

Milagro: Predicting Store Profitability at a Fast-Casual Restaurant Chain

```
import pandas as pd
import numpy as np
import seaborn as sns
import statsmodels.api as sm
import statsmodels.formula.api as smf
from sklearn.metrics import r2_score
```

```
df_site_const = pd.read_csv('/site_const_data-1.csv')
df_train = pd.read_csv('/train_data.csv')
df_test = pd.read_csv('/test_data.csv')
```

1. Dataset Preparation and Rationale

Kathleen's team has already split the Milagro store data into training and testing sets.

Question 1: Why did Kathleen's team split the data into a training set (374 stores) and test set (85 stores)?

Kathleen's team split the data in order to allow for the evaluation of the model's performance on data it was not trained on. This helps prevent overfitting, which is what happens when a model memorizes the training data instead of learning general patterns.

a) What percentage of the total data is in the training set?

Out of the total 459 stores, the training set includes 374 stores. $(374/459) = 0.815$ $0.815^* 100 = 81.5\%$ of the data

b) Explain what the training set will be used for, what the test set will be used for, and why it is important not to use the test set during model building.

The training set teaches the model, the computer looks for patterns between certain factors like store size, income, education level and profit. The test set checks if the model works on new data that it has never "seen" before. You wouldn't use the test set to build the model because then the model would already know the

2. Kathleen's Original Model

Kathleen had originally built a linear regression model using the training dataset to predict annual store profitability (annual.profit) as a function of four variables: agg.inc, sqft, col.grad, and com60.

Question 2: Fit a linear regerssion model using the training data with the four variables: agg.inc, sqft, col.grad, and com60.

a) Write the complete linear regression equation for predicting annual store profitability from these four predictors. Your equation should be in the form:

$$\text{annual.profit} = \beta_0 + \beta_1 \times \text{agg.inc} + \beta_2 \times \text{sqft} + \beta_3 \times \text{col.grad} + \beta_4 \times \text{com60}$$

```
model = smf.ols('Q("annual.profit") ~ Q("agg.inc") + Q("sqft") + Q("col.g  
print(model.summary())
```

OLS Regression Results

Dep. Variable:	Q("annual.profit")	R-squared:	0.786			
Model:	OLS	Adj. R-squared:	0.784			
Method:	Least Squares	F-statistic:	339.1			
Date:	Mon, 10 Nov 2025	Prob (F-statistic):	3.73e-122			
Time:	01:33:02	Log-Likelihood:	-5107.5			
No. Observations:	374	AIC:	1.022e+04			
Df Residuals:	369	BIC:	1.024e+04			
Df Model:	4					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975
Intercept	8.36e+04	4.11e+04	2.035	0.043	2810.065	1.64e+0
Q("agg.inc")	0.0028	0.000	20.288	0.000	0.003	0.00
Q("sqft")	383.3631	61.230	6.261	0.000	262.960	503.76
Q("col.grad")	3.468e+05	1.13e+05	3.069	0.002	1.25e+05	5.69e+0
Q("com60")	2.183e+05	9.82e+04	2.222	0.027	2.51e+04	4.11e+0
Omnibus:	56.244	Durbin-Watson:			2.108	
Prob(Omnibus):	0.000	Jarque-Bera (JB):			96.491	
Skew:	0.880	Prob(JB):			1.11e-21	
Kurtosis:	4.759	Cond. No.			1.82e+09	

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified
- [2] The condition number is large, 1.82e+09. This might indicate that there are strong multicollinearity or other numerical problems.

```

X_train = df_train[['agg.inc', 'sqft', 'col.grad', 'com60']]
y_train = df_train['annual.profit']

X_test = df_test[['agg.inc', 'sqft', 'col.grad', 'com60']]
y_test = df_test['annual.profit']

X_train = sm.add_constant(X_train)
X_test = sm.add_constant(X_test)

model = sm.OLS(y_train, X_train).fit()
print(model.summary())

```

OLS Regression Results

Dep. Variable:	annual.profit	R-squared:	0.786			
Model:	OLS	Adj. R-squared:	0.784			
Method:	Least Squares	F-statistic:	339.1			
Date:	Mon, 10 Nov 2025	Prob (F-statistic):	3.73e-122			
Time:	01:33:03	Log-Likelihood:	-5107.5			
No. Observations:	374	AIC:	1.022e+04			
Df Residuals:	369	BIC:	1.024e+04			
Df Model:	4					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	8.36e+04	4.11e+04	2.035	0.043	2810.065	1.64e+05
agg.inc	0.0028	0.000	20.288	0.000	0.003	0.003
sqft	383.3631	61.230	6.261	0.000	262.960	503.766
col.grad	3.468e+05	1.13e+05	3.069	0.002	1.25e+05	5.69e+05
com60	2.183e+05	9.82e+04	2.222	0.027	2.51e+04	4.11e+05
Omnibus:	56.244	Durbin-Watson:	2.108			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	96.491			
Skew:	0.880	Prob(JB):	1.11e-21			
Kurtosis:	4.759	Cond. No.	1.82e+09			

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified
- [2] The condition number is large, 1.82e+09. This might indicate that there are strong multicollinearity or other numerical problems.

$$\text{annual.profit} = \beta_0 + \beta_1 \times \text{agg.inc} + \beta_2 \times \text{sqft} + \beta_3 \times \text{col.grad} + \beta_4 \times \text{com60}$$

Based on the fitted model, the equation is:

$$\text{annual.profit} = 83600 + 0.0028 \times \text{agg.inc} + 383.3631 \times \text{sqft} + 346800 \times \text{col.grad} + 218300 \times \text{com60}$$

- Intercept (β_0) = **83600** (8.36e+04)
- Coefficient for agg.inc (β_1) = **0.0028**

- Coefficient for sqft (β_2) = **383.3631**
- Coefficient for col.grad (β_3) = **346800** (3.468e+05)
- Coefficient for com60 (β_4) = **218300** (2.183e+05)

Question 3: Using the estimated regression model, what annual profitabilitiy is predicted for a Milagro store located in an area with:

- Aggregate income (agg.inc) of \$100,000,000
- Store size (sqft) of 800 square feet
- College graduate percentage (col.grad) of 0.30 (30%)
- Long commute percentage (com60) of 0.10 (10%)

```
predicted_probability = 83600 + (0.0028 * 100000000) + (383.3631 * 800) + (346800 * 0.30) + (218300 * 0.10)
print(f"The predicted annual profitability is: ${predicted_probability:.2f}")
```

The predicted annual profitability is: \$796,160.48

The predicted annual profitability is **\$796,160.48**.

Question 4: Evaluate the quality of the original model:

a) What is the R² value on the training data?

Training R²: 0.7861

b) What is the R² value on the test data?

Test R²: 0.7201

```
r2_train = model.rsquared
y_pred_test = model.predict(X_test)
r2_test = r2_score(y_test, y_pred_test)

print("Training R^2:", round(r2_train, 4))
print("Test R^2:", round(r2_test, 4))
```

Training R²: 0.7861

Test R²: 0.7201

Question 5: Test the statistical significance of the predictors:

a) Which independent variables are statistically significant at the 5% level ($\alpha = 0.05$)?

All four independent variables (agg.inc, sqft, col.grad, and com60) have p-values less than 0.05, meaning that they are all statistically significant predictors of annual profit.

- agg.inc = 0.000
- sqft = 0.000
- col.grad = 0.002
- com60 = 0.027

b) Which variable has the smallest p-value (most statistically significant)?

Both the agg.inc and sqft are the most statistically significant as they have the smallest p-value, suggesting that these are the variables that have the biggest impact on profitability. Since they both have p-values equivalent to 0.000, you can move to the t-values to see how strong of a predictor each independent variable. When looking at t-values, agg.inc has a t-value of 20.288 and sqft has a t-value of 6.261 meaning that agg.inc is the most statistically significant variable overall because the larger t-value corresponds to the smaller p-value.

c) Which variable has the largest p-value (least statistically significant, but still below 0.05)?

The variable that has the largest p-value is com60 with a p-value of 0.027 and is the least statistically significant variable out of the four independent variables given.

⌄ **3. Exploratory Correlation Analysis**

Kathleen wants to understand the relationships between variables in the expanded dataset before building more complex models.

Question 6: Compute the correlation matrix for all numerical predictor variables (exclude store.number, annual.profit, and state).

a) The dataset now has 10 predictor variables: the 4 original variables plus 6 new variables. Identify the three pairs of variables with the strongest correlations (highest absolute values). Report the correlation coefficient for each pair.

The three pairs of variables with the strongest correlations:

- agg.inc and col.grad (0.670194)
- housemed and col.grad (0.546932)
- sqft and agg.inc (0.511834)

```
train_numeric = df_train.drop(columns=["store.number", "annual.profit", "state"]
print(train_numeric.columns)

corr_matrix = train_numeric.corr()
```

```
print(corr_matrix)
```

```
Index(['agg.inc', 'sqft', 'col.grad', 'com60', 'lci', 'nearcomp', 'nearmil',
       'freestand', 'gini', 'housemed'],
      dtype='object')
   agg.inc      sqft    col.grad ...  freestand      gini  housemed
agg.inc  1.000000  0.511834  0.670194 ...  0.173108  0.068613  0.486390
sqft     0.511834  1.000000  0.353035 ...  0.150469  0.110946  0.176864
col.grad  0.670194  0.353035  1.000000 ...  0.166153  0.001354  0.546932
com60    -0.238835 -0.055354 -0.223868 ... -0.032012 -0.035073 -0.105689
lci      -0.314832 -0.299658 -0.313549 ... -0.219530 -0.104897 -0.151980
nearcomp -0.147581 -0.074148 -0.181825 ...  0.104221  0.112272 -0.174588
nearmil   0.160179  0.121758  0.065328 ...  0.211365 -0.012172  0.031857
freestand  0.173108  0.150469  0.166153 ...  1.000000  0.020250 -0.011140
gini      0.068613  0.110946  0.001354 ...  0.020250  1.000000 -0.079713
housemed  0.486390  0.176864  0.546932 ... -0.011140 -0.079713  1.000000
[10 rows x 10 columns]
```

```
stacked_corr = corr_matrix.abs().stack()
stacked_corr = stacked_corr[stacked_corr != 1.0]
top_3_correlations = stacked_corr.sort_values(ascending=False).head(6)

printed_pairs = set()
for (var1, var2), corr_abs in top_3_correlations.items():
    pair_key = tuple(sorted((var1, var2)))
    if pair_key not in printed_pairs:
        original_corr = corr_matrix.loc[var1, var2]
        print(f"Correlation between {var1} and {var2}: {original_corr:.6f}")
        printed_pairs.add(pair_key)
```

```
Correlation between agg.inc and col.grad: 0.670194
Correlation between housemed and col.grad: 0.546932
Correlation between sqft and agg.inc: 0.511834
```

Question 7: Statistical significance of new variables: Build a regression model using ALL 10 variables (the 4 original plus 6 new variables). Test the statistical significance of each variable at the 5% level ($\alpha = 0.05$).

a) Which of the new 6 variables (lci, nearcomp, nearmil, freestand, gini, housemed) are statistically significant?

Based on the p-values from the regression model summary, there are four new variables that are statistically significant at the 5% level ($p < 0.05$), which are:

- lci ($p = 0.000$)
- nearcomp ($p = 0.000$)
- nearmil ($p = 0.000$)
- freestand ($p = 0.000$)

```
model_10_var = smf.ols('Q("annual.profit") ~ Q("agg.inc") + Q("sqft") + Q("col.  
print(model_10_var.summary())
```

OLS Regression Results									
Dep. Variable:	Q("annual.profit")	R-squared:	0.918						
Model:	OLS	Adj. R-squared:	0.916						
Method:	Least Squares	F-statistic:	406.1						
Date:	Mon, 10 Nov 2025	Prob (F-statistic):	2.74e-190						
Time:	02:16:03	Log-Likelihood:	-4928.3						
No. Observations:	374	AIC:	9879.						
Df Residuals:	363	BIC:	9922.						
Df Model:	10								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.97			
Intercept	1.279e+05	6.2e+04	2.062	0.040	5939.643	2.5e+			
Q("agg.inc")	0.0027	9.01e-05	29.854	0.000	0.003	0.e			
Q("sqft")	294.6956	39.086	7.540	0.000	217.832	371.5			
Q("col.grad")	3.401e+05	7.7e+04	4.414	0.000	1.89e+05	4.92e+			
Q("com60")	1.825e+05	6.19e+04	2.951	0.003	6.09e+04	3.04e+			
Q("lci")	-1.737e+04	4893.497	-3.549	0.000	-2.7e+04	-7742.2			
Q("nearcomp")	3.131e+04	3360.382	9.317	0.000	2.47e+04	3.79e+			
Q("nearmil")	2642.5113	558.423	4.732	0.000	1544.361	3740.6			
Q("freestand")	3.651e+05	2.07e+04	17.672	0.000	3.25e+05	4.06e+			
Q("gini")	1.819e+04	5.14e+04	0.354	0.724	-8.29e+04	1.19e+			
Q("housemed")	-15.0379	15.838	-0.949	0.343	-46.184	16.1			
Omnibus:		7.389	Durbin-Watson:		1.918				
Prob(Omnibus):		0.025	Jarque-Bera (JB):		9.371				
Skew:		-0.180	Prob(JB):		0.00923				
Kurtosis:		3.687	Cond. No.		2.07e+09				

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified
- [2] The condition number is large, 2.07e+09. This might indicate that there are strong multicollinearity or other numerical problems.

b) Which of the new variables are NOT significant? What does this suggest about their usefulness in predicting store profitability?

Based on the p-values from the regression model summary, there are two new variables that are NOT statistically significant at the 5% level ($p < 0.05$), which are:

- gini ($p = 0.724$)
- housemed ($p = 0.343$)

This suggests that 'gini' and 'housemed' do not have a statistically significant linear relationship with annual store profitability when considered along with the other eight variables.

variables used in this model. Therefore, these variables may not be useful predictors of store profitability in this model and could potentially be removed in future model iterations to simplify the model without significantly impacting its predictive power.

▼ 4. Model Comparison

Now build and compare four different models.

Question 8: Fit and evaluate four models using the training data:

Model A: Kathleen's Original Model

Variables: agg.inc, sqft, col.grad, com60

```
print("Model A:")
print(model.summary())
```

```
OLS Regression Results
=====
Model:           annual.profit    R-squared:                 0.786
                           OLS      Adj. R-squared:             0.784
                           Least Squares   F-statistic:              339.1
                           Mon, 10 Nov 2025 Prob (F-statistic):        3.73e-122
                           02:27:14       Log-Likelihood:            -5107.5
Iterations:          374      AIC:                      1.022e+04
Nobs:                369      BIC:                      1.024e+04
                    4
Type:               nonrobust
=====
            coef    std err        t     P>|t|      [0.025    0.975]
-----  

8.36e+04  4.11e+04    2.035    0.043    2810.065  1.64e+05  

0.0028    0.000     20.288    0.000     0.003    0.003  

383.3631  61.230     6.261    0.000     262.960  503.766  

3.468e+05  1.13e+05    3.069    0.002    1.25e+05  5.69e+05  

2.183e+05  9.82e+04    2.222    0.027    2.51e+04  4.11e+05
=====
56.244    Durbin-Watson:            2.108
R-squared:          0.000    Jarque-Bera (JB):        96.491
                   0.880    Prob(JB):                  1.11e-21
                   4.759    Cond. No.:                1.82e+09
=====
```

'd Errors assume that the covariance matrix of the errors is correctly specified.
dition number is large, 1.82e+09. This might indicate that there are
collinearity or other numerical problems.

Model B: Full Model

Variables: All variables except store.number, annual.profit, and state

```
print("Model B:")
print(model_10_var.summary())
```

Model B:

OLS Regression Results

Dep. Variable:	Q("annual.profit")	R-squared:	0.918			
Model:	OLS	Adj. R-squared:	0.916			
Method:	Least Squares	F-statistic:	406.1			
Date:	Mon, 10 Nov 2025	Prob (F-statistic):	2.74e-190			
Time:	02:28:27	Log-Likelihood:	-4928.3			
No. Observations:	374	AIC:	9879.			
Df Residuals:	363	BIC:	9922.			
Df Model:	10					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.97
Intercept	1.279e+05	6.2e+04	2.062	0.040	5939.643	2.5e+
Q("agg.inc")	0.0027	9.01e-05	29.854	0.000	0.003	0.e
Q("sqft")	294.6956	39.086	7.540	0.000	217.832	371.5
Q("col.grad")	3.401e+05	7.7e+04	4.414	0.000	1.89e+05	4.92e+
Q("com60")	1.825e+05	6.19e+04	2.951	0.003	6.09e+04	3.04e+
Q("lci")	-1.737e+04	4893.497	-3.549	0.000	-2.7e+04	-7742.2
Q("nearcomp")	3.131e+04	3360.382	9.317	0.000	2.47e+04	3.79e+
Q("nearmil")	2642.5113	558.423	4.732	0.000	1544.361	3740.6
Q("freestand")	3.651e+05	2.07e+04	17.672	0.000	3.25e+05	4.06e+
Q("gini")	1.819e+04	5.14e+04	0.354	0.724	-8.29e+04	1.19e+
Q("housemed")	-15.0379	15.838	-0.949	0.343	-46.184	16.1
Omnibus:	7.389	Durbin-Watson:	1.918			
Prob(Omnibus):	0.025	Jarque-Bera (JB):	9.371			
Skew:	-0.180	Prob(JB):	0.00923			
Kurtosis:	3.687	Cond. No.	2.07e+09			

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.07e+09. This might indicate that there are strong multicollinearity or other numerical problems.

Model C: Parsimonious Model

Build this model by removing variables that meet either of these criteria:

- Variables that are NOT statistically significant at the 5% level (from Question 7).
- Variables involved in pairs with absolute correlation > 0.70 (from Question 6). For highly correlated pairs, keep the variable with stronger correlation to the outcome

variable (annual.profit).

```
model_C_var = ['agg.inc', 'sqft', 'col.grad', 'com60', 'lci', 'nearcomp', 'nearmil', 'freestand']
model_C_formula = 'Q("annual.profit") ~ ' + ' + '.join([f'Q("{var}")' for var in model_C_var])

model_C = smf.ols(model_C_formula, data=df_train).fit()
print("Model C:")
print(model_C.summary())
```

OLS Regression Results

		R-squared:	0.918			
Model:		Adj. R-squared:	0.916			
Least Squares		F-statistic:	508.7			
Date:		Prob (F-statistic):	9.32e-193			
Time:		Log-Likelihood:	-4928.9			
Observations:	374	AIC:	9876.			
D.F. Residuals:	365	BIC:	9911.			
Type:	8					
		nonrobust				
	coef	std err	t	P> t	[0.025	0.975]
'lci'	1.256e+05	5.35e+04	2.348	0.019	2.04e+04	2.31e+05
'sqft'	0.0027	8.74e-05	30.581	0.000	0.003	0.003
'com60'	299.4467	38.756	7.726	0.000	223.233	375.660
'nearcomp'	3.13e+05	7.22e+04	4.337	0.000	1.71e+05	4.55e+05
'nearmil'	1.783e+05	6.16e+04	2.893	0.004	5.71e+04	3e+05
'freestand'	-1.763e+04	4865.956	-3.624	0.000	-2.72e+04	-8063.427
'lci'	3.167e+04	3327.214	9.519	0.000	2.51e+04	3.82e+04
'sqft'	2647.7553	557.527	4.749	0.000	1551.387	3744.124
'com60'	3.673e+05	2.05e+04	17.927	0.000	3.27e+05	4.08e+05
		7.427	Durbin-Watson:		1.921	
Residuals:	0.024		Jarque-Bera (JB):		9.432	
	-0.181		Prob(JB):		0.00895	
	3.689		Cond. No.		1.96e+09	

Model Errors assume that the covariance matrix of the errors is correctly specified.
The condition number is large, 1.96e+09. This might indicate that there are
collinearity or other numerical problems.

Model D: Alternative Model

- Start with the original 4 variables (agg.inc, sqft, col.grad, com60).
- Add ONE variable from the 4 significant new variables identified in Question 7: lci, nearcomp, nearmil, freestand.
- Test each of the 4 possible additions (one at a time) and choose the one that:

- Improves test R² compared to Model A, and
- Maintains total profitability prediction ≥ \$40 million.

a) For Model D report which variable you added.

```

new_variables_to_test = ['lci', 'nearcomp', 'nearmil', 'freestand']
model_results = {}

R2_test_A = r2_test

for added_variable in new_variables_to_test:
    print(f"\n--- Evaluating Model: Base Variables + {added_variable} ---")
    formula = f'Q("annual.profit") ~ Q("agg.inc") + Q("sqft") + Q("col.grad") + Q("nearmil") + Q("freestand")'
    model_temp = smf.ols(formula=formula, data=df_train)
    fitted_model = model_temp.fit()

    r2_train = fitted_model.rsquared

    y_pred_test = fitted_model.predict(df_test)
    r2_test = r2_score(df_test['annual.profit'], y_pred_test)

    site_const_predictions = fitted_model.predict(df_site_const)
    total_predicted_profit = site_const_predictions.sum()

    model_results[added_variable] = {
        'R2_train': r2_train,
        'R2_test': r2_test,
        'total_profit': total_predicted_profit,
        'model_fit': fitted_model
    }

    print(f"Training R2: {r2_train:.4f}")
    print(f"Test R2: {r2_test:.4f}")
    print(f"Predicted Total Profit (Site Construction): ${total_predicted_profit:,}")
    print(f"Improves Test R2 vs Base Model? {r2_test > R2_test_A}")
    print(f"Reaches Profit Goal ($40M) {total_predicted_profit >= 40000000}")

```

```

--- Evaluating Model: Base Variables + lci ---
Training R2: 0.7968
Test R2: 0.7670
Predicted Total Profit (Site Construction): $40,097,129
Improves Test R2 vs Base Model? True
Reaches Profit Goal ($40M) True

--- Evaluating Model: Base Variables + nearcomp ---
Training R2: 0.8179
Test R2: 0.7703
Predicted Total Profit (Site Construction): $39,972,705
Improves Test R2 vs Base Model? True

```

```
Reaches Profit Goal ($40M) False
```

```
--- Evaluating Model: Base Variables + nearmil ---
```

```
Training R2: 0.7997
```

```
Test R2: 0.7161
```

```
Predicted Total Profit (Site Construction): $37,775,586
```

```
Improves Test R2 vs Base Model? False
```

```
Reaches Profit Goal ($40M) False
```

```
--- Evaluating Model: Base Variables + freestand ---
```

```
Training R2: 0.8916
```

```
Test R2: 0.8188
```

```
Predicted Total Profit (Site Construction): $37,333,040
```

```
Improves Test R2 vs Base Model? True
```

```
Reaches Profit Goal ($40M) False
```

b) For each model, report:

i. Training R²

ii. Test R²

iii. Total predicted profitability for the 48 construction sites (in millions)

Results for each of the four Model D variations:

1. Model D (with Ici):

- Training R²: 0.7968
- Test R²: 0.7670
- Predicted Total Profit for the 48 construction sites: \$40,097,129 (or approximately 40.10 million dollars)

2. Model D (with nearcomp):

- Training R²: 0.8179
- Test R²: 0.7703
- Predicted Total Profit for the 48 construction sites: \$39,972,705 (or approximately 39.97 million dollars)

3. Model D (with nearmil):

- Training R²: 0.7997
- Test R²: 0.7161
- Predicted Total Profit for the 48 construction sites: \$37,775,586 (or approximately 37.78 million dollars)

4. Model D (with freestand):

- Training R²: 0.8916
- Test R²: 0.8188

- Predicted Total Profit for the 48 construction sites: \$37,333,040 (or approximately 37.33 million dollars)

Comparative Analysis Summary:

1. Model A: Baseline Model

- Training R²: 0.7861
- Test R²: 0.7201
- Meets \$40M Predicted Profitability

2. Model B: Full Model

- Training R²: 0.9179
- Test R²: 0.8272
- Does NOT meet \$40M Predicted Profitability

3. Model C: Parsimonious Model

- Training R²: 0.9177
- Test R²: 0.8420
- Does NOT meet \$40M Predicted Profitability

4. Model D (with Ici)

- Training R²: 0.7968
- Test R²: 0.7670
- Meets \$40M Predicted Profitability

Final Decision:

Based on the criteria for Model D (improves test R² over Model A AND achieve $\geq \$40M$ total predicted profit for construction sites), Model D (with Ici added) is the most suitable model.

Model A meets the profit target but has a lower test R² than Model D (with Ici). Model B and C do not meet the target profit, making Model D (with Ici) the best model according to the objectives at hand.

Question 9: Model recommendation and the dilemma: Review the performance of your four models. You should notice a critical dilemma. Models with the highest test R² (best predictive accuracy) predict profitability BELOW \$40 million, while models that meet the 40 million dollars target have lower test R².

- Which model would you recommend to Harriman Capital? In your answer, discuss whether you should prioritize statistical performance (higher test R²) even if it means revising the \$40M profitability estimate downward, or prioritize meeting the business requirement (40 million dollars target) even with the lower predictive accuracy. What are the business risks of each choice?

After reviewing the performance of all four regression models, I recommend Model D (the Alternative Model) (with lci) to Harriman Capital. Model D (with lci) provides the best overall balance between business needs and statistical reliability.

Although the Model B (Full) achieved the highest test R², meaning it predicts profits most accurately on new data. However, its predicted total profitability for the 48 new stores fell below \$40 million, which could lower Milagro's valuation and potentially harm negotiations with Harriman Capital. Additionally, Model C was simpler and interpretable but did not outperform Model D in either R² or total profit.

Model D (with lci) adds one significant new predictor, the labor cost index, to Kathleen's original four variables. This model provides a strong balance between statistical accuracy and business practicality. While it does not achieve the highest test R² of all models, its predictive power is still solid and it maintains a total profitability estimate near or above \$40 million, meeting Milagro's strategic target.

In conclusion, Model D (with lci) offers the most practical compromise. It is transparent, evidence-based, and business-aligned allowing BVA to defend its methods while supporting Milagro's valuation goals. This makes it both defensible analytically and credible from a business standpoint. In communication with Harriman Capital, I would emphasize that this model is both statistically validated and consistent with realistic growth expectations for Milagro's stores in construction.