

- [Part II: Intermediate Level - Statistical Foundations](#)
 - [Chapter 6: Statistical Concepts for Machine Learning](#)
 - [Chapter 7: Advanced Evaluation Metrics and Techniques](#)

Part II: Intermediate Level - Statistical Foundations

Chapter 6: Statistical Concepts for Machine Learning

User: I feel like I have a good grasp of the basics now, but I keep hearing about statistical concepts that seem important for machine learning. Things like confidence intervals, correlation coefficients, and statistical tests. Can you help me understand how these connect to what we've learned?

Expert: Absolutely! You're ready to dive into the statistical foundations that make machine learning more rigorous and reliable. These concepts will help you better understand your data, validate your models, and make more confident decisions. Let's start with the fundamentals.

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy import stats
6 from sklearn.datasets import make_regression
7 from sklearn.linear_model import LinearRegression
8 from sklearn.metrics import mean_squared_error, r2_score
9 import warnings
10 warnings.filterwarnings('ignore')
11 def statistical_foundations_intro():
12     """
13         Introduction to key statistical concepts for ML
14     """
15
16     print("STATISTICAL FOUNDATIONS FOR MACHINE LEARNING")
17     print("=" * 50)
18
19     print("Why Statistics Matter in ML:")
20     print("✓ Understand data distributions and patterns")
21     print("✓ Quantify uncertainty in predictions")
22     print("✓ Test hypotheses about model performance")
23     print("✓ Make confident business decisions")
24     print("✓ Detect and avoid statistical pitfalls")
25
26     return True
27 statistical_foundations_intro()

```

User: That makes sense! I want to understand how to be more confident about my model's predictions and avoid making wrong conclusions. Where should we start?

Expert: Perfect mindset! Let's start with **Pearson's Correlation Coefficient** - one of the most fundamental statistical measures for understanding relationships between variables.

```

1 def pearson_correlation_deep_dive():
2     """
3         Comprehensive explanation of Pearson's correlation
4     """
5     print("PEARSON'S CORRELATION COEFFICIENT")
6     print("=" * 40)
7
8     # Generate sample data with different correlation strengths
9     np.random.seed(42)
10    n_samples = 100
11
12    # Perfect positive correlation

```

```

13     x1 = np.random.normal(0, 1, n_samples)
14     y1 = 2 * x1 + np.random.normal(0, 0.1, n_samples) # Almost
15     perfect
16
17     # Moderate positive correlation
18     x2 = np.random.normal(0, 1, n_samples)
19     y2 = x2 + np.random.normal(0, 1, n_samples)
20
21     # No correlation
22     x3 = np.random.normal(0, 1, n_samples)
23     y3 = np.random.normal(0, 1, n_samples)
24
25     # Negative correlation
26     x4 = np.random.normal(0, 1, n_samples)
27     y4 = -1.5 * x4 + np.random.normal(0, 0.5, n_samples)
28
29     # Calculate Pearson correlations
30     corr1, p_value1 = stats.pearsonr(x1, y1)
31     corr2, p_value2 = stats.pearsonr(x2, y2)
32     corr3, p_value3 = stats.pearsonr(x3, y3)
33     corr4, p_value4 = stats.pearsonr(x4, y4)
34
35     print("CORRELATION COEFFICIENT INTERPRETATION:")
36     print("-" * 40)
37     print("Range: -1 to +1")
38     print("+1: Perfect positive correlation")
39     print(" 0: No linear correlation")
40     print("-1: Perfect negative correlation")
41     print()
42
43     print("STRENGTH GUIDELINES:")
44     print("0.0 - 0.3: Weak correlation")
45     print("0.3 - 0.7: Moderate correlation")
46     print("0.7 - 1.0: Strong correlation")
47     print()
48
49     # Visualize different correlations
50     fig, axes = plt.subplots(2, 2, figsize=(12, 10))
51
52     # Strong positive
53     axes[0, 0].scatter(x1, y1, alpha=0.6, color='red')
54     axes[0, 0].set_title(f'Strong Positive\nnr = {corr1:.3f}, p = {p_value1:.3f}')
55     axes[0, 0].set_xlabel('X')
56     axes[0, 0].set_ylabel('Y')
57     axes[0, 0].grid(True, alpha=0.3)
58
59     # Moderate positive
60     axes[0, 1].scatter(x2, y2, alpha=0.6, color='blue')
61     axes[0, 1].set_title(f'Moderate Positive\nnr = {corr2:.3f}, p = {p_value2:.3f}')
62     axes[0, 1].set_xlabel('X')
63     axes[0, 1].set_ylabel('Y')
64     axes[0, 1].grid(True, alpha=0.3)
65
66     # No correlation
67     axes[1, 0].scatter(x3, y3, alpha=0.6, color='green')
68     axes[1, 0].set_title(f'No Correlation\nnr = {corr3:.3f}, p = {p_value3:.3f}')
69     axes[1, 0].set_xlabel('X')
70     axes[1, 0].set_ylabel('Y')
71     axes[1, 0].grid(True, alpha=0.3)
72
73     # Negative correlation
74     axes[1, 1].scatter(x4, y4, alpha=0.6, color='purple')
75     axes[1, 1].set_title(f'Negative Correlation\nnr = {corr4:.3f}, p = {p_value4:.3f}')
76     axes[1, 1].set_xlabel('X')
77     axes[1, 1].set_ylabel('Y')
78     axes[1, 1].grid(True, alpha=0.3)

```

```

{p_value2:.3f}'")
61     axes[0, 1].set_xlabel('X')
62     axes[0, 1].set_ylabel('Y')
63     axes[0, 1].grid(True, alpha=0.3)
64
65     # No correlation
66     axes[1, 0].scatter(x3, y3, alpha=0.6, color='green')
67     axes[1, 0].set_title(f'No Correlation\nnr = {corr3:.3f}, p =
{p_value3:.3f}')
68     axes[1, 0].set_xlabel('X')
69     axes[1, 0].set_ylabel('Y')
70     axes[1, 0].grid(True, alpha=0.3)
71
72     # Negative correlation
73     axes[1, 1].scatter(x4, y4, alpha=0.6, color='purple')
74     axes[1, 1].set_title(f'Strong Negative\nnr = {corr4:.3f}, p =
{p_value4:.3f}')
75     axes[1, 1].set_xlabel('X')
76     axes[1, 1].set_ylabel('Y')
77     axes[1, 1].grid(True, alpha=0.3)
78
79     plt.tight_layout()
80     plt.show()
81
82     print("BUSINESS INTERPRETATION:")
83     print("-" * 25)
84     print(f"Strong Positive (r={corr1:.3f}): As X increases, Y
increases predictably")
85     print(f"Moderate Positive (r={corr2:.3f}): X and Y tend to
increase together")
86     print(f"No Correlation (r={corr3:.3f}): X and Y are
independent")
87     print(f"Strong Negative (r={corr4:.3f}): As X increases, Y
decreases predictably")
88
89     print("\nP-VALUES EXPLAINED:")
90     print("-" * 20)
91     print("p < 0.05: Statistically significant correlation")
92     print("p ≥ 0.05: Could be due to random chance")
93
94     return corr1, corr2, corr3, corr4
95 correlations = pearson_correlation_deep_dive()

```

User: This is really helpful! I can see how correlation strength affects the scatter plots. But what about confidence intervals? I hear this term a lot but I'm not sure what it means practically.

Expert: Excellent question! **Confidence intervals** are crucial for understanding the uncertainty in our estimates. Let me show you how they work in the context of machine learning.

```
1 def confidence_intervals_explained():
2     """
3     Comprehensive explanation of confidence intervals
4     """
5     print("CONFIDENCE INTERVALS IN MACHINE LEARNING")
6     print("-" * 45)
7
8     print("WHAT IS A CONFIDENCE INTERVAL?")
9     print("-" * 35)
10    print("A range of values that likely contains the true
11        parameter")
12    print("95% CI means: If we repeated this study 100 times,")
13    print("95 of those intervals would contain the true value")
14    print()
15
16    # Generate sample data
17    np.random.seed(42)
18    n_samples = 50
19    true_slope = 2.5
20    true_intercept = 10
21
22    x = np.random.uniform(1, 10, n_samples)
23    noise = np.random.normal(0, 2, n_samples)
24    y = true_intercept + true_slope * x + noise
25
26    # Fit linear regression
27    from sklearn.linear_model import LinearRegression
28    model = LinearRegression()
29    X = x.reshape(-1, 1)
30    model.fit(X, y)
31
32    # Get predictions
33    y_pred = model.predict(X)
34
35    # Calculate confidence intervals manually
36    def calculate_prediction_intervals(X, y, model,
37        confidence=0.95):
38        """
39            Calculate prediction intervals for linear regression
40        """
41        n = len(y)
42        y_pred = model.predict(X)
43
44        # Calculate residuals and standard error
45        residuals = y - y_pred
46        mse = np.mean(residuals**2)
47        se = np.sqrt(mse)
```

```

47      # Degrees of freedom
48      dof = n - 2 # n - (number of parameters)
49
50      # t-critical value
51      alpha = 1 - confidence
52      t_crit = stats.t.ppf(1 - alpha/2, dof)
53
54      # Prediction intervals
55      margin_error = t_crit * se * np.sqrt(1 + 1/n + (X.flatten()
56      - np.mean(X.flatten()))**2 / np.sum((X.flatten() -
57      np.mean(X.flatten()))**2))
58
59      lower_bound = y_pred - margin_error
60      upper_bound = y_pred + margin_error
61
62      return lower_bound, upper_bound, se, t_crit
63
64      print("CONFIDENCE INTERVAL CALCULATION:")
65      print("-" * 35)
66      print(f"Sample size: {n_samples}")
67      print(f"Degrees of freedom: {n_samples - 2}")
68      print(f"Standard error: {se:.3f}")
69      print(f"t-critical value (95% CI): {t_crit:.3f}")
70      print(f"True slope: {true_slope:.2f}")
71      print(f"Estimated slope: {model.coef_[0]:.2f}")
72      print(f"True intercept: {true_intercept:.2f}")
73      print(f"Estimated intercept: {model.intercept_:.2f}")
74
75      # Visualize confidence intervals
76      plt.figure(figsize=(12, 8))
77
78      # Main plot with confidence intervals
79      plt.subplot(2, 2, 1)
80      plt.scatter(x, y, alpha=0.6, color='blue', label='Data')
81      plt.plot(x, y_pred, 'r-', label='Regression Line')
82      plt.fill_between(x, lower_bound, upper_bound, alpha=0.3,
83          color='red', label='95% Prediction Interval')
84      plt.xlabel('X')
85      plt.ylabel('Y')
86      plt.title('Linear Regression with 95% Prediction Intervals')
87      plt.legend()
88      plt.grid(True, alpha=0.3)
89
90      # Demonstrate multiple samples
91      plt.subplot(2, 2, 2)
92      # Simulate multiple experiments

```

```

93     slopes = []
94     intercepts = []
95
96     for i in range(100):
97         # Generate new sample each time
98         x_sim = np.random.uniform(1, 10, n_samples)
99         y_sim = true_intercept + true_slope * x_sim +
100            np.random.normal(0, 2, n_samples)
101
102     # Fit model
103     model_sim = LinearRegression()
104     model_sim.fit(x_sim.reshape(-1, 1), y_sim)
105
106     slopes.append(model_sim.coef_[0])
107     intercepts.append(model_sim.intercept_)
108
109     plt.hist(slopes, bins=20, alpha=0.7, color='green')
110     plt.axvline(true_slope, color='red', linestyle='--',
111      linewidth=2, label=f'True Slope ({true_slope})')
112     plt.axvline(np.mean(slopes), color='blue', linestyle='--',
113      linewidth=2, label=f'Mean Estimate ({np.mean(slopes):.2f})')
114     plt.xlabel('Estimated Slope')
115     plt.ylabel('Frequency')
116     plt.title('Distribution of Slope Estimates (100 experiments)')
117     plt.legend()
118     plt.grid(True, alpha=0.3)
119
120     # Calculate confidence interval for slope
121     slope_ci_lower = np.percentile(slopes, 2.5)
122     slope_ci_upper = np.percentile(slopes, 97.5)
123
124     # Business interpretation
125     plt.subplot(2, 1, 2)
126     plt.axis('off')
127
128     interpretation_text = f"""
129 BUSINESS INTERPRETATION OF CONFIDENCE INTERVALS:
130 ✓ SLOPE ESTIMATE: {model.coef_[0]:.2f}
131 ✓ 95% CI for slope: [{slope_ci_lower:.2f}, {slope_ci_upper:.2f}]
132 WHAT THIS MEANS:
133 • We're 95% confident the true relationship between X and Y
134   has a slope between {slope_ci_lower:.2f} and {slope_ci_upper:.2f}
135 • For every 1-unit increase in X, Y increases by approximately
136   {model.coef_[0]:.2f} units ( $\pm$ {(slope_ci_upper-
137     slope_ci_lower)/2:.2f})

```

```

130 BUSINESS DECISIONS:
141 • If planning based on this relationship, account for uncertainty
142 • Wider intervals = more uncertainty = more conservative planning
143 • Narrower intervals = more confidence = can make bolder decisions
145 SAMPLE SIZE EFFECT:
146 • Larger samples → Narrower confidence intervals
147 • More data → More confident estimates
148     """
149
150     plt.text(0.05, 0.95, interpretation_text,
151               transform=plt.gca().transAxes,
152               fontsize=10, verticalalignment='top',
153               fontfamily='monospace',
154               bbox=dict(boxstyle='round', facecolor='lightblue',
155                     alpha=0.8))
156
157     plt.tight_layout()
158     plt.show()
159
160     print("PRACTICAL APPLICATIONS:")
161     print("-" * 25)
162     print("• Revenue forecasting: 'Sales will be $100K ± $15K'")
163     print("• A/B testing: 'Treatment effect is 5% ± 2%'")
164     print("• Model performance: 'Accuracy is 85% ± 3%'")
165     print("• Risk assessment: 'Default rate is 2% ± 0.5%'")
166
167     return slopes, intercepts, slope_ci_lower, slope_ci_upper
168     slopes, intercepts, ci_lower, ci_upper =
169     confidence_intervals_explained()

```

User: This is incredibly insightful! I love how you showed that confidence intervals help us understand the uncertainty in our estimates. Now I'm curious about statistical tests. When would I use something like a t-test in machine learning?

Expert: Great question! Statistical tests are essential for making data-driven decisions and validating our models. Let me show you the **t-test** and other important statistical tests in ML contexts.

```

1 def statistical_tests_for_ml():
2     """
3     Statistical tests commonly used in machine learning
4     """
5     print("STATISTICAL TESTS IN MACHINE LEARNING")
6     print("=" * 40)
7
8     print("WHY WE NEED STATISTICAL TESTS:")
9     print("-" * 35)

```

```
10     print("• Validate model improvements")
11     print("• Compare different algorithms")
12     print("• Test feature importance")
13     print("• Validate A/B test results")
14     print("• Detect significant differences")
15
16     # Generate sample data for demonstrations
17     np.random.seed(42)
18
19     print("\n1. ONE-SAMPLE T-TEST")
20     print("-" * 20)
21     print("Question: Is our model's accuracy significantly
different from a baseline?")
22
23     # Simulate model accuracy scores from cross-validation
24     baseline_accuracy = 0.75
25     model_accuracies = np.random.normal(0.82, 0.05, 30) # 30 CV
folds
26
27     # Perform one-sample t-test
28     t_stat, p_value = stats.ttest_1samp(model_accuracies,
baseline_accuracy)
29
30     print(f"Baseline accuracy: {baseline_accuracy:.1%}")
31     print(f"Our model accuracy: {np.mean(model_accuracies):.1%} ±
{np.std(model_accuracies):.1%}")
32     print(f"t-statistic: {t_stat:.3f}")
33     print(f"p-value: {p_value:.6f}")
34
35     if p_value < 0.05:
36         print("✅ RESULT: Significant improvement over baseline!")
37     else:
38         print("❌ RESULT: No significant improvement over
baseline")
39
40     print("\n2. TWO-SAMPLE T-TEST")
41     print("-" * 20)
42     print("Question: Is Model A significantly better than Model
B?")
43
44     # Simulate performance of two models
45     model_a_scores = np.random.normal(0.85, 0.04, 25)
46     model_b_scores = np.random.normal(0.82, 0.045, 25)
47
48     # Perform independent t-test
49     t_stat_2, p_value_2 = stats.ttest_ind(model_a_scores,
model_b_scores)
50
51     print(f"Model A accuracy: {np.mean(model_a_scores):.1%} ±
{np.std(model_a_scores):.1%}")
```

```

52     print(f"Model B accuracy: {np.mean(model_b_scores):.1%} ±
53         {np.std(model_b_scores):.1%}")
54     print(f"Difference: {np.mean(model_a_scores) -
55         np.mean(model_b_scores):.1%}")
56     print(f"t-statistic: {t_stat_2:.3f}")
57     print(f"p-value: {p_value_2:.6f}")
58
59     if p_value_2 < 0.05:
60         print("✅ RESULT: Model A is significantly better than
61             Model B!")
62     else:
63         print("❌ RESULT: No significant difference between
64             models")
65
66     print("\n3. PAIRED T-TEST")
67     print("-" * 16)
68     print("Question: Does feature engineering significantly improve
69             performance?")
70
71     # Simulate before/after feature engineering on same datasets
72     before_scores = np.random.normal(0.78, 0.06, 20)
73     after_scores = before_scores + np.random.normal(0.05, 0.02, 20)
74     # Improvement
75
76     # Perform paired t-test
77     t_stat_paired, p_value_paired = stats.ttest_rel(after_scores,
78             before_scores)
79
80     print(f"Before feature engineering:
81         {np.mean(before_scores):.1%} ± {np.std(before_scores):.1%}")
82     print(f"After feature engineering: {np.mean(after_scores):.1%}
83         ± {np.std(after_scores):.1%}")
84     print(f"Average improvement: {np.mean(after_scores) -
85             before_scores:.1%}")
86     print(f"t-statistic: {t_stat_paired:.3f}")
87     print(f"p-value: {p_value_paired:.6f}")
88
89     if p_value_paired < 0.05:
90         print("✅ RESULT: Feature engineering significantly
91             improves performance!")
92     else:
93         print("❌ RESULT: No significant improvement from feature
94             engineering")
95
96     # Visualize the tests
97     fig, axes = plt.subplots(2, 2, figsize=(15, 10))
98
99     # One-sample t-test visualization
100    axes[0, 0].hist(model_accuracies, bins=15, alpha=0.7,
101        color='blue', density=True)

```

```

89     axes[0, 0].axvline(baseline_accuracy, color='red', linestyle='--',
90                          linewidth=2, label=f'Baseline ({baseline_accuracy:.1%})')
91     axes[0, 0].axvline(np.mean(model_accuracies), color='green',
92                          linestyle='--', linewidth=2, label=f'Our Model
93                                         ({np.mean(model_accuracies):.1%})')
94
95     axes[0, 0].set_xlabel('Accuracy')
96     axes[0, 0].set_ylabel('Density')
97     axes[0, 0].set_title('One-Sample t-test: Model vs Baseline')
98     axes[0, 0].legend()
99     axes[0, 0].grid(True, alpha=0.3)
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127 EFFECT SIZES & PRACTICAL SIGNIFICANCE:
128 Cohen's d interpretation:
130 • Small effect: d = 0.2

```

```

131 • Medium effect: d = 0.5
132 • Large effect: d = 0.8
133 Model A vs B: d = {effect_size_ab:.2f}
135 Before vs After: d = {effect_size_paired:.2f}
136 REMEMBER:
138 ✓ Statistical significance ≠ Practical significance
139 ✓ Large samples can make tiny differences "significant"
140 ✓ Always consider business impact
141 ✓ Effect size matters more than p-value for decisions
142 BUSINESS DECISION FRAMEWORK:
144 1. Is the difference statistically significant? (p < 0.05)
145 2. Is the effect size meaningful? (d > 0.2)
146 3. Is the business impact worth the cost?
147 """
148
149     axes[1, 1].text(0.05, 0.95, effect_text, transform=axes[1,
1].transAxes,
150                     fontsize=9, verticalalignment='top',
151                     fontfamily='monospace',
152                     bbox=dict(boxstyle='round',
153                               facecolor='lightyellow', alpha=0.8))
154
155
156     return {
157         'one_sample': (t_stat, p_value),
158         'two_sample': (t_stat_2, p_value_2),
159         'paired': (t_stat_paired, p_value_paired),
160         'effect_sizes': (effect_size_ab, effect_size_paired)
161     }
162 test_results = statistical_tests_for_ml()

```

User: This is fantastic! I can see how statistical tests help us make confident decisions about model improvements. But I'm also curious about the error metrics we use. I know about MSE and RMSE, but can you explain the differences and when to use each one?

Expert: Excellent question! Understanding different error metrics is crucial for choosing the right evaluation approach. Let me show you **MSE**, **RMSE**, and **MAE** with their practical implications.

```

1 def error_metrics_comprehensive():
2     """
3     Comprehensive explanation of regression error metrics
4     """
5     print("REGRESSION ERROR METRICS: MSE, RMSE, MAE")
6     print("=" * 45)
7

```

```

8     # Generate sample data with different error patterns
9     np.random.seed(42)
10    n_samples = 100
11
12    # True values
13    y_true = np.random.uniform(10, 100, n_samples)
14
15    # Different prediction scenarios
16    # Scenario 1: Good predictions with small errors
17    y_pred_good = y_true + np.random.normal(0, 5, n_samples)
18
19    # Scenario 2: Predictions with some outliers
20    y_pred_outliers = y_true + np.random.normal(0, 3, n_samples)
21    outlier_indices = np.random.choice(n_samples, 5, replace=False)
22    y_pred_outliers[outlier_indices] += np.random.choice([-30, 30],
23                                            5)
23
24    # Scenario 3: Consistently biased predictions
25    y_pred_biased = y_true * 1.2 + np.random.normal(0, 2,
26                                                    n_samples)
26
27    def calculate_all_metrics(y_true, y_pred):
28        """Calculate all error metrics"""
29        mse = np.mean((y_true - y_pred) ** 2)
30        rmse = np.sqrt(mse)
31        mae = np.mean(np.abs(y_true - y_pred))
32
33        # Additional metrics
34        mape = np.mean(np.abs((y_true - y_pred) / y_true)) * 100
35        max_error = np.max(np.abs(y_true - y_pred))
36
37        return {
38            'MSE': mse,
39            'RMSE': rmse,
40            'MAE': mae,
41            'MAPE': mape,
42            'Max_Error': max_error
43        }
44
45    # Calculate metrics for all scenarios
46    metrics_good = calculate_all_metrics(y_true, y_pred_good)
47    metrics_outliers = calculate_all_metrics(y_true,
48                                              y_pred_outliers)
49    metrics_biased = calculate_all_metrics(y_true, y_pred_biased)
50
51    print("METRIC DEFINITIONS:")
52    print("-" * 20)
53    print("MSE (Mean Squared Error): Average of squared
      differences")
54    print("  • Formula: (1/n) Σ(actual - predicted)²")

```

```

54     print("  • Units: Original units squared")
55     print("  • Heavily penalizes large errors")
56     print()
57     print("RMSE (Root Mean Squared Error): Square root of MSE")
58     print("  • Formula:  $\sqrt{MSE}$ ")
59     print("  • Units: Same as original data")
60     print("  • Interpretable scale, penalizes large errors")
61     print()
62     print("MAE (Mean Absolute Error): Average of absolute
differences")
63     print("  • Formula:  $(1/n) \sum |actual - predicted|$ ")
64     print("  • Units: Same as original data")
65     print("  • Robust to outliers")
66
67     # Create comparison table
68     print(f"\nMETRIC COMPARISON:")
69     print("-" * 18)
70     print(f"{'Scenario':<15} {'MSE':<8} {'RMSE':<8} {'MAE':<8}
{'MAPE':<8} {'Max_Err':<8}")
71     print("-" * 65)
72     print(f"{'Good Model':<15} {metrics_good['MSE']:<8.1f}
{metrics_good['RMSE']:<8.1f} {metrics_good['MAE']:<8.1f}
{metrics_good['MAPE']:<8.1f}% {metrics_good['Max_Error']:<8.1f}")
73     print(f"{'With Outliers':<15} {metrics_outliers['MSE']:<8.1f}
{metrics_outliers['RMSE']:<8.1f} {metrics_outliers['MAE']:<8.1f}
{metrics_outliers['MAPE']:<8.1f}% {metrics_outliers['Max_Error']:<8.1f}")
74     print(f"{'Biased Model':<15} {metrics_biased['MSE']:<8.1f}
{metrics_biased['RMSE']:<8.1f} {metrics_biased['MAE']:<8.1f}
{metrics_biased['MAPE']:<8.1f}% {metrics_biased['Max_Error']:<8.1f}")
75
76     # Visualize the differences
77     fig, axes = plt.subplots(2, 3, figsize=(18, 12))
78
79     scenarios = [
80         (y_pred_good, 'Good Model', 'green'),
81         (y_pred_outliers, 'With Outliers', 'red'),
82         (y_pred_biased, 'Biased Model', 'blue')
83     ]
84
85     for i, (y_pred, title, color) in enumerate(scenarios):
86         # Scatter plot
87         axes[0, i].scatter(y_true, y_pred, alpha=0.6, color=color)
88         axes[0, i].plot([y_true.min(), y_true.max()],
[y_true.min(), y_true.max()], 'k--', alpha=0.7)
89         axes[0, i].set_xlabel('True Values')
90         axes[0, i].set_ylabel('Predicted Values')
91         axes[0, i].set_title(f'{title}\nRMSE:
{np.sqrt(np.mean((y_true - y_pred)**2)):.1f}, MAE:

```

```

    {np.mean(np.abs(y_true - y_pred)):.1f}')}
92     axes[0, i].grid(True, alpha=0.3)
93
94     # Residual plot
95     residuals = y_true - y_pred
96     axes[1, i].scatter(y_pred, residuals, alpha=0.6,
97     color=color)
97     axes[1, i].axhline(y=0, color='k', linestyle='--',
98     alpha=0.7)
98     axes[1, i].set_xlabel('Predicted Values')
99     axes[1, i].set_ylabel('Residuals')
100    axes[1, i].set_title(f'Residuals - {title}')
101    axes[1, i].grid(True, alpha=0.3)
102
103    plt.tight_layout()
104    plt.show()
105
106    print("\nWHEN TO USE EACH METRIC:")
107    print("-" * 27)
108    print("USE MSE WHEN:")
109    print("  ✓ You want to heavily penalize large errors")
110    print("  ✓ Optimizing model parameters (differentiable)")
111    print("  ✓ Mathematical convenience in algorithms")
112    print()
113    print("USE RM
114 <user>: continue
115 <Claude Sonnet 4>: ````python
116    print("USE RMSE WHEN:")
117    print("  ✓ You want interpretable error in original units")
118    print("  ✓ Large errors are more costly than small ones")
119    print("  ✓ Reporting to stakeholders (easier to understand)")
120    print("  ✓ Comparing models on the same scale as your data")
121    print()
122    print("USE MAE WHEN:")
123    print("  ✓ You want equal weight for all errors")
124    print("  ✓ Data has outliers that shouldn't dominate")
125    print("  ✓ Robust evaluation is more important")
126    print("  ✓ All errors have similar business impact")
127
128    # Demonstrate sensitivity to outliers
129    print(f"\nOUTLIER SENSITIVITY ANALYSIS:")
130    print("-" * 32)
131
132    # Compare metrics with and without outliers
133    y_clean = y_true[:95]  # Remove outlier points
134    y_pred_clean = y_pred_outliers[:95]
135    y_pred_with_outliers = y_pred_outliers
136
137    metrics_clean = calculate_all_metrics(y_clean, y_pred_clean)
138    metrics_with_outliers = calculate_all_metrics(y_true,

```

```

y_pred_with_outliers)

139     print("Impact of 5 outliers on metrics:")
140     print(f"RMSE: {metrics_clean['RMSE']:.1f} →
141         {metrics_with_outliers['RMSE']:.1f}
142         ({((metrics_with_outliers['RMSE'])/metrics_clean['RMSE'])-1)*100:+.0f}%")
143     print(f"MAE: {metrics_clean['MAE']:.1f} →
144         {metrics_with_outliers['MAE']:.1f}
145         ({((metrics_with_outliers['MAE'])/metrics_clean['MAE'])-1)*100:+.0f}%")
146     print("→ RMSE is more sensitive to outliers than MAE")

147     return {
148         'good': metrics_good,
149         'outliers': metrics_outliers,
150         'biased': metrics_biased
151     }
152 error_metrics = error_metrics_comprehensive()

```

User: This comparison really helps me understand when to use each metric! Now I'm curious about something I keep hearing about - sensitivity and recall. Are these the same thing? And how do they relate to the precision and recall we learned about earlier?

Expert: Great question! **Sensitivity and recall are actually the same thing** - just different names used in different contexts. Let me clarify this and show you the complete picture of classification metrics.

```

1 def sensitivity_recall_comprehensive():
2     """
3     Comprehensive explanation of sensitivity, recall, and related
4     metrics
5     """
6     print("SENSITIVITY, RECALL, AND CLASSIFICATION METRICS")
7     print("=" * 50)
8
8     print("TERMINOLOGY CLARIFICATION:")
9     print("-" * 28)
10    print("SENSITIVITY = RECALL = TRUE POSITIVE RATE")
11    print("• Medical/Scientific context: 'Sensitivity'")
12    print("• Machine Learning context: 'Recall'")
13    print("• Signal Detection context: 'True Positive Rate'")
14    print("• Same formula: TP / (TP + FN)")
15    print()
16
17    # Generate sample medical diagnosis data
18    np.random.seed(42)
19    n_patients = 1000
20

```

```
21     # True disease status (1 = has disease, 0 = healthy)
22     true_disease = np.random.choice([0, 1], n_patients, p=[0.9,
23                                     0.1]) # 10% have disease
24
25     # Simulate test results with realistic performance
26     predicted_disease = []
27     for i in range(n_patients):
28         if true_disease[i] == 1: # Actually has disease
29             # 85% chance test detects it (sensitivity = 0.85)
30             pred = np.random.choice([0, 1], p=[0.15, 0.85])
31         else: # Actually healthy
32             # 90% chance test correctly identifies healthy
33             (specificity = 0.90)
34             pred = np.random.choice([0, 1], p=[0.90, 0.10])
35         predicted_disease.append(pred)
36
37     predicted_disease = np.array(predicted_disease)
38
39     # Calculate confusion matrix
40     from sklearn.metrics import confusion_matrix,
41     classification_report
42     cm = confusion_matrix(true_disease, predicted_disease)
43     tn, fp, fn, tp = cm.ravel()
44
45     # Calculate all metrics
46     sensitivity_recall = tp / (tp + fn) # Same as recall
47     specificity = tn / (tn + fp)
48     precision_ppv = tp / (tp + fp) # Positive Predictive Value
49     npv = tn / (tn + fn) # Negative Predictive Value
50     accuracy = (tp + tn) / (tp + tn + fp + fn)
51     f1_score = 2 * (precision_ppv * sensitivity_recall) /
52     (precision_ppv + sensitivity_recall)
53
54     print("COMPLETE METRICS BREAKDOWN:")
55     print("-" * 30)
56     print(f"True Positives (TP): {tp:3d} - Correctly identified
disease")
57     print(f"True Negatives (TN): {tn:3d} - Correctly identified
healthy")
58     print(f"False Positives (FP): {fp:3d} - Incorrectly flagged as
disease")
59     print(f"False Negatives (FN): {fn:3d} - Missed disease cases")
60     print()
61
62     print("SENSITIVITY/RECALL FAMILY:")
63     print("-" * 28)
64     print(f"Sensitivity/Recall: {sensitivity_recall:.3f} - Of
diseased patients, % correctly identified")
65     print(f"Specificity: {specificity:.3f} - Of healthy
patients, % correctly identified")
```

```

62     print(f"Precision/PPV:           {precision_ppv:.3f} - Of
positive tests, % actually have disease")
63     print(f"Negative Pred. Value:   {npv:.3f} - Of negative tests,
% actually healthy")
64     print(f"Accuracy:                 {accuracy:.3f} - Overall
correct predictions")
65     print(f"F1-Score:                  {f1_score:.3f} - Harmonic mean
of precision & recall")
66
67     # Visualize the confusion matrix and metrics
68     fig, axes = plt.subplots(2, 3, figsize=(18, 12))
69
70     # Confusion Matrix
71     im = axes[0, 0].imshow(cm, interpolation='nearest',
cmap='Blues')
72     axes[0, 0].set_title('Confusion Matrix\n(Medical Test
Example)')
73
74     # Add text annotations
75     thresh = cm.max() / 2.
76     for i in range(2):
77         for j in range(2):
78             axes[0, 0].text(j, i, f'{cm[i, j]}\n{cm[i,
j]/n_patients:.1%})',
79                           ha="center", va="center",
80                           color="white" if cm[i, j] > thresh else
"black",
81                           fontsize=12, fontweight='bold')
82
83     axes[0, 0].set_ylabel('True Label')
84     axes[0, 0].set_xlabel('Predicted Label')
85     axes[0, 0].set_xticks([0, 1])
86     axes[0, 0].set_xticklabels(['Healthy', 'Disease'])
87     axes[0, 0].set_yticks([0, 1])
88     axes[0, 0].set_yticklabels(['Healthy', 'Disease'])
89
90     # Sensitivity vs Specificity Trade-off
91     # Simulate different threshold values
92     thresholds = np.linspace(0, 1, 100)
93     sensitivities = []
94     specificities = []
95
96     for threshold in thresholds:
97         # Simulate how changing threshold affects predictions
98         sens = 1 - threshold * 0.8 # As threshold increases,
sensitivity decreases
99         spec = threshold * 0.9      # As threshold increases,
specificity increases
100        sensitivities.append(max(0, min(1, sens)))
101        specificities.append(max(0, min(1, spec)))

```

```

102
103     axes[0, 1].plot(1 - np.array(specificities), sensitivities, 'b-
104     ', linewidth=2)
105     axes[0, 1].plot([0, 1], [0, 1], 'k--', alpha=0.5)
106     axes[0, 1].scatter(1 - specificity, sensitivity_recall,
107     color='red', s=100, zorder=5)
108     axes[0, 1].set_xlabel('1 - Specificity (False Positive Rate)')
109     axes[0, 1].set_ylabel('Sensitivity (True Positive Rate)')
110     axes[0, 1].set_title('ROC Curve\n(Receiver Operating
111     Characteristic)')
112     axes[0, 1].grid(True, alpha=0.3)
113     axes[0, 1].annotate(f'Current Test\n(Sens=
114     {sensitivity_recall:.2f}, Spec={specificity:.2f})',
115             xy=(1-specificity, sensitivity_recall),
116             xytext=(0.6, 0.3),
117             arrowprops=dict(arrowstyle='->',
118             color='red'))
119
120     # Precision-Recall Trade-off
121     precisions = []
122     recalls = []
123
124     for threshold in thresholds:
125         # Simulate precision-recall trade-off
126         recall = 1 - threshold * 0.8
127         precision = 0.3 + threshold * 0.6 # Higher threshold →
128         higher precision
129         recalls.append(max(0, min(1, recall)))
130         precisions.append(max(0, min(1, precision)))
131
132         axes[0, 2].plot(recalls, precisions, 'g-', linewidth=2)
133         axes[0, 2].scatter(sensitivity_recall, precision_ppv,
134         color='red', s=100, zorder=5)
135         axes[0, 2].set_xlabel('Recall (Sensitivity)')
136         axes[0, 2].set_ylabel('Precision (PPV)')
137         axes[0, 2].set_title('Precision-Recall Curve')
138         axes[0, 2].grid(True, alpha=0.3)
139         axes[0, 2].annotate(f'Current Test\n(Prec={precision_ppv:.2f},
140         Rec={sensitivity_recall:.2f})',
141             xy=(sensitivity_recall, precision_ppv),
142             xytext=(0.3, 0.8),
143             arrowprops=dict(arrowstyle='->',
144             color='red'))
145
146     # Context-specific interpretations
147     contexts = [
148         ('Medical Screening', 'High Sensitivity', "Don't miss any
149         diseases"),
150         ('Spam Detection', 'High Precision', "Don't block
151         important emails"),

```

```

139         ('Fraud Detection', 'Balanced F1', 'Balance catching fraud
140         vs false alarms')
141
142     context_colors = ['red', 'blue', 'green']
143
144     axes[1, 0].axis('off')
145     axes[1, 0].set_title('Context-Specific Priorities',
146                           fontweight='bold')
147
148     for i, (context, priority, reason) in enumerate(contexts):
149         y_pos = 0.8 - i * 0.25
150         axes[1, 0].text(0.1, y_pos, f'{context}:',
151                         fontweight='bold', fontsize=12, color=context_colors[i])
152         axes[1, 0].text(0.1, y_pos-0.05, f'Priority: {priority}',
153                         fontsize=10)
154         axes[1, 0].text(0.1, y_pos-0.10, f'Reason: {reason}',
155                         fontsize=10, style='italic')
156
157     # Metrics comparison chart
158     metrics_names = ['Sensitivity\n(Recall)', 'Specificity',
159                       'Precision\n(PPV)', 'NPV', 'Accuracy', 'F1-Score']
160     metrics_values = [sensitivity_recall, specificity,
161                        precision_ppv, npv, accuracy, f1_score]
162
163     bars = axes[1, 1].bar(metrics_names, metrics_values,
164                           color=['red', 'blue', 'green', 'orange',
165                               'purple', 'brown'],
166                               alpha=0.7)
167     axes[1, 1].set_ylabel('Score')
168     axes[1, 1].set_title('All Classification Metrics')
169     axes[1, 1].set_ylim(0, 1)
170     axes[1, 1].tick_params(axis='x', rotation=45)
171     axes[1, 1].grid(True, alpha=0.3)
172
173     # Add value labels on bars
174     for bar, value in zip(bars, metrics_values):
175         axes[1, 1].text(bar.get_x() + bar.get_width()/2,
176                         bar.get_height() + 0.01,
177                         f'{value:.3f}', ha='center', va='bottom',
178                         fontweight='bold')
179
180     # Business decision framework
181     axes[1, 2].axis('off')
182
183     decision_text = """
184 BUSINESS DECISION FRAMEWORK:
185 MEDICAL DIAGNOSIS:
186 ✓ High Sensitivity (Recall): {sensitivity_recall:.1%}
187     → Catch most disease cases

```

```

180 ✓ Accept lower Specificity: {specificity:.1%}
181   → Some false alarms OK
182 SPAM DETECTION:
184 ✓ High Precision: Minimize false positives
185 ✓ Accept lower Recall: Some spam gets through
186   → Better than blocking important emails
188 FRAUD DETECTION:
189 ✓ Balance Sensitivity & Precision
190 ✓ Use F1-Score: {f1_score:.1%}
191   → Optimize overall performance
193 KEY INSIGHT:
194 The "best" metric depends on the cost of
195 different types of errors in your specific
196 business context.
197 """
198
199     axes[1, 2].text(0.05, 0.95, decision_text, transform=axes[1,
2].transAxes,
200                     fontsize=10, verticalalignment='top',
201                     fontfamily='monospace',
202                     bbox=dict(boxstyle='round',
203                               facecolor='lightcyan', alpha=0.8))
204
205
206     print("\nPRACTICAL GUIDELINES:")
207     print("-" * 22)
208     print("PRIORITIZE SENSITIVITY/RECALL WHEN:")
209     print("  • Missing positives is very costly (medical
diagnosis)")
210     print("  • You want to catch as many cases as possible")
211     print("  • False negatives are worse than false positives")
212     print()
213     print("PRIORITIZE SPECIFICITY WHEN:")
214     print("  • False positives are very costly")
215     print("  • You need to be sure about negative cases")
216     print("  • Resources for follow-up are limited")
217     print()
218     print("PRIORITIZE PRECISION WHEN:")
219     print("  • Acting on positive predictions is expensive")
220     print("  • You need high confidence in positive predictions")
221     print("  • False alarms cause significant problems")
222
223     return {
224         'confusion_matrix': cm,
225         'sensitivity': sensitivity_recall,
226         'specificity': specificity,
227         'precision': precision_ppv,
228         'npv': npv,

```

```
229         'accuracy': accuracy,
230         'f1_score': f1_score
231     }
232 classification_metrics = sensitivity_recall_comprehensive()
```

User: This is incredibly comprehensive! I love how you showed that sensitivity and recall are the same thing, just used in different contexts. Now I'm really curious about ROC curves - I see them mentioned a lot but I'm not sure how to interpret them properly.

Expert: Excellent! **ROC curves** (Receiver Operating Characteristic) are one of the most important tools for evaluating classification models. Let me give you a complete understanding of how to create, interpret, and use them effectively.

```
1 def roc_curve_comprehensive():
2     """
3     Comprehensive explanation of ROC curves and AUC
4     """
5     print("ROC CURVES AND AUC: COMPLETE GUIDE")
6     print("=" * 40)
7
8     print("WHAT IS A ROC CURVE?")
9     print("-" * 22)
10    print("• Plots True Positive Rate (Sensitivity) vs False
11        Positive Rate (1-Specificity)")
12    print("• Shows performance across all classification
13        thresholds")
14    print("• Helps choose optimal threshold for your business
15        needs")
16    print("• AUC (Area Under Curve) summarizes overall
17        performance")
18    print()
19
20    # Generate sample data with different model qualities
21    np.random.seed(42)
22    n_samples = 1000
23
24    # True labels (20% positive class)
25    y_true = np.random.choice([0, 1], n_samples, p=[0.8, 0.2])
26
27    # Simulate different model prediction probabilities
28
29    # Model 1: Excellent model
30    y_prob_excellent = []
31    for label in y_true:
32        if label == 1:
33            prob = np.random.beta(8, 2) # High probabilities for
34            positive class
```

```

30     else:
31         prob = np.random.beta(2, 8) # Low probabilities for
32         # negative class
33         y_prob_excellent.append(prob)
34         y_prob_excellent = np.array(y_prob_excellent)
35
36     # Model 2: Good model
37     y_prob_good = []
38     for label in y_true:
39         if label == 1:
40             prob = np.random.beta(5, 3)
41         else:
42             prob = np.random.beta(3, 5)
43             y_prob_good.append(prob)
44             y_prob_good = np.array(y_prob_good)
45
46     # Model 3: Poor model
47     y_prob_poor = []
48     for label in y_true:
49         if label == 1:
50             prob = np.random.beta(3, 4)
51         else:
52             prob = np.random.beta(4, 3)
53             y_prob_poor.append(prob)
54             y_prob_poor = np.array(y_prob_poor)
55
56     # Model 4: Random model
57     y_prob_random = np.random.uniform(0, 1, n_samples)
58
59     # Calculate ROC curves
60     from sklearn.metrics import roc_curve, auc
61
62     models = {
63         'Excellent Model': y_prob_excellent,
64         'Good Model': y_prob_good,
65         'Poor Model': y_prob_poor,
66         'Random Model': y_prob_random
67     }
68
69     colors = ['green', 'blue', 'orange', 'red']
70
71     fig, axes = plt.subplots(2, 3, figsize=(18, 12))
72
73     # Plot ROC curves
74     axes[0, 0].plot([0, 1], [0, 1], 'k--', alpha=0.5, label='Random
75     Classifier')
76
77     roc_data = {}
78
79     for i, (model_name, y_prob) in enumerate(models.items()):
80
81         fpr, tpr, thresholds = roc_curve(y_true, y_prob)
82
83         auc_value = auc(fpr, tpr)
84
85         axes[0, i].plot(fpr, tpr, color=colors[i], label=f'{model_name} (AUC = {auc_value:.2f})')
86
87         roc_data[model_name] = {'fpr': fpr, 'tpr': tpr, 'thresholds': thresholds}
88
89     # Add overall legend for ROC curves
90     axes[0, 0].legend()
91
92     # Add overall title for the first subplot
93     axes[0, 0].set_title('ROC Curves for Four Models')
94
95     # Set overall x-axis label for the bottom row
96     axes[1, 0].set_xlabel('False Positive Rate (FPR)')
97
98     # Set overall y-axis label for the left column
99     axes[0, 0].set_ylabel('True Positive Rate (TPR)')
```

```

78     fpr, tpr, thresholds = roc_curve(y_true, y_prob)
79     roc_auc = auc(fpr, tpr)
80
81     axes[0, 0].plot(fpr, tpr, color=colors[i], linewidth=2,
82                      label=f'{model_name} (AUC = {roc_auc:.3f})')
83
84     roc_data[model_name] = {
85         'fpr': fpr,
86         'tpr': tpr,
87         'thresholds': thresholds,
88         'auc': roc_auc
89     }
90
91     axes[0, 0].set_xlabel('False Positive Rate (1 - Specificity)')
92     axes[0, 0].set_ylabel('True Positive Rate (Sensitivity)')
93     axes[0, 0].set_title('ROC Curves Comparison')
94     axes[0, 0].legend(loc='lower right')
95     axes[0, 0].grid(True, alpha=0.3)
96
97     # Show probability distributions
98     axes[0, 1].hist(y_prob_excellent[y_true == 0], bins=30,
99                      alpha=0.7,
100                     label='Negative Class', color='red',
101                     density=True)
102     axes[0, 1].hist(y_prob_excellent[y_true == 1], bins=30,
103                      alpha=0.7,
104                     label='Positive Class', color='green',
105                     density=True)
106     axes[0, 1].axvline(0.5, color='black', linestyle='--',
107                         alpha=0.7, label='Threshold = 0.5')
108     axes[0, 1].set_xlabel('Predicted Probability')
109     axes[0, 1].set_ylabel('Density')
110     axes[0, 1].set_title('Excellent Model: Probability
111 Distributions')
112     axes[0, 1].legend()
113     axes[0, 1].grid(True, alpha=0.3)
114
115     # Find optimal threshold (closest to top-left corner)
116     distances = np.sqrt((excellent_fpr - 0)**2 + (excellent_tpr -
117 1)**2)
118     optimal_idx = np.argmin(distances)
119     optimal_threshold = excellent_thresholds[optimal_idx]
120     optimal_fpr = excellent_fpr[optimal_idx]
121     optimal_tpr = excellent_tpr[optimal_idx]

```

```

120
121     axes[0, 2].plot(excellent_fpr, excellent_tpr, 'g-',
122     linewidth=2, label='ROC Curve')
123     axes[0, 2].scatter(optimal_fpr, optimal_tpr, color='red',
124     s=100, zorder=5,
125                 label=f'Optimal Threshold =
126 {optimal_threshold:.3f}')
127     axes[0, 2].plot([0, 1], [0, 1], 'k--', alpha=0.5)
128     axes[0, 2].set_xlabel('False Positive Rate')
129     axes[0, 2].set_ylabel('True Positive Rate')
130     axes[0, 2].set_title('Threshold Selection')
131     axes[0, 2].legend()
132     axes[0, 2].grid(True, alpha=0.3)
133
134     # AUC interpretation
135     axes[1, 0].axis('off')
136
137     auc_text = """
138 AUC INTERPRETATION GUIDE:
139 AUC = 1.0: Perfect classifier
140 AUC = 0.9–1.0: Excellent
141 AUC = 0.8–0.9: Good
142 AUC = 0.7–0.8: Fair
143 AUC = 0.6–0.7: Poor
144 AUC = 0.5: Random (no skill)
145 AUC < 0.5: Worse than random
146 OUR MODELS:
147 • Excellent: {roc_data['Excellent Model']['auc']:.3f}
148 • Good: {roc_data['Good Model']['auc']:.3f}
149 • Poor: {roc_data['Poor Model']['auc']:.3f}
150 Random: {roc_data['Random Model']['auc']:.3f}
151 BUSINESS MEANING:
152 AUC = Probability that model ranks
153 a random positive example higher
154 than a random negative example.
155 """
156
157     axes[1, 0].text(0.05, 0.95, auc_text, transform=axes[1,
158 0].transAxes,
159                 fontsize=11, verticalalignment='top',
160                 fontfamily='monospace',
161                 bbox=dict(boxstyle='round',
162 facecolor='lightgreen', alpha=0.8))
163
164     # Threshold analysis table
165     axes[1, 1].axis('off')
166
167     # Calculate metrics at different thresholds for excellent model
168     thresholds_to_analyze = [0.3, 0.5, 0.7, optimal_threshold]

```

```

167     threshold_analysis = []
168     for thresh in thresholds_to_analyze:
169         y_pred = (y_prob_excellent >= thresh).astype(int)
170
171         from sklearn.metrics import confusion_matrix,
172         precision_score, recall_score
173         cm = confusion_matrix(y_true, y_pred)
174         tn, fp, fn, tp = cm.ravel()
175
176         precision = precision_score(y_true, y_pred)
177         recall = recall_score(y_true, y_pred)
178         specificity = tn / (tn + fp)
179
180         threshold_analysis.append({
181             'threshold': thresh,
182             'precision': precision,
183             'recall': recall,
184             'specificity': specificity,
185             'fpr': 1 - specificity
186         })
187
188     # Create table
189     table_data = []
190     table_data.append(['Threshold', 'Precision', 'Recall',
191 'Specificity', 'FPR'])
192
193     for analysis in threshold_analysis:
194         row = [
195             f'{analysis["threshold"]:.3f}',
196             f'{analysis["precision"]:.3f}',
197             f'{analysis["recall"]:.3f}',
198             f'{analysis["specificity"]:.3f}',
199             f'{analysis["fpr"]:.3f}'
200         ]
201         table_data.append(row)
202
203     table = axes[1, 1].table(cellText=table_data[1:],
204                             colLabels=table_data[0],
205                             cellLoc='center',
206                             loc='center',
207                             bbox=[0, 0.3, 1, 0.6])
208
209     table.auto_set_font_size(False)
210     table.set_fontsize(10)
211     table.scale(1, 2)
212
213     # Color code the header and optimal row
214     for i in range(5):
215         table[(0, i)].set_facecolor('#4CAF50')
216         table[(0, i)].set_text_props(weight='bold', color='white')

```

```
215
216     # Highlight optimal threshold row
217     optimal_row = 4 # Last row is optimal
218     for i in range(5):
219         table[(optimal_row, i)].set_facecolor('#FFE082')
220
221     axes[1, 1].set_title('Threshold Analysis for Excellent Model',
222                           pad=20, fontweight='bold')
223
224     # Business decision guide
225     axes[1, 2].axis('off')
226
227     decision_guide = f"""
228 BUSINESS DECISION GUIDE:
229 CHOOSE THRESHOLD BASED ON:
230 1. COST OF FALSE POSITIVES:
231     High cost → Higher threshold
232     Low cost → Lower threshold
233 2. COST OF FALSE NEGATIVES:
234     High cost → Lower threshold
235     Low cost → Higher threshold
236 3. BUSINESS CONTEXT:
237 MEDICAL SCREENING:
238 • Low threshold (0.3)
239 • High sensitivity, catch all cases
240 • Accept more false alarms
241 SPAM DETECTION:
242 • High threshold (0.7)
243 • High precision, avoid blocking
244 • Accept some spam getting through
245 FRAUD DETECTION:
246 • Optimal threshold ({optimal_threshold:.3f})
247 • Balance sensitivity & precision
248 • Minimize total cost
249 REMEMBER:
250 ROC curve shows ALL possible
251 trade-offs. Choose the point
252 that best fits your business needs!
253 """
254
255     axes[1, 2].text(0.05, 0.95, decision_guide, transform=axes[1,
256                           2].transAxes,
257                           fontsize=9, verticalalignment='top',
258                           fontfamily='monospace',
259                           bbox=dict(boxstyle='round',
260                                     facecolor='lightblue', alpha=0.8))
261
262     plt.tight_layout()
263     plt.show()
264
```

```

269     print("KEY INSIGHTS:")
270     print("-" * 15)
271     print("✓ ROC curve shows performance across all thresholds")
272     print("✓ AUC summarizes overall discriminative ability")
273     print("✓ Closer to top-left corner = better performance")
274     print("✓ Diagonal line = random performance")
275     print("✓ Choose threshold based on business costs")
276     print("✓ ROC curves are threshold-independent evaluation")
277
278     return roc_data, threshold_analysis
280 roc_results, threshold_results = roc_curve_comprehensive()

```

Chapter 7: Advanced Evaluation Metrics and Techniques

User: The ROC curves really helped me understand model evaluation better! But I feel like there are still some advanced techniques I should know about. What about cross-validation? And are there other important evaluation methods I should be aware of?

Expert: Absolutely! You're ready for the advanced evaluation techniques that separate good data scientists from great ones. Let's dive into **cross-validation**, **bootstrap sampling**, and other sophisticated evaluation methods.

```

1 def cross_validation_comprehensive():
2     """
3     Comprehensive guide to cross-validation techniques
4     """
5     print("CROSS-VALIDATION: ADVANCED MODEL EVALUATION")
6     print("=" * 50)
7
8     print("WHY CROSS-VALIDATION?")
9     print("-" * 23)
10    print("• Single train/test split can be misleading")
11    print("• Results depend on which data points end up in test
set")
12    print("• Cross-validation gives more robust performance
estimates")
13    print("• Helps detect overfitting and model instability")
14    print()
15
16    # Generate sample data
17    from sklearn.datasets import make_classification
18    from sklearn.model_selection import cross_val_score,
StratifiedKFold, TimeSeriesSplit
19    from sklearn.linear_model import LogisticRegression
20    from sklearn.ensemble import RandomForestClassifier

```

```

21     from sklearn.metrics import accuracy_score, precision_score,
22     recall_score, f1_score
23
24     np.random.seed(42)
25     X, y = make_classification(n_samples=1000, n_features=20,
26                                n_informative=10,
27                                n_redundant=10,
28                                n_clusters_per_class=1, random_state=42)
29
30     # Models to compare
31     models = {
32         'Logistic Regression': LogisticRegression(random_state=42),
33         'Random Forest': RandomForestClassifier(n_estimators=100,
34                                                random_state=42)
35     }
36
37     print("TYPES OF CROSS-VALIDATION:")
38     print("-" * 30)
39
40     # 1. K-Fold Cross-Validation
41     print("1. K-FOLD CROSS-VALIDATION")
42     print("    • Split data into k equal folds")
43     print("    • Train on k-1 folds, test on 1 fold")
44     print("    • Repeat k times, average results")
45     print()
46
47     from sklearn.model_selection import KFold
48
49     kf = KFold(n_splits=5, shuffle=True, random_state=42)
50
51     fig, axes = plt.subplots(3, 2, figsize=(15, 18))
52
53     # Visualize K-Fold splits
54     fold_colors = ['red', 'blue', 'green', 'orange', 'purple']
55
56     for i, (train_idx, test_idx) in enumerate(kf.split(X)):
57         if i < 5:  # Show first 5 folds
58             y_pos = [i] * len(train_idx) + [i] * len(test_idx)
59             colors = ['blue'] * len(train_idx) + ['red'] *
60             len(test_idx)
61             x_pos = list(train_idx) + list(test_idx)
62
63             if i == 0:  # Only show for first subplot
64                 axes[0, 0].scatter(x_pos[:len(train_idx)],
65                                     y_pos[:len(train_idx)],
66                                     c='blue', alpha=0.6, s=1,
67                                     label='Training')
68                 axes[0, 0].scatter(x_pos[len(train_idx):],
69                                     y_pos[len(train_idx):],
70                                     c='red', alpha=0.6, s=1,
71                                     label='Test')

```

```

    label='Testing')
63 <user>: continue
64 <Claude Sonnet 4>: ```python
65         else:
66             axes[0, 0].scatter(x_pos[:len(train_idx)],
67                                 c='blue', alpha=0.6, s=1)
68             axes[0, 0].scatter(x_pos[len(train_idx):],
69                                 y_pos[len(train_idx):],
70                                 c='red', alpha=0.6, s=1)
71
72             axes[0, 0].set_xlabel('Sample Index')
73             axes[0, 0].set_ylabel('Fold Number')
74             axes[0, 0].set_title('5-Fold Cross-Validation Splits')
75             axes[0, 0].legend()
76             axes[0, 0].grid(True, alpha=0.3)
77
78     # Perform cross-validation for all models
79     cv_results = {}
80
81     for model_name, model in models.items():
82         scores = cross_val_score(model, X, y, cv=5,
83         scoring='accuracy')
84         cv_results[model_name] = {
85             'scores': scores,
86             'mean': scores.mean(),
87             'std': scores.std()
88         }
89
90     # Plot cross-validation results
91     model_names = list(cv_results.keys())
92     means = [cv_results[name]['mean'] for name in model_names]
93     stds = [cv_results[name]['std'] for name in model_names]
94
95     bars = axes[0, 1].bar(model_names, means, yerr=stds, capsize=5,
96                           alpha=0.7, color=['blue', 'green'])
97     axes[0, 1].set_ylabel('Accuracy')
98     axes[0, 1].set_title('Cross-Validation Results (5-Fold)')
99     axes[0, 1].grid(True, alpha=0.3)
100
101     # Add value labels
102     for i, (bar, mean, std) in enumerate(zip(bars, means, stds)):
103         axes[0, 1].text(bar.get_x() + bar.get_width()/2,
104                         bar.get_height() + std + 0.01,
105                         f'{mean:.3f}±{std:.3f}', ha='center',
106                         va='bottom', fontweight='bold')
107
108     print("5-Fold CV Results:")
109     for name in model_names:
110         mean_acc = cv_results[name]['mean']

```

```

107     std_acc = cv_results[name]['std']
108     print(f" {name}: {mean_acc:.3f} ± {std_acc:.3f}")
109
110 # 2. Stratified K-Fold
111 print("\n2. STRATIFIED K-FOLD CROSS-VALIDATION")
112 print("    • Maintains class distribution in each fold")
113 print("    • Important for imbalanced datasets")
114 print("    • Ensures each fold is representative")
115 print()
116
117 skf = StratifiedKFold(n_splits=5, shuffle=True,
118 random_state=42)
119
120 # Compare class distributions
121 fold_distributions = []
122 for i, (train_idx, test_idx) in enumerate(skf.split(X, y)):
123     train_dist = np.bincount(y[train_idx]) / len(train_idx)
124     test_dist = np.bincount(y[test_idx]) / len(test_idx)
125     fold_distributions.append({
126         'fold': i+1,
127         'train_class_0': train_dist[0],
128         'train_class_1': train_dist[1],
129         'test_class_0': test_dist[0],
130         'test_class_1': test_dist[1]
131     })
132
133 # Visualize class distributions
134 folds = [d['fold'] for d in fold_distributions]
135 train_0 = [d['train_class_0'] for d in fold_distributions]
136 train_1 = [d['train_class_1'] for d in fold_distributions]
137 test_0 = [d['test_class_0'] for d in fold_distributions]
138 test_1 = [d['test_class_1'] for d in fold_distributions]
139
140 x = np.arange(len(folds))
141 width = 0.35
142
143 axes[1, 0].bar(x - width/2, train_0, width, label='Train Class
144 0', alpha=0.7, color='lightblue')
145 axes[1, 0].bar(x - width/2, train_1, width, bottom=train_0,
146 label='Train Class 1', alpha=0.7, color='lightcoral')
147 axes[1, 0].bar(x + width/2, test_0, width, label='Test Class
148 0', alpha=0.7, color='blue')
149 axes[1, 0].bar(x + width/2, test_1, width, bottom=test_0,
150 label='Test Class 1', alpha=0.7, color='red')
151
152 axes[1, 0].set_xlabel('Fold')
153 axes[1, 0].set_ylabel('Class Proportion')
154 axes[1, 0].set_title('Stratified K-Fold: Class Distribution')
155 axes[1, 0].set_xticks(x)
156 axes[1, 0].set_xticklabels([f'Fold {i}' for i in folds])

```

```

152     axes[1, 0].legend()
153     axes[1, 0].grid(True, alpha=0.3)
154
155     # 3. Time Series Cross-Validation
156     print("3. TIME SERIES CROSS-VALIDATION")
157     print("    • Respects temporal order of data")
158     print("    • No future data leakage")
159     print("    • Expanding or sliding window approach")
160     print()
161
162     # Generate time series data
163     np.random.seed(42)
164     n_time_points = 100
165     time_X = np.random.randn(n_time_points, 5)
166     time_y = np.random.randn(n_time_points)
167
168     tscv = TimeSeriesSplit(n_splits=5)
169
170     # Visualize time series splits
171     for i, (train_idx, test_idx) in enumerate(tscv.split(time_X)):
172         y_pos = [i] * len(range(n_time_points))
173         colors = ['lightgray'] * n_time_points
174
175         for idx in train_idx:
176             colors[idx] = 'blue'
177         for idx in test_idx:
178             colors[idx] = 'red'
179
180         axes[1, 1].scatter(range(n_time_points), y_pos, c=colors,
181                         alpha=0.7, s=10)
182
183         axes[1, 1].set_xlabel('Time Point')
184         axes[1, 1].set_ylabel('Fold Number')
185         axes[1, 1].set_title('Time Series Cross-Validation')
186
187         # Add legend
188         from matplotlib.patches import Patch
189         legend_elements = [Patch(facecolor='blue', label='Training'),
190                            Patch(facecolor='red', label='Testing'),
191                            Patch(facecolor='lightgray', label='Not
192                             Used')]
193
194         axes[1, 1].legend(handles=legend_elements)
195         axes[1, 1].grid(True, alpha=0.3)
196
197     # 4. Leave-One-Out Cross-Validation (LOOCV)
198     print("4. LEAVE-ONE-OUT CROSS-VALIDATION (LOOCV)")
199     print("    • Each sample is test set once")
200     print("    • n-fold CV where n = number of samples")
201     print("    • Maximum use of data, but computationally
202          expensive")

```

```
199     print()
200
201     # Demonstrate with smaller dataset
202     from sklearn.model_selection import LeaveOneOut
203
204     # Use smaller subset for LOOCV demonstration
205     X_small = X[:50]
206     y_small = y[:50]
207
208     loo = LeaveOneOut()
209     loo_scores =
210         cross_val_score(LogisticRegression(random_state=42), X_small,
211                         y_small, cv=loo)
212
213     print(f"LOOCV Results on 50 samples:")
214     print(f"  Mean Accuracy: {loo_scores.mean():.3f}")
215     print(f"  Std Accuracy: {loo_scores.std():.3f}")
216     print(f"  Number of folds: {len(loo_scores)}")
217
218     # 5. Nested Cross-Validation
219     print(f"\n5. NESTED CROSS-VALIDATION")
220     print("  • Outer loop: Model evaluation")
221     print("  • Inner loop: Hyperparameter tuning")
222     print("  • Unbiased estimate of model performance")
223     print()
224
225     # Nested CV example
226     def nested_cross_validation(X, y, model, param_grid,
227         outer_cv=5, inner_cv=3):
228         outer_scores = []
229
230         outer_kf = StratifiedKFold(n_splits=outer_cv, shuffle=True,
231             random_state=42)
232
233         for train_idx, test_idx in outer_kf.split(X, y):
234             X_train, X_test = X[train_idx], X[test_idx]
235             y_train, y_test = y[train_idx], y[test_idx]
236
237             # Inner CV for hyperparameter tuning
238             inner_cv = StratifiedKFold(n_splits=inner_cv,
239                 shuffle=True, random_state=42)
240             grid_search = GridSearchCV(model, param_grid,
241                 cv=inner_cv, scoring='accuracy')
242             grid_search.fit(X_train, y_train)
243
244             # Test best model on outer test set
245             best_model = grid_search.best_estimator_
246             score = best_model.score(X_test, y_test)
```

```

243         outer_scores.append(score)
244
245     return np.array(outer_scores)
246
247     # Perform nested CV
248     param_grid = {'C': [0.1, 1, 10], 'max_iter': [100, 1000]}
249     nested_scores = nested_cross_validation(X, y,
250                                              LogisticRegression(random_state=42), param_grid)
251
252     print(f"Nested CV Results:")
253     print(f"  Mean Accuracy: {nested_scores.mean():.3f}")
254     print(f"  Std Accuracy: {nested_scores.std():.3f}")
255
256     # Visualize different CV methods comparison
257     cv_methods = ['5-Fold', 'Stratified 5-Fold', 'LOOCV (50
258     samples)', 'Nested CV']
259     cv_means = [
260         cv_results['Logistic Regression']['mean'],
261         cv_results['Logistic Regression']['mean'], # Similar for
262         stratified
263         loo_scores.mean(),
264         nested_scores.mean()
265     ]
266     cv_stds = [
267         cv_results['Logistic Regression']['std'],
268         cv_results['Logistic Regression']['std'],
269         loo_scores.std(),
270         nested_scores.std()
271     ]
272
273     bars = axes[2, 0].bar(cv_methods, cv_means, yerr=cv_stds,
274                           capsizes=5,
275                           alpha=0.7, color=['blue', 'green',
276                           'orange', 'red'])
277     axes[2, 0].set_ylabel('Accuracy')
278     axes[2, 0].set_title('Cross-Validation Methods Comparison')
279     axes[2, 0].tick_params(axis='x', rotation=45)
280     axes[2, 0].grid(True, alpha=0.3)
281
282     # Add value labels
283     for bar, mean, std in zip(bars, cv_means, cv_stds):
284         axes[2, 0].text(bar.get_x() + bar.get_width()/2,
285                         bar.get_height() + std + 0.01,
286                         f'{mean:.3f}±{std:.3f}', ha='center',
287                         va='bottom', fontweight='bold', fontsize=9)
288
289     # Cross-validation best practices
290     axes[2, 1].axis('off')
291
292     best_practices = """

```

```

286 CROSS-VALIDATION BEST PRACTICES:
288 ✓ CHOOSE THE RIGHT METHOD:
289   • Standard: K-Fold (k=5 or k=10)
290   • Imbalanced data: Stratified K-Fold
291   • Time series: TimeSeriesSplit
292   • Small datasets: LOOCV
293   • Hyperparameter tuning: Nested CV
295 ✓ CONSIDERATIONS:
296   • More folds = less bias, more variance
297   • Fewer folds = more bias, less variance
298   • Computational cost increases with folds
300 ✓ COMMON MISTAKES TO AVOID:
301   • Data leakage between folds
302   • Not stratifying imbalanced data
303   • Using future data in time series
304   • Tuning hyperparameters on test set
305 ✓ REPORTING RESULTS:
307   • Always report mean ± standard deviation
308   • Consider confidence intervals
309   • Test statistical significance of differences
310   • Validate on truly held-out test set
311   """
312
313     axes[2, 1].text(0.05, 0.95, best_practices, transform=axes[2,
314                               1].transAxes,
315                               fontsize=10, verticalalignment='top',
316                               fontfamily='monospace',
317                               bbox=dict(boxstyle='round',
318                               facecolor='lightcyan', alpha=0.8))
319
320     plt.tight_layout()
321     plt.show()
322
320     return cv_results, nested_scores
322 cv_results, nested_scores = cross_validation_comprehensive()

```

User: Cross-validation makes so much sense now! I can see how it gives us much more confidence in our model evaluation. But what about bootstrap sampling? I've heard it's another important technique for understanding uncertainty.

Expert: Excellent question! **Bootstrap sampling** is a powerful statistical technique that complements cross-validation beautifully. It's particularly useful for estimating confidence intervals and understanding the stability of your model performance.

```

1 def bootstrap_sampling_comprehensive():
2     """
3     Comprehensive guide to bootstrap sampling for model evaluation

```

```
4     """
5     print("BOOTSTRAP SAMPLING FOR MODEL EVALUATION")
6     print("=" * 45)
7
8     print("WHAT IS BOOTSTRAP SAMPLING?")
9     print("-" * 32)
10    print("• Sampling WITH replacement from your dataset")
11    print("• Creates multiple 'bootstrap samples' of same size as
original")
12    print("• Each sample is slightly different due to random
sampling")
13    print("• Estimates sampling distribution of any statistic")
14    print("• Provides confidence intervals without assumptions")
15    print()
16
17    # Generate sample data
18    from sklearn.datasets import make_regression
19    from sklearn.linear_model import LinearRegression
20    from sklearn.metrics import mean_squared_error, r2_score
21
22    np.random.seed(42)
23    X, y = make_regression(n_samples=200, n_features=5, noise=10,
random_state=42)
24
25    # Original model performance
26    model = LinearRegression()
27    model.fit(X, y)
28    y_pred = model.predict(X)
29    original_r2 = r2_score(y, y_pred)
30    original_mse = mean_squared_error(y, y_pred)
31
32    print("BOOTSTRAP PROCEDURE:")
33    print("-" * 21)
34    print("1. Sample n observations WITH replacement")
35    print("2. Train model on bootstrap sample")
36    print("3. Evaluate model performance")
37    print("4. Repeat B times (typically B = 1000)")
38    print("5. Analyze distribution of performance metrics")
39    print()
40
41    # Perform bootstrap sampling
42    def bootstrap_model_evaluation(X, y, model, n_bootstrap=1000):
43        """
44            Perform bootstrap evaluation of model performance
45        """
46        n_samples = len(X)
47        bootstrap_scores = {
48            'r2_scores': [],
49            'mse_scores': [],
50            'coefficients': [],
```

```

51         'intercepts': []
52     }
53
54     for i in range(n_bootstrap):
55         # Bootstrap sample
56         bootstrap_indices = np.random.choice(n_samples,
57         size=n_samples, replace=True)
58         X_bootstrap = X[bootstrap_indices]
59         y_bootstrap = y[bootstrap_indices]
60
61         # Train model on bootstrap sample
62         bootstrap_model = LinearRegression()
63         bootstrap_model.fit(X_bootstrap, y_bootstrap)
64
65         # Evaluate on bootstrap sample
66         y_pred_bootstrap = bootstrap_model.predict(X_bootstrap)
67         r2_bootstrap = r2_score(y_bootstrap, y_pred_bootstrap)
68         mse_bootstrap = mean_squared_error(y_bootstrap,
69         y_pred_bootstrap)
70
71         # Store results
72         bootstrap_scores['r2_scores'].append(r2_bootstrap)
73         bootstrap_scores['mse_scores'].append(mse_bootstrap)
74
75         bootstrap_scores['coefficients'].append(bootstrap_model.coef_)
76
77         bootstrap_scores['intercepts'].append(bootstrap_model.intercept_)
78
79     return bootstrap_scores
80
81     print("Performing 1000 bootstrap samples...")
82     bootstrap_results = bootstrap_model_evaluation(X, y, model,
83     n_bootstrap=1000)
84
85     # Convert to numpy arrays for easier analysis
86     r2_scores = np.array(bootstrap_results['r2_scores'])
87     mse_scores = np.array(bootstrap_results['mse_scores'])
88     coefficients = np.array(bootstrap_results['coefficients'])
89     intercepts = np.array(bootstrap_results['intercepts'])
90
91     # Calculate confidence intervals
92     def calculate_confidence_intervals(data,
93     confidence_level=0.95):
94         """
95             Calculate confidence intervals using percentile method
96         """
97         alpha = 1 - confidence_level
98         lower_percentile = (alpha/2) * 100
99         upper_percentile = (1 - alpha/2) * 100

```

```

95         lower_bound = np.percentile(data, lower_percentile)
96         upper_bound = np.percentile(data, upper_percentile)
97
98     return lower_bound, upper_bound
99
100    # Calculate CIs for performance metrics
101    r2_ci_lower, r2_ci_upper =
102        calculate_confidence_intervals(r2_scores)
103    mse_ci_lower, mse_ci_upper =
104        calculate_confidence_intervals(mse_scores)
105
106    print("BOOTSTRAP RESULTS:")
107    print("-" * 18)
108    print(f"Original R2 Score: {original_r2:.4f}")
109    print(f"Bootstrap R2 Mean: {np.mean(r2_scores):.4f}")
110    print(f"Bootstrap R2 Std: {np.std(r2_scores):.4f}")
111    print(f"95% CI for R2: [{r2_ci_lower:.4f}, {r2_ci_upper:.4f}]")
112
113    print()
114    print(f"Original MSE: {original_mse:.2f}")
115    print(f"Bootstrap MSE Mean: {np.mean(mse_scores):.2f}")
116    print(f"Bootstrap MSE Std: {np.std(mse_scores):.2f}")
117    print(f"95% CI for MSE: [{mse_ci_lower:.2f}, {mse_ci_upper:.2f}]")
118
119    # Visualize bootstrap results
120    fig, axes = plt.subplots(3, 2, figsize=(15, 18))
121
122    # R2 distribution
123    axes[0, 0].hist(r2_scores, bins=50, alpha=0.7, color='blue',
124    density=True)
125    axes[0, 0].axvline(original_r2, color='red', linestyle='--',
126    linewidth=2, label=f'Original R2 ({original_r2:.3f})')
127    axes[0, 0].axvline(np.mean(r2_scores), color='green',
128    linestyle='--', linewidth=2, label=f'Bootstrap Mean
129    ({np.mean(r2_scores):.3f})')
130    axes[0, 0].axvline(r2_ci_lower, color='orange', linestyle=':',
131    alpha=0.7, label=f'95% CI')
132    axes[0, 0].axvline(r2_ci_upper, color='orange', linestyle=':',
133    alpha=0.7)
134
135    axes[0, 0].set_xlabel('R2 Score')
136    axes[0, 0].set_ylabel('Density')
137    axes[0, 0].set_title('Bootstrap Distribution of R2 Scores')
138    axes[0, 0].legend()
139    axes[0, 0].grid(True, alpha=0.3)
140
141    # MSE distribution
142    axes[0, 1].hist(mse_scores, bins=50, alpha=0.7, color='red',
143    density=True)
144    axes[0, 1].axvline(original_mse, color='blue', linestyle='--',
145    linewidth=2, label=f'Original MSE ({original_mse:.2f})')
146    axes[0, 1].axvline(np.mean(mse_scores), color='green',
147    linestyle='--', linewidth=2, label=f'Bootstrap Mean
148    ({np.mean(mse_scores):.2f})')
149    axes[0, 1].axvline(mse_ci_lower, color='orange', linestyle=':',
150    alpha=0.7, label=f'95% CI')
151    axes[0, 1].axvline(mse_ci_upper, color='orange', linestyle=':',
152    alpha=0.7)
153
154    axes[0, 1].set_xlabel('MSE')
155    axes[0, 1].set_ylabel('Density')
156    axes[0, 1].set_title('Bootstrap Distribution of MSE')
157    axes[0, 1].legend()
158    axes[0, 1].grid(True, alpha=0.3)
159
160    # Correlation matrix heatmap
161    correlation_matrix = np.corrcoef(r2_scores, mse_scores)
162    heatmap = sns.heatmap(correlation_matrix, annot=True,
163    cmap='coolwarm', square=True, xticklabels=['R2'],
164    yticklabels=['MSE'])
165    heatmap.set_title('Correlation Matrix Heatmap')
166
167    # Final summary table
168    summary_table = pd.DataFrame({
169        'Metric': ['Original', 'Bootstrap'],
170        'Value': [original_r2, np.mean(r2_scores)],
171        'CI_Lower': [r2_ci_lower, r2_ci_upper],
172        'CI_Upper': [r2_ci_upper, r2_ci_lower]
173    })
174    summary_table.set_index('Metric', inplace=True)
175    print(summary_table)
176
177    # Save figures
178    fig.savefig('bootstrap_results.png')
179    heatmap.savefig('correlation_matrix.png')
180
181    # Clean up
182    plt.close('all')
183
```

```

    linewidth=2, label=f'Original MSE ({original_mse:.1f})')
134     axes[0, 1].axvline(np.mean(mse_scores), color='green',
135     linestyle='--', linewidth=2, label=f'Bootstrap Mean
136     ({np.mean(mse_scores):.1f})')
137     axes[0, 1].axvline(mse_ci_lower, color='orange', linestyle='--',
138     alpha=0.7, label=f'95% CI')
139     axes[0, 1].axvline(mse_ci_upper, color='orange', linestyle='--',
140     alpha=0.7)
141     axes[0, 1].set_xlabel('MSE')
142     axes[0, 1].set_ylabel('Density')
143     axes[0, 1].set_title('Bootstrap Distribution of MSE')
144     axes[0, 1].legend()
145     axes[0, 1].grid(True, alpha=0.3)
146
147     # Coefficient stability
148     feature_names = [f'Feature_{i+1}' for i in
149     range(coefficients.shape[1])]
150
151     # Box plot of coefficients
152     axes[1, 0].boxplot(coefficients, labels=feature_names)
153     axes[1, 0].set_ylabel('Coefficient Value')
154     axes[1, 0].set_title('Bootstrap Distribution of Model
155     Coefficients')
156     axes[1, 0].tick_params(axis='x', rotation=45)
157     axes[1, 0].grid(True, alpha=0.3)
158
159     # Add original coefficients as red dots
160     for i, coef in enumerate(model.coef_):
161         axes[1, 0].scatter(i+1, coef, color='red', s=50, zorder=5)
162
163     # Coefficient confidence intervals
164     coef_stats = []
165     for i in range(coefficients.shape[1]):
166         coef_mean = np.mean(coefficients[:, i])
167         coef_std = np.std(coefficients[:, i])
168         coef_ci_lower, coef_ci_upper =
169         calculate_confidence_intervals(coefficients[:, i])
170
171         coef_stats.append({
172             'feature': feature_names[i],
173             'original': model.coef_[i],
174             'bootstrap_mean': coef_mean,
175             'bootstrap_std': coef_std,
176             'ci_lower': coef_ci_lower,
177             'ci_upper': coef_ci_upper
178         })
179
180     # Coefficient comparison plot
181     x_pos = np.arange(len(feature_names))
182     bootstrap_means = [stat['bootstrap_mean'] for stat in
183     coef_stats]

```

```

    coef_stats]
176     bootstrap_stds = [stat['bootstrap_std'] for stat in coef_stats]
177     original_coefs = [stat['original'] for stat in coef_stats]
178
179     axes[1, 1].errorbar(x_pos, bootstrap_means,
180                         yerr=bootstrap_stds,
181                         fmt='o', capsized=5, label='Bootstrap Mean ±
182                         Std')
183     axes[1, 1].scatter(x_pos, original_coefs, color='red', s=50,
184                         label='Original Coefficients', zorder=5)
185     axes[1, 1].set_xlabel('Features')
186     axes[1, 1].set_ylabel('Coefficient Value')
187     axes[1, 1].set_title('Coefficient Stability Analysis')
188     axes[1, 1].set_xticks(x_pos)
189     axes[1, 1].set_xticklabels(feature_names, rotation=45)
190     axes[1, 1].legend()
191     axes[1, 1].grid(True, alpha=0.3)
192
193     # Bootstrap vs Cross-Validation comparison
194     print(f"\nBOOTSTRAP vs CROSS-VALIDATION:")
195     print("-" * 35)
196
197     # Perform 5-fold CV for comparison
198     from sklearn.model_selection import cross_val_score
199     cv_r2_scores = cross_val_score(model, X, y, cv=5, scoring='r2')
200     cv_mse_scores = -cross_val_score(model, X, y, cv=5,
201                                     scoring='neg_mean_squared_error')
202
203     print(f"Cross-Validation R² Mean: {np.mean(cv_r2_scores):.4f} ±
204           {np.std(cv_r2_scores):.4f}")
205     print(f"Bootstrap R² Mean:           {np.mean(r2_scores):.4f} ±
206           {np.std(r2_scores):.4f}")
207     print()
208     print(f"Cross-Validation MSE Mean: {np.mean(cv_mse_scores):.2f} ±
209           {np.std(cv_mse_scores):.2f}")
210     print(f"Bootstrap MSE Mean:         {np.mean(mse_scores):.2f} ±
211           {np.std(mse_scores):.2f}")
212
213     # Comparison visualization
214     methods = ['Cross-Validation', 'Bootstrap']
215     r2_means_comp = [np.mean(cv_r2_scores), np.mean(r2_scores)]
216     r2_stds_comp = [np.std(cv_r2_scores), np.std(r2_scores)]
217
218     bars = axes[2, 0].bar(methods, r2_means_comp,
219                           yerr=r2_stds_comp,
220                           capsized=5, alpha=0.7, color=['blue',
221                                         'green'])
222     axes[2, 0].set_ylabel('R² Score')
223     axes[2, 0].set_title('Bootstrap vs Cross-Validation: R²
224                           Comparison')

```

```
215     axes[2, 0].grid(True, alpha=0.3)
216
217     # Add value labels
218     for bar, mean, std in zip(bars, r2_means_comp, r2_stds_comp):
219         axes[2, 0].text(bar.get_x() + bar.get_width()/2,
220                         bar.get_height() + std + 0.01,
221                         f'{mean:.3f}±{std:.3f}', ha='center',
222                         va='bottom', fontweight='bold')
223
224
225     applications_text = """
226 BOOTSTRAP APPLICATIONS & BEST PRACTICES:
227 ✓ WHEN TO USE BOOTSTRAP:
228     • Estimate confidence intervals for any metric
229     • Assess model stability and robustness
230     • Compare models statistically
231     • Understand parameter uncertainty
232     • Small datasets where CV is unstable
233
234 ✓ ADVANTAGES:
235     • Works with any statistic or metric
236     • No distributional assumptions
237     • Provides full sampling distribution
238     • Easy to implement and interpret
239
240 ✓ LIMITATIONS:
241     • Assumes sample represents population
242     • Can be computationally expensive
243     • May not work well with very small samples
244     • Doesn't account for model selection bias
245
246 ✓ BEST PRACTICES:
247     • Use  $B \geq 1000$  bootstrap samples
248     • Report confidence intervals
249     • Check for stability across runs
250     • Combine with other validation methods
251     • Consider stratified bootstrap for classification
252
253 ✓ BUSINESS VALUE:
254     • Quantify uncertainty in predictions
255     • Make risk-informed decisions
256     • Communicate model reliability
257     • Support regulatory requirements
258     """
259
260
261     axes[2, 1].text(0.05, 0.95, applications_text,
262                     transform=axes[2, 1].transAxes,
263                     fontsize=9, verticalalignment='top',
264                     fontfamily='monospace',
265                     bbox=dict(boxstyle='round',
266                               facecolor='lightyellow', alpha=0.8))
267
```

```
265     plt.tight_layout()
266     plt.show()
267
268     return bootstrap_results, coef_stats
269 bootstrap_results, coefficient_stats =
    bootstrap_sampling_comprehensive()
```
