

- Preprocess data and make it amenable to machine learning
- Train and test out of the box models

In [20]: # Evaluate the model on the test dataset using the square root of the mean squared error (RMSE) met
from sklearn.metrics import mean_squared_error

```
y_pred = logistic_regression_model.predict(X_test_scaled_const)
```

```
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
```

```
print("RMSE:", rmse)
```

RMSE: 1434.3398492274823

Your Response: We may be off by 1434 CHF on average because the RMSE is 1434. This means that the model's predictions are off by 1434 CHF on average.

3. (2 pts) Report the R^2 score on the test dataset and interpret it.

In [21]: **from** sklearn.metrics **import** r2_score

```
r2 = r2_score(y_test, y_pred)
```

```
print("R2 score on test dataset:", r2)
```

R2 score on test dataset: 0.8122727616953431

Your Response: The R^2 score on test set is 0.812, which means that 81.2% of the variance in the target variable (Price) can be explained by the features in the model.

Comments: If the student reports the R^2 score on the training set, 3 points will be deducted. The remaining 1 point will be deducted if the interpretation is incorrect.

4. (2 pts) Which features are statistically significant at a 5% significance level?

In [22]: p_values = logistic_regression_model.pvalues

```
significant_features = p_values[p_values < 0.05].index
```

```
# comment:
```

```
# Technically, the constant term is not a feature, but we don't deduct points for this.
```

```
# If the student reports the X1, X2 etc instead of the actual column names, 1 point should be deducted
```

```
significant_features = significant_features.drop('const')
```

```
print("Significant features:", significant_features)
```

Significant features: Index(['Age', 'Mileage', 'HP', 'CC', 'Weight', 'FuelType_Diesel',
1
1),
dtype='object')

5. (2 pts) Determine which two feature have the highest coefficient? What does it imply?

Your Response:

- Age and Weight.
- These two features have the most significant impact on the price of the car.

comments:

- constant term is not a feature, 1 point should be deducted if the student includes it in the answer.

2.3 (2 pts): Improvement Discussion

- Suggest a few additional features that could potentially explain this remaining variance in the data (at least 2 features).

Your Response:

possible answers: damage condition, optional upgrade, maintenance history, Previous owner, etc.

But car brand is a **wrong** answer because this dataset only contains data for Toyota Corolla cars! Car model is fine because Toyota Corolla has different models. If the student suggests car brand, 1 point will be deducted. The remaining 1 point depends on if there is another valid feature suggestion.

2.4 (2 pts): Identifying Confounding Variables

The feature "Weight" shows a very low p-value and a high coefficient, but it doesn't seem to be a major factor for customers buying a second-hand car. You go to your mentor Tim to discuss this issue. Indeed, Tim suggests that never in his career has he seen a customer who asked for the weight of a car before buying it. You suspect that there might be a confounding variable that is correlated with the car's weight and significantly influences its price.

- Suggest a possible confounding variable that may be correlated with the car's weight and significantly influences its price (it doesn't need to be a variable in the dataset). Explain why this variable could be a confounding variable.

Your Response:

possible answers: The Toyota Corolla has different versions, such as the Corolla Hatchback, Corolla Sedan, and Corolla Touring Sports. The weight of the car might be correlated with the version of the car. A heavier version is usually a bigger car with more features, which could explain the higher price.

2.5 (2 pts): Adding an Inverse Mileage Term

From the previous scatter plot, the relationship between car price and mileage appears non-linear, with a steep price drop initially and then a flattening. A suitable approach to model this behavior is by incorporating an inverse term of mileage.

- Add the inverse mileage term to the model and retrain it using the code provided. Print the model summary and interpret the effect of the inverse mileage term.

In [23]: X['Mileage_inverse'] = 1/X['Mileage']

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
scaler = StandardScaler()
```

```
X_train = scaler.fit_transform(X_train)
```

```
X_test = scaler.transform(X_test)
```

```
X_train = sm.add_constant(X_train)
```

```
X_test = sm.add_constant(X_test)
```

```
logistic_regression_model = sm.OLS(y_train, X_train).fit()
```

```
print(logistic_regression_model.summary())
```

```
y_pred = logistic_regression_model.predict(X_test)
```

```
r2 = r2_score(y_test, y_pred)
```

```
print("R2 score on test dataset:", r2)
```

```
# comment:
```

```
# The student can report either the R2 score showed in the summary or the one calculated using the
```

```
# The R2 score should be pretty much the same as the one without the inverse feature. This means that
```

```
# is not bringing any additional information to the model.
```

```
=====
```

```
OLS Regression Results
```

	coef	std err	t	P> t	[0.025	0.975]
const	9441.4983	36.685	257.367	0.000	9369.520	9513.476
x1	-1974.6675	51.769	-38.144	0.000	-2076.241	-1873.094
x2	-511.6765	52.819	-9.687	0.000	-615.310	-408.043
x3	472.3240	94.256	5.011	0.000	287.389	657.259
x4	-477.8440	111.891	-4.301	0.000	-695.811	-259.878
x5	1166.8577	76.853	15.343	0.000	1017.636	1316.079
x6	441.3814	170.678	2.566	0.010	106.516	776.247
x7	373.8717	166.515	2.210	0.000	164.083	582.000
x8	13.5550	37.881	0.366	0.715	-59.200	86.310
x9	-12.8931	39.026	-0.330	0.741	-89.464	63.677
x10	-150.5418	437.842	-0.344	0.731	-1088.084	706.901
x11	-17.8593	201.884	-0.088	0.946	-531.533	495.815
x12	-219.9402	441.262	-0.498	0.618	-1085.722	645.842
x13	-78.1493	37.611	-1.865	0.062	-143.943	3.645
Omnibus:	167.538	Durbin-Watson:			2.073	
Prob(Omnibus):	0.000	Jarque-Bera (JB):			2236.151	
Skew:	-0.001	Prob(JB):			0.000	
Kurtosis:	9.837	Cond. No.			34.4	

```
=====
```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

R2 score on test dataset: 0.8132778753638919

Part 3 Supervised Learning (40 pts)

After completing your analysis, you're satisfied with the results. You handed the Jupyter notebook over to your mentor.

(Fun fact: The name "Jupyter" is derived from Julia, Python, and R—three programming languages that the platform was originally designed for.)

Your mentor Tim is very impressed with your work and asks you the following question:

"This looks great! It will be very useful for our sales team. While looking at the results, I realized that there might be one thing that we can improve. For companies like us, it is important to sell the cars quickly. If we are patient, we might be able to sell the car for a higher price, but that's not always the best strategy. We need to consider the maintenance costs for the car, the cash flow and the fact that the price of the car decreases over time."

He then continues: "Three months is a sweet spot for us. If we can sell the car within the first three months, it is great. If not, it is worth considering lowering the price to sell it faster and increase our cash flow. I can ask Ivan from Sales to collect data in the last few months on whether the car was sold within the first three months or not. This would be great if you could have a model that tells us if the car will be sold in the first three months or not."

This sparks your interest, and soon Ivan has provided you with the new data containing an additional column `sold_within_3_months` which is a binary variable indicating whether the car was sold within the first three months or not.

Note: The data for this part is in the file `Task3_ToyotaCorolla_sales_3months.csv` and it has already unified the currency and distance units.

In [24]: data_df = pd.read_csv('data/Task3_ToyotaCorolla_sales_3months.csv', index_col=0)

```
print(data_df.head())
```

	Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight	\
0	13500	23	46986	Diesel	90	1	0	2000	3	1165	
1	13750	23	72937	Diesel	90	1	0	2000	3	1165	
2	13950	24	41711	Diesel	90	1	0	2000	3	1165	
3	14950	26	49000	Diesel	90	0	0	2000	3	1165	
4	13750	30	38500	Diesel	90	0	0	2000	3	1170	

```
data_df['sold_within_3_months']
```

```
0      0
```

```
1      0
```

```
2      0
```

```
3      0
```

```
4      0
```

3.1 (2 pts): Preprocess the Data

- (1 pts) How many cars in the dataset were sold in the first three months, and how many were not?

In [25]: print(data_df['sold_within_3_months'].value_counts())

```
print(data_df.head())
```

```
0      888
```

```
1      556
```

Name: count, dtype: int64

- (1 pts) Preprocess the categorical variables to one-hot encoding using the `pd.get_dummies()` function.

In [26]: # 7000

```
data_df = pd.get_dummies(data_df, columns=['FuelType', 'MetColor', 'Automatic', 'Doors'], drop_first=True)
```

```
# comment:
```

```
# deduct 1 point if not using drop_first=True
```

3.2 (20 pts): Logistic Regression Model

1. (2 pts) Split the data into features (X) and target (y) variables. The target variable is the 'sold_within_3_months' column. The `Price` column should be included as a feature.

In [27]: y_df = data_df['sold_within_3_months']

```
X_df = data_df.drop(columns=['sold_within_3_months'])
```

2. (2 pts) Then split the data into train test sets using a 80-20 split. Use `random_state=42` for reproducibility.

In [28]: X_train, X_test, y_train, y_test = train_test_split(X_df, y_df, test_size=0.2, random_state=42)

```
X_train, X_test, y_train, y_test = train_test_split(X_train, y_train, test_size=0.2, random_state=42)
```

```
print(f"No. of training samples: {len(X_train)} No. of testing samples: {len(X_test)}")
```

No. of training samples: 1148 No. of testing samples: 288

3. (2 pts) Standardize the features using `StandardScaler` from `sklearn.preprocessing` and then add a constant column using `sm.add_constant()`.

In [29]: #

```
scaler = StandardScaler()
```

```
X_train = scaler.fit_transform(X_train)
```

```
X_test = scaler.transform(X_test)
```

```
X_train = sm.add_constant(X_train)
```

```
X_test = sm.add_constant(X_test)
```

4. (2 pts) Fit a logistic regression model on the training dataset.

In [30]: logistic_regression_model = sm.Logit(y_train, X_train).fit()

```
print(logistic_regression_model.summary())
```

Warning: Maximum number of iterations has been exceeded.

Current function value: 0.165914

Iterations: 35

```
=====
```

```
Logit Regression Results
```

	coef	std err	z	P> z	[0.025	0.975]
const	-0.3529	19.815	-0.019	0.985	-37.621	36.915
x1	-0.9849	0.882	-1.1325	0.000	-11.713	-0.257
x2	-0.6112	0.319	-1.918	0.055	-1.236	0.013
x3	-0.2687	0.211	-1.235	0.217	-0.674	0.153
x4	-1.5443	0.654	-2.360	0.018	-2.827	-0.262
x5	1.5921	0.683	2.329	0.020	0.252	2.932
x6	0.1231	0.541	0.228	0.828	-0.937	1.184
x7	-2.9212	1.026	-1.970	0.049	-4.033	-0.010
x8	0.8738	0.589	0.125	0.900	-1.081	1.229
x9	0.1427	0.129	1.105	0.269	-0.110	0.396
x10	-0.1535	0.135	-1.134	0.257	-0.419	0.112
x11	-6.2822	5483.785	-0.001	0.999	-1.06e+04	1.06e+04
x12	-3.7854	3212.435	-0.001	0.999	-6299.962	6292.551
x13	-6.3895	5448.479	-0.001	0.999	-1.07e+04	1.07e+04

```
=====
```

Possibly complete quasi-separation: A fraction 0.20 of observations can be perfectly predicted. This might indicate that there is complete quasi-separation. In this case some parameters will not be identified.

/Users/saibo/miniconda3/envs/ada_hwl_dryrun3.10/site-packages/statsmodels/base/model.py:687: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle_retvals

warnings.warn("Maximum Likelihood optimization failed to converge. Check mle_retvals",

5. (2 pts) Evaluate the model on the test dataset using the accuracy score metric. Report the accuracy score.

In [31]: y_pred_p_logistic = logistic_regression_model.predict(X_test)

```
y_test_pred_logistic = np.where(y_pred_p_logistic > 0.5, 1, 0)
```

```
print(accuracy_fn(y_test, y_test_pred_logistic))
```

```
# Accuracy between 0.91-0.95 gives full points
```

0.9385555555555556

6. (2 pts) Calculate the precision, recall, and F1-score.

In [32]: precision = precision_fn(y_test, y_test_pred_logistic)

```
recall = recall_fn(y_test, y_test_pred_logistic)
```

```
f1 = f1_score_fn(y_test, y_test_pred_logistic)
```

```
print("Precision:", precision)
```

```
print("Recall:", recall)
```

```
print("F1-Score:", f1)
```

```
# comment:
```

```
# Any difference within 1% is considered correct. Otherwise, 1 point should be deducted for each me
```

Precision: 0.9375

Recall: 0.9482758620688655

F1-Score: 0.9428571428571428

7. (2 pts) Suppose that your company is running short on cash flow and needs to sell the cars quickly. How should you adjust the threshold for the logistic regression model to ensure that the company can sell the cars as quickly as possible?

- A. Increase the threshold
- B. Decrease the threshold

In a more general sense, how does the choice of threshold affect the precision and recall of the model?

Your Response:

- Increase the threshold will make the model predict fewer positive cases (cars will be sold within 3 months). For all the positive cases that the model predicts, the sales team will apply a discount to sell the cars quickly. So we should decrease the threshold to ensure that the company can sell the cars as quickly as possible.

However, this question has been asked in a confusing way. The student may answer the question in a different way. So we give full points for both increasing and decreasing the threshold as long as the explanation is reasonable.

8. (6 pts) Use binary search to find the optimal threshold that maximizes the F1-score. Implement a binary search algorithm to find the threshold that maximizes the f1-score of the logistic regression model on the training set. The search interval should be between 0 and 1, and the stopping criterion is 10 iterations. What is the optimal threshold and what difference does the optimal threshold make in the F1-score?

In [33]: def get_f1_score(y_test, y_pred_p, threshold):

```
y_pred = np.where(y_pred_p > threshold, 1, 0)
```

```
return f1_score_fn(y_test, y_pred)
```

```
def binary_search_threshold(y_test, y_pred_p, low, high, iterations):
```

```
for i in range(iterations):
```

```
mid = (low + high) / 2
```

```
f1_score_mid = get_f1_score(y_test, y_pred_p, mid)
```

```
f1_score_high = get_f1_score(y_test, y_pred_p, high)
```

```
if f1_score_mid > f1_score_high:
```

```
high = mid
```

```
else:
```

```
low = mid
```

```
return mid
```

```
threshold = binary_search_threshold(y_test, y_pred_p_logistic, 0, 1, 10)
```

```
print("Threshold:", threshold)
```

```
y_test_pred_logistic = np.where(y_pred_p_logistic > threshold, 1, 0)
```

```
optimal_f1 = f1_score_fn(y_test, y_test_pred_logistic)
```

```
print("Optimal F1-Score:", optimal_f1)
```

```
# comment:
```

```
# This question is flawed. Any attempt to solve this question should be given full points.
```

Threshold: 0.4677734375

Optimal F1-Score: 0.90032467532467539417

3.3 (18 pts) Decision Tree Model

Use a Decision Tree model from `Sklearn` to predict whether a car will be sold within the first three months.

Follow these steps to complete the task:

1. (2 pts) Train a Decision Tree Classifier to predict the target variable (`sold_within_3_months`). You can reuse the train and test sets from the previous step. Set `random_state=42` for reproducibility in `DecisionTreeClassifier`.

In [34]: **from** sklearn.tree **import** DecisionTreeClassifier

```
decision_tree_classifier = DecisionTreeClassifier(random_state=42)
```

```
decision_tree_classifier.fit(X_train, y_train)
```

Out [34]:

```
DecisionTreeClassifier(random_state=42)
```

2. (2 pts) Evaluate the model on the test set and report the depth of the tree.

In [35]: # This recall looks not very promising, you wonder if you had a bug in your case.

```
# you decide to calculate a baseline which all predictions are randomly 0 or 1 with the same probab
```

```
y_test_pred_rf = decision_tree_classifier.predict(X_test)
```

```
acc = accuracy_fn(y_test, y_test_pred_rf)
```

```
depth = decision_tree_classifier.get_depth()
```

```
# comments:
```

```
# Any difference within 1% is considered correct. Otherwise, 1 point should be deducted for each me
```

```
print(f"Accuracy: {acc}, Depth: {depth}")
```

```
# In difference is acceptable, otherwise 1 point should be deducted for each metric, max 2 points c
```

```
# depth should be 16, if not, check what could be wrong, if issue comes from incorrect applicatio
```

Accuracy: 0.9201388888888888, Depth: 16

3. (2 pts) Visualize the Decision Tree

In [36]: # Visualize the decision tree

```
from sklearn.tree import plot_tree
```

```
plt.figure(figsize=(20, 10))
```

```
plot_tree(decision_tree_classifier, filled=True)
```

```
plt.show()
```

4. (2 pts) Retrain the Decision Tree Classifier with a maximum depth of 8 and evaluate it on the test set. Compare and explain the results.

In [37]: # train a decision tree with max_depth=8

```
decision_tree_classifier_depth_5 = DecisionTreeClassifier(max_depth=8, random_state=42)
```

```
decision_tree_classifier_depth_5.fit(X_train, y_train)
```

```
y_test_pred_rf_depth_5 = decision_tree_classifier_depth_5.predict(X_test)
```

```
print(accuracy_fn(y_test, y_test_pred_rf_depth_5))
```

```
# comments:
```

```
# the tree with max_depth=8 should have a better performance than the one with the default depth=16
```

```
# this is because the tree with max_depth=8 is less likely to overfit the training data.
```

0.9236111111111112

5. (6 pts) Train a Decision Tree Classifier for each depth from 1 to D where D is the maximum depth of the Decision Tree Classifier seen in the previous step. Evaluate each model on the test set and plot the accuracy of the models as a function of the depth and find the optimal depth.

In [38]: #

```
depths = []
```

```
for d in depths:
```

```
decision_tree_classifier = DecisionTreeClassifier(max_depth=d, random_state=42)
```

```
decision_tree_classifier.fit(X_train, y_train)
```

```
y_test_pred_rf = decision_tree_classifier.predict(X_test)
```

```
acc = accuracy_fn(y_test, y_test_pred_rf)
```

```
accuracies.append(acc)
```

```
plt.plot(depths, accuracies)
```

```
plt.xlabel('Depth')
```

```
plt.ylabel('Accuracy')
```

```
plt.title('Accuracy vs Depth')
```

```
plt.show()
```

```
# The optimal depth is 1, which is very counter-intuitive. You decide to investigate further by plo
```

Accuracy vs Depth



6. (4 pts) Train a decision tree of depth = 1, visualize the tree and explain what is the decision rule at the root node.

In [39]: decision_tree_classifier_depth_1 = DecisionTreeClassifier(max_depth=1, random_state=42)

```
decision_tree_classifier_depth_1.fit(X_train, y_train)
```

```
plt.figure(figsize=(20, 10))
```

```
plot_tree(decision_tree_classifier_depth_1, filled=True)
```

```
plt.show()
```

```
# comment: Price is used as the first split of the tree, if it is lower than -0.036, the car is cl
```

x[1] <= -0.036
gini = 0.474
samples = 1148
value = [442, 706]

gini = 0.115
samples = 717
value = [44, 673]

gini = 0.141
samples = 431
value = [398, 33]

Part 4 Propensity Score Matching (10 pts)

Your mentor is thrilled with the progress, and he has asked Ivan to put the model into production. Based on the model's prediction, the sales manager Ivan will decide whether to lower the car's price by 5%.

A new quarter has passed, and Ivan has collected updated sales data, which includes the following columns:

- `Price`: The initial price of the car.
- `Pred_Prob`: The predicted probability of the car being sold within the first three months.
- `Discounted_Price`: Whether the discount was applied (Yes=1, No=0).
- `Applied_Discount`: The car's final price, calculated as `Price * 0.95` if the discount was applied; otherwise, it's equal to `Price`.
- `Sold_within_3_months`: Whether the car was sold within the first three months (Yes=1, No=0).

Your task is to estimate the causal effect of the discount on sales within the first three months using propensity score matching.

In [40]: data_df = pd.read_csv('data/Task4_ToyotaCorolla_discount_sales.csv', index_col=0)

```
print(data_df.head())
```

	Price	Pred_Prob	Applied_Discount	Discounted_Price	Sold_within_3_months
0	12750	0.01	1	11475.0	1
1	21950	0.00	1	19755.0	1
2	9550	0.79	0	9550.0	1
3	9530	0.91	1	8957.0	0
4	9450	0.97	0	9450.0	0

- 4.1 (1 pts): How many samples are in the treated group, and how many are in the control group?

In [41]: print(data_df['Applied_Discount'].value_counts())

```
# 118 are in the treated group and 82 are in the control group
```

```
Applied_Discount
```

```
0      118
```

```
1       82
```

Name: count, dtype: int64

4.2 (5 pts): Propensity Score Matching

- The propensity score is the predicted probability of the car being sold within the first three months from the logistic regression model. I.e. `Pred_Prob` column in the `Task4_ToyotaCorolla_discount_sales.csv` file. Create pairs of matched samples as follows:
 - For each treated sample (discount applied), find a control sample (discount not applied) with a difference in propensity score of less than 0.05.
 - If there is more than one control sample for a treated sample, choose the control sample with the smallest difference in propensity score.
 - If there is no control sample satisfying the condition, discard the treated sample.

Notice that your output should be 1-to-1 matching, meaning that each treated sample must be matched with at most one control sample. So each sample can only appear once in the pairs.

In [42]: # 7000

```
treatment_df = data_df[data_df['Applied_Discount'] == 1]
```

```
control_df = data_df[data_df['Applied_Discount'] == 0]
```

```
print(f"There are {treatment_df.shape[0]} samples in the treated group.")
```

```
print(f"There are {control_df.shape[0]} samples in the control group.")
```

```
import networkx as nx
```

```
def get_similarity(propensity_score1, propensity_score2):
```

```
    Calculate similarity for instances with given propensity scores'''
```

```
    return 1 - np.abs(propensity_score1 - propensity_score2)
```

```
G = nx.Graph()
```

```
for treatment_id, treatment_row in treatment_df.iterrows():
```

```
    for control_id, control_row in control_df.iterrows():
```

```
        similarity = get_similarity(control_row['Pred_Prob'], treatment_row['Pred_Prob'])
```

```
        if similarity > 0.95:
```

```
            G.add_weighted_edges_from([(control_id, treatment_id, similarity)])
```

```
matching = nx.max_weight_matching(G, maxcardinality=True)
```

```
print(f"We have {len(matching)} successful matches!")
```

```
# comment:
```

```
# The networkx library is the correct library to use for this task and it was used in the lab sessi
```

```
# Technically, the student can use other libraries to solve this task, but most often they use a gr
```

```
# We don't deduct points for the argument "maxcardinality=True" in the nx.max_weight_matching func
```

There are 82 samples in the treated group.

There are 118 samples in the control group.

We have 49 successful matches!

4.3 (4 pts): Average Treatment Effect (ATE)

Now let's estimate the effect of the discount on sales.

For each matched pair, there is one treated sample and one control sample. They may have different outcomes and we can calculate the average treatment effect (ATE) as