

Chapter 2: End-to-End Machine Learning Project

Summary from “Hands-On Machine Learning”

Study Notes

June 29, 2025

Contents

1 Chapter Overview	2
2 The 8-Step ML Project Workflow	2
3 Step 1: Frame the Problem	2
3.1 Problem Definition	2
3.2 Performance Measures	3
4 Step 2: Get the Data	3
4.1 Data Loading	3
4.2 Train/Test Split	3
5 Step 3: Explore the Data	4
5.1 Geographic Visualization	4
5.2 Correlation Analysis	4
6 Step 4: Prepare the Data	4
6.1 Handle Missing Values	4
6.2 Feature Scaling	4
6.3 Feature Engineering	5
7 Step 5: Select and Train Models	5
7.1 Model Comparison	5
7.2 Cross-Validation	6
8 Step 6: Fine-Tune the Model	6
8.1 Grid Search	6
8.2 Feature Importance	7
9 Step 7: Evaluate on Test Set	7
10 Complete Pipeline	7
11 Key Takeaways	8
12 Common Pitfalls	8
13 Conclusion	8

1 Chapter Overview

This chapter demonstrates the complete workflow of a machine learning project using the California housing dataset. It covers all essential steps from problem framing to model deployment.

Key Learning Objectives:

- Understand the complete ML project workflow
- Learn data exploration and visualization techniques
- Master data preprocessing and feature engineering
- Compare multiple ML algorithms systematically
- Implement proper model evaluation and validation

2 The 8-Step ML Project Workflow

1. **Frame the Problem** - Define objectives and constraints
2. **Get the Data** - Collect and load the dataset
3. **Explore the Data** - Visualize and understand patterns
4. **Prepare the Data** - Clean and preprocess
5. **Select Models** - Try different algorithms
6. **Fine-tune** - Optimize hyperparameters
7. **Present Solution** - Communicate results
8. **Deploy & Monitor** - Production deployment

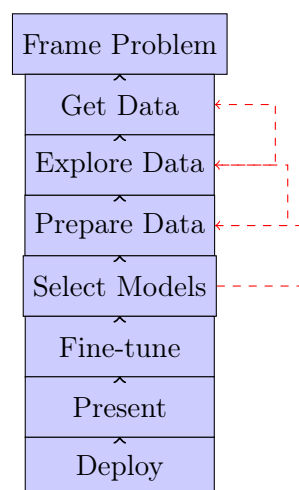


Figure 1: End-to-End Machine Learning Project Workflow

3 Step 1: Frame the Problem

3.1 Problem Definition

The California housing project predicts median house values using features like population, median income, and location.

Key Questions:

- What is the business objective?
- How will the solution be used?
- Supervised, unsupervised, or reinforcement learning?
- Classification or regression?
- How to measure performance?

3.2 Performance Measures

For regression problems:

$$\text{RMSE} = \sqrt{\frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2} \quad (1)$$

$$\text{MAE} = \frac{1}{m} \sum_{i=1}^m |h(x^{(i)}) - y^{(i)}| \quad (2)$$

4 Step 2: Get the Data

4.1 Data Loading

Listing 1: Loading Dataset

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # Load the data
6 housing = pd.read_csv('housing.csv')
7
8 # Quick exploration
9 print(housing.head())
10 print(housing.info())
11 print(housing.describe())
```

4.2 Train/Test Split

Listing 2: Data Splitting

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.model_selection import StratifiedShuffleSplit
3
4 # Simple random split
5 train_set, test_set = train_test_split(
6     housing, test_size=0.2, random_state=42)
7
8 # Stratified split (better for small datasets)
9 split = StratifiedShuffleSplit(n_splits=1, test_size=0.2, random_state
10                                =42)
11 for train_index, test_index in split.split(housing, housing["income_cat
12                                             "]):
13     strat_train_set = housing.loc[train_index]
14     strat_test_set = housing.loc[test_index]
```

5 Step 3: Explore the Data

5.1 Geographic Visualization

California housing data shows clear geographic patterns with coastal areas having higher prices.

5.2 Correlation Analysis

Listing 3: Correlation Analysis

```

1 # Compute correlation matrix
2 corr_matrix = housing.corr()
3
4 # Focus on target variable
5 print(corr_matrix["median_house_value"].sort_values(ascending=False))
6
7 # Scatter matrix for key attributes
8 from pandas.plotting import scatter_matrix
9 attributes = ["median_house_value", "median_income",
10              "total_rooms", "housing_median_age"]
11 scatter_matrix(housing[attributes], figsize=(12, 8))

```

6 Step 4: Prepare the Data

6.1 Handle Missing Values

Listing 4: Missing Data Strategies

```

1 # Three options:
2 # 1. Drop rows: housing.dropna(subset=["total_bedrooms"])
3 # 2. Drop column: housing.drop("total_bedrooms", axis=1)
4 # 3. Fill values: housing["total_bedrooms"].fillna(median)
5
6 # Using sklearn
7 from sklearn.impute import SimpleImputer
8 imputer = SimpleImputer(strategy="median")
9 housing_num = housing.drop("ocean_proximity", axis=1)
10 imputer.fit(housing_num)
11 X = imputer.transform(housing_num)

```

6.2 Feature Scaling

Two main approaches:

$$\text{Min-Max: } x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (3)$$

$$\text{Standardization: } x_{scaled} = \frac{x - \mu}{\sigma} \quad (4)$$

Listing 5: Feature Scaling

```

1 from sklearn.preprocessing import StandardScaler, MinMaxScaler
2
3 # Standardization
4 std_scaler = StandardScaler()
5 housing_scaled = std_scaler.fit_transform(housing_num)

```

```

6
7 # Min-Max scaling
8 minmax_scaler = MinMaxScaler()
9 housing_scaled = minmax_scaler.fit_transform(housing_num)

```

6.3 Feature Engineering

Listing 6: Creating New Features

```

1 # Combination attributes
2 housing["rooms_per_household"] = housing["total_rooms"]/housing["
   households"]
3 housing["bedrooms_per_room"] = housing["total_bedrooms"]/housing["
   total_rooms"]
4 housing["population_per_household"] = housing["population"]/housing["
   households"]
5
6 # Custom transformer
7 from sklearn.base import BaseEstimator, TransformerMixin
8
9 class CombinedAttributesAdder(BaseEstimator, TransformerMixin):
10     def __init__(self, add_bedrooms_per_room=True):
11         self.add_bedrooms_per_room = add_bedrooms_per_room
12
13     def fit(self, X, y=None):
14         return self
15
16     def transform(self, X):
17         rooms_per_household = X[:, rooms_ix] / X[:, households_ix]
18         population_per_household = X[:, population_ix] / X[:,
            households_ix]
19         if self.add_bedrooms_per_room:
20             bedrooms_per_room = X[:, bedrooms_ix] / X[:, rooms_ix]
21             return np.c_[X, rooms_per_household,
                population_per_household, bedrooms_per_room]
22         else:
23             return np.c_[X, rooms_per_household,
                population_per_household]

```

7 Step 5: Select and Train Models

7.1 Model Comparison

Table 1: Algorithm Performance Comparison

Model	RMSE	Speed
Linear Regression	68,628	Fast
Decision Tree	0 (overfitting)	Medium
Random Forest	49,682	Slow
SVM	111,094	Very Slow

Listing 7: Training Multiple Models

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.tree import DecisionTreeRegressor
3 from sklearn.ensemble import RandomForestRegressor
4 from sklearn.metrics import mean_squared_error
5
6 # Linear Regression
7 lin_reg = LinearRegression()
8 lin_reg.fit(housing_prepared, housing_labels)
9 predictions = lin_reg.predict(housing_prepared)
10 lin_mse = mean_squared_error(housing_labels, predictions)
11 lin_rmse = np.sqrt(lin_mse)
12
13 # Random Forest
14 forest_reg = RandomForestRegressor()
15 forest_reg.fit(housing_prepared, housing_labels)
```

7.2 Cross-Validation

Listing 8: K-Fold Cross-Validation

```
1 from sklearn.model_selection import cross_val_score
2
3 def display_scores(scores):
4     print("Scores:", scores)
5     print("Mean:", scores.mean())
6     print("Standard deviation:", scores.std())
7
8 # 10-fold cross-validation
9 forest_scores = cross_val_score(forest_reg, housing_prepared,
10                                housing_labels,
11                                scoring="neg_mean_squared_error", cv=10)
12 forest_rmse_scores = np.sqrt(-forest_scores)
13 display_scores(forest_rmse_scores)
```

8 Step 6: Fine-Tune the Model

8.1 Grid Search

Listing 9: Hyperparameter Tuning

```
1 from sklearn.model_selection import GridSearchCV
2
3 param_grid = [
4     {'n_estimators': [3, 10, 30], 'max_features': [2, 4, 6, 8]},
5     {'bootstrap': [False], 'n_estimators': [3, 10], 'max_features': [2,
6     3, 4]},
7 ]
8
9 forest_reg = RandomForestRegressor()
10 grid_search = GridSearchCV(forest_reg, param_grid, cv=5,
11                             scoring='neg_mean_squared_error')
12 grid_search.fit(housing_prepared, housing_labels)
13
14 print("Best parameters:", grid_search.best_params_)
```

8.2 Feature Importance

Random Forest can tell us which features are most important:

Listing 10: Feature Importance Analysis

```

1 feature_importances = grid_search.best_estimator_.feature_importances_
2 extra_attribs = ["rooms_per_hhold", "pop_per_hhold", "bedrooms_per_room"]
3
4 cat_encoder = full_pipeline.named_transformers_["cat"]
5 cat_one_hot_attribs = list(cat_encoder.categories_[0])
6 attributes = num_attribs + extra_attribs + cat_one_hot_attribs
7 print(sorted(zip(feature_importances, attributes), reverse=True))

```

Top features typically include:

1. Median income (0.334)
2. Ocean proximity (0.364 combined)
3. Longitude/Latitude (0.088/0.069)
4. Rooms per household (custom feature)

9 Step 7: Evaluate on Test Set

Listing 11: Final Evaluation

```

1 # Use best model from grid search
2 final_model = grid_search.best_estimator_
3
4 # Prepare test data
5 X_test = strat_test_set.drop("median_house_value", axis=1)
6 y_test = strat_test_set["median_house_value"].copy()
7 X_test_prepared = full_pipeline.transform(X_test)
8
9 # Final predictions
10 final_predictions = final_model.predict(X_test_prepared)
11 final_mse = mean_squared_error(y_test, final_predictions)
12 final_rmse = np.sqrt(final_mse)
13
14 print(f"Final RMSE: {final_rmse}")
15
16 # Confidence interval
17 from scipy import stats
18 squared_errors = (final_predictions - y_test) ** 2
19 confidence_interval = np.sqrt(stats.t.interval(
20     0.95, len(squared_errors) - 1,
21     loc=squared_errors.mean(),
22     scale=stats.sem(squared_errors)))

```

10 Complete Pipeline

Listing 12: Full Processing Pipeline

```

1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer

```

```
3 from sklearn.preprocessing import OneHotEncoder
4
5 # Numerical pipeline
6 num_pipeline = Pipeline([
7     ('imputer', SimpleImputer(strategy="median")),
8     ('attrs_adder', CombinedAttributesAdder()),
9     ('std_scaler', StandardScaler()),
10 ])
11
12 # Complete pipeline
13 num_attribs = list(housing_num)
14 cat_attribs = ["ocean_proximity"]
15
16 full_pipeline = ColumnTransformer([
17     ("num", num_pipeline, num_attribs),
18     ("cat", OneHotEncoder(), cat_attribs),
19 ])
20
21 # Transform data
22 housing_prepared = full_pipeline.fit_transform(housing)
```

11 Key Takeaways

Best Practices:

1. Create test set before exploring data
2. Use stratified sampling for small datasets
3. Explore data with visualizations
4. Handle missing values appropriately
5. Scale features for distance-based algorithms
6. Try multiple algorithms
7. Use cross-validation for evaluation
8. Fine-tune hyperparameters systematically
9. Analyze feature importance
10. Evaluate on test set only once

12 Common Pitfalls

Avoid These Mistakes:

- Looking at test data early
- Data snooping during tuning
- Ignoring data leakage
- Not handling categorical variables
- Forgetting feature scaling
- Overfitting to validation set
- Ignoring business context

13 Conclusion

Chapter 2 provides a systematic framework for ML projects. The key insight is that machine learning involves much more than just algorithms—proper data handling, systematic evaluation,

and following a structured methodology are crucial for success.

The California housing example demonstrates how to:

- Frame problems correctly
- Handle real-world data issues
- Compare algorithms fairly
- Tune models systematically
- Evaluate performance rigorously

This workflow can be adapted to various regression and classification problems across different domains.