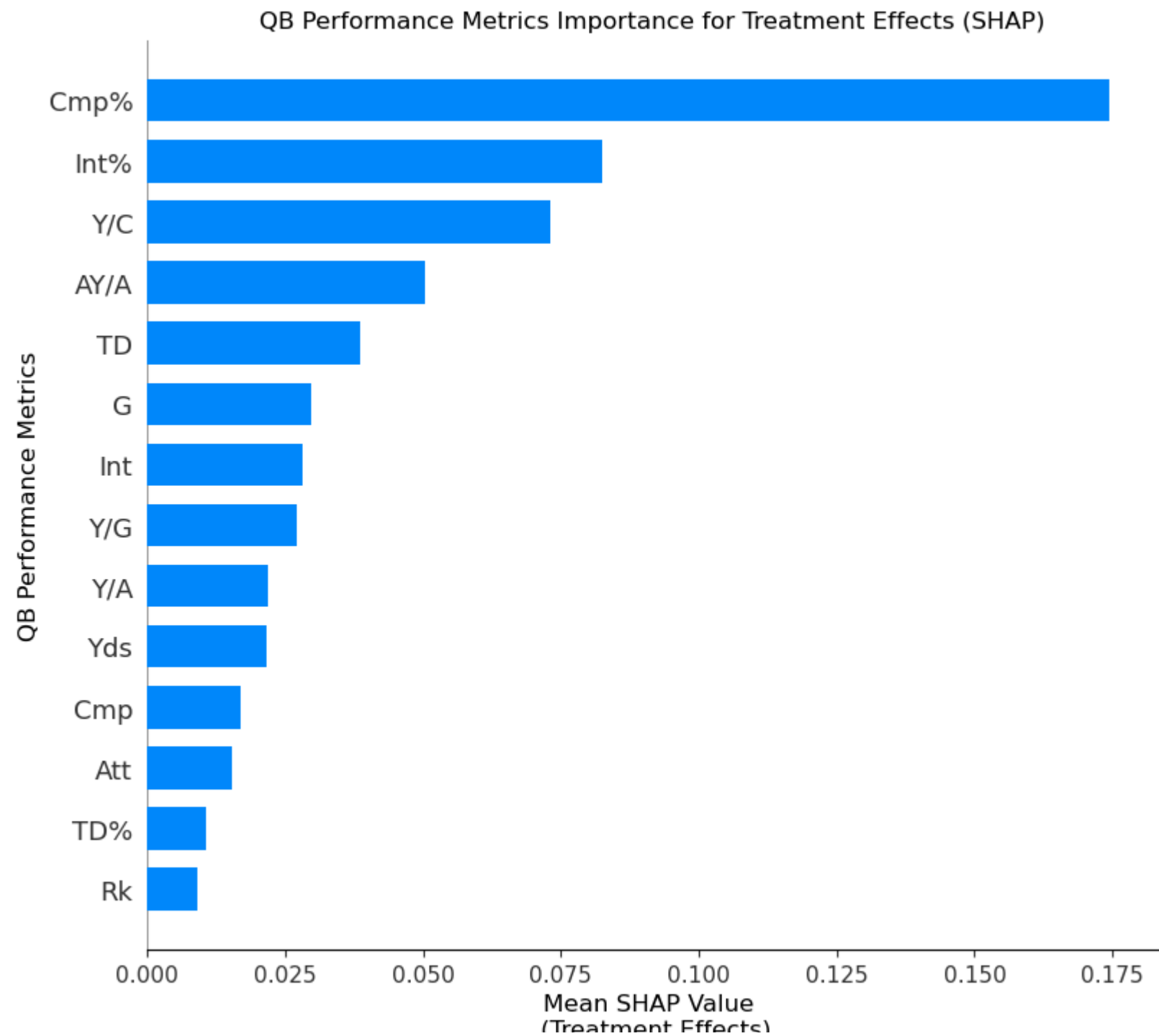
```
In [30]: 1 #make a copy to avoid SettingWithCopyWarning
2 qb_data_filtered = qb_data_filtered.copy()
3
4 #Convert "Y/G" Column to numerical as it is stored as a string
5 qb_data_filtered["Y/G"] = pd.to_numeric(qb_data_filtered["Y/G"])
6
7 #split data into high yards per game and low yards per game using the median
8 high_yards_per_game = qb_data_filtered["Y/G"] > qb_data_filtered["Y/G"].median()
9 low_yards_per_game = ~high_yards_per_game
10
11 #get treatment effects of both groups
12 tau_high = cf.effect(X[high_yards_per_game])
13 tau_low = cf.effect(X[low_yards_per_game])
14
15 #plot treatment effect for each group using histogram
16 plt.hist(tau_high, alpha=0.5, color="blue", label="High Y/G")
17 plt.hist(tau_low, alpha=0.5, color="orange", label="Low Y/G")
18
19 # Plot ATE lines for each group
20 plt.axvline(np.mean(tau_high), color="blue", linestyle="--", label="High Y/G ATE")
21 plt.axvline(np.mean(tau_low), color="orange", linestyle="--", label="Low Y/G ATE")
22
23 #plot labels
24 plt.legend()
25 plt.title("Treatment Effects by Yards/Game Subgroup")
26 plt.xlabel("Estimated Effect on QB Rating")
27 plt.ylabel("Frequency")
28
29 #store as pdf
30 plt.savefig(FIG+"TreatmentEffectsSubgroup.pdf")
31
32 #output diagram
33 plt.show()
```



SHAP Summary

```
In [31]: 1 import shap
2
3 # Ensure X is all numeric
4 X = X.apply(pd.to_numeric, errors="coerce").dropna()
5
6 #create a function to get treatment effect from cf model
7 def model_for_shap(X_input):
8     return cf.const_marginal_effect(X_input).flatten()
9
10 #create a background data set
11 background = shap.sample(X, 100, random_state=0)
12
13 # Create SHAP explainer
14 explainer = shap.Explainer(model_for_shap, background)
15
16 # Compute SHAP values
17 shap_values = explainer(X)
18
19 # Plot SHAP summary bar chart
20 shap.summary_plot(
21     shap_values.values,
22     features=X,
23     feature_names=X.columns,
24     plot_type="bar",
25     show=False)
26
27 #plot axis
28 plt.title("QB Performance Metrics Importance for Treatment Effects (SHAP)")
29 plt.xlabel("Mean SHAP Value \n (Treatment Effects)", fontsize=12)
30 plt.ylabel("QB Performance Metrics", fontsize = 12)
31
32
33 #store as pdf
34 plt.savefig(FIG+"SHAPvalues.pdf")
35
36 #out
37 plt.show()
38
```

PermutationExplainer explainer: 122it [01:19, 1.35it/s]



In [32]:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 from sklearn.linear_model import LinearRegression
4
5 # Convert "Cmp%" values to numeric as they are stored as strings
6 cmp_percent = pd.to_numeric(qb_data_filtered.loc[X.index, "Cmp%"], errors='coerce').values
7
8 # Get predicted treatment effects
9 te_predicts = cf.const_marginal_effect(X).flatten()
10
11 # Remove any missing values so we can plot properly
12 valid_data = ~np.isnan(cmp_percent) & ~np.isnan(te_predicts)
13 cmp_percent = cmp_percent[valid_data]
14
15 #reverse the data as retrieved treatment effect is reversed
16 te_predicts = 1 - te_predicts[valid_data]
17
18 # generate a simple linear regression model to show relation between "cmp%" and predicted treatment effect
19 model = LinearRegression()
20 cmp_resaped = cmp_percent.reshape(-1, 1)
21 model.fit(cmp_resaped, te_predicts)
22 line_x = np.linspace(cmp_percent.min(), cmp_percent.max(), 100).reshape(-1, 1)
23 line_y = model.predict(line_x)
24
25 # Plot the scatter graph and the trend line
26 plt.figure(figsize=(8, 6))
27 plt.scatter(cmp_percent, te_predicts, label="QBs")
28 plt.plot(line_x, line_y, color="red", linewidth=2, label="General Correlation")
29
30 # plot axis, title and legend
31 plt.xlabel("Completion Percentage (Cmp%)")
32 plt.ylabel("Predicted Treatment Effect")
33 plt.title("Correlation between Cmp% & Predicted Treatment Effect")
34 plt.legend()
35
36 # store as pdf
37 plt.savefig(FIG + "Cmp_vs_TreatmentEffect_MPL.pdf")
38
39 #output
```



```
40 plt.show()
```

