

Data Analysis with Python

Cheat Sheet: Model Evaluation and Refinement

Process	Description	Code Example
Splitting data for training and testing	The process involves first separating the target attribute from the rest of the data. Treat the target attribute as the output and the rest of the data as input. Now split the input and output datasets into training and testing subsets.	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 15. 15 16. 16 17. 17 18. 18 19. 19 20. 20 21. 21 22. 22 23. 23 24. 24 25. 25 26. 26 27. 27 28. 28 29. 29 30. 30 31. 31 32. 32 33. 33 34. 34 35. 35 36. 36 37. 37 38. 38 39. 39 40. 40 41. 41 42. 42 43. 43 44. 44 45. 45 46. 46 47. 47 48. 48 49. 49 50. 50 51. 51 52. 52 53. 53 54. 54 55. 55 56. 56 57. 57 58. 58 59. 59 60. 60 61. 61 62. 62 63. 63 64. 64 65. 65 66. 66 67. 67 68. 68 69. 69 70. 70 71. 71 72. 72 73. 73 74. 74 75. 75 76. 76 77. 77 78. 78 79. 79 80. 80 81. 81 82. 82 83. 83 84. 84 85. 85 86. 86 87. 87 88. 88 89. 89 90. 90 91. 91 92. 92 93. 93 94. 94 95. 95 96. 96 97. 97 98. 98 99. 99 100. 100</pre> <pre>1. from sklearn.model_selection import train_test_split 2. x_data = df[['target_attribute']] 3. y_data = df[['target_attribute']] 4. x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.3, random_state=0) 5. print(x_train.shape, y_train.shape, x_test.shape, y_test.shape)</pre> <div>Copy</div>
Cross validation score	Without sufficient data, you go for cross validation, which involves creating different subsets of training and testing data multiple times and evaluating performance across all of them using the $R^2$ value.	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 15. 15 16. 16 17. 17 18. 18 19. 19 20. 20 21. 21 22. 22 23. 23 24. 24 25. 25 26. 26 27. 27 28. 28 29. 29 30. 30 31. 31 32. 32 33. 33 34. 34 35. 35 36. 36 37. 37 38. 38 39. 39 40. 40 41. 41 42. 42 43. 43 44. 44 45. 45 46. 46 47. 47 48. 48 49. 49 50. 50 51. 51 52. 52 53. 53 54. 54 55. 55 56. 56 57. 57 58. 58 59. 59 60. 60 61. 61 62. 62 63. 63 64. 64 65. 65 66. 66 67. 67 68. 68 69. 69 70. 70 71. 71 72. 72 73. 73 74. 74 75. 75 76. 76 77. 77 78. 78 79. 79 80. 80 81. 81 82. 82 83. 83 84. 84 85. 85 86. 86 87. 87 88. 88 89. 89 90. 90 91. 91 92. 92 93. 93 94. 94 95. 95 96. 96 97. 97 98. 98 99. 99 100. 100</pre> <pre>1. from sklearn.model_selection import cross_val_score 2. from sklearn.linear_model import LinearRegression 3. x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.3, random_state=0) 4. # Create a linear regression model 5. model = LinearRegression() 6. # Fit the model 7. model.fit(x_train, y_train) 8. # Predict on the test set 9. y_pred = model.predict(x_test) 10. # Calculate the cross-validation score 11. score = cross_val_score(model, x_data, y_data, cv=5) 12. print(score)</pre> <div>Copy</div>
Cross validation prediction	Use a cross validated model to create prediction of the output.	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 15. 15 16. 16 17. 17 18. 18 19. 19 20. 20 21. 21 22. 22 23. 23 24. 24 25. 25 26. 26 27. 27 28. 28 29. 29 30. 30 31. 31 32. 32 33. 33 34. 34 35. 35 36. 36 37. 37 38. 38 39. 39 40. 40 41. 41 42. 42 43. 43 44. 44 45. 45 46. 46 47. 47 48. 48 49. 49 50. 50 51. 51 52. 52 53. 53 54. 54 55. 55 56. 56 57. 57 58. 58 59. 59 60. 60 61. 61 62. 62 63. 63 64. 64 65. 65 66. 66 67. 67 68. 68 69. 69 70. 70 71. 71 72. 72 73. 73 74. 74 75. 75 76. 76 77. 77 78. 78 79. 79 80. 80 81. 81 82. 82 83. 83 84. 84 85. 85 86. 86 87. 87 88. 88 89. 89 90. 90 91. 91 92. 92 93. 93 94. 94 95. 95 96. 96 97. 97 98. 98 99. 99 100. 100</pre> <pre>1. from sklearn.model_selection import cross_val_predict 2. from sklearn.linear_model import LinearRegression 3. x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.3, random_state=0) 4. # Create a linear regression model 5. model = LinearRegression() 6. # Fit the model 7. model.fit(x_train, y_train) 8. # Predict on the test set 9. y_pred = model.predict(x_test) 10. # Calculate the cross-validation score 11. score = cross_val_score(model, x_data, y_data, cv=5) 12. print(score)</pre> <div>Copy</div>
Ridge Regression and Prediction	To create a better fitting polynomial regression model, like one that avoids overfitting to the training data, we use the Ridge regression model with a parameter alpha that is used to modify the effect of higher-order parameters on the model prediction.	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 15. 15 16. 16 17. 17 18. 18 19. 19 20. 20 21. 21 22. 22 23. 23 24. 24 25. 25 26. 26 27. 27 28. 28 29. 29 30. 30 31. 31 32. 32 33. 33 34. 34 35. 35 36. 36 37. 37 38. 38 39. 39 40. 40 41. 41 42. 42 43. 43 44. 44 45. 45 46. 46 47. 47 48. 48 49. 49 50. 50 51. 51 52. 52 53. 53 54. 54 55. 55 56. 56 57. 57 58. 58 59. 59 60. 60 61. 61 62. 62 63. 63 64. 64 65. 65 66. 66 67. 67 68. 68 69. 69 70. 70 71. 71 72. 72 73. 73 74. 74 75. 75 76. 76 77. 77 78. 78 79. 79 80. 80 81. 81 82. 82 83. 83 84. 84 85. 85 86. 86 87. 87 88. 88 89. 89 90. 90 91. 91 92. 92 93. 93 94. 94 95. 95 96. 96 97. 97 98. 98 99. 99 100. 100</pre> <pre>1. from sklearn.linear_model import Ridge 2. from sklearn.preprocessing import PolynomialFeatures 3. x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.3, random_state=0) 4. # Create a polynomial regression model 5. model = Ridge(alpha=1.0) 6. # Fit the model 7. model.fit(x_train, y_train) 8. # Predict on the test set 9. y_pred = model.predict(x_test) 10. # Calculate the cross-validation score 11. score = cross_val_score(model, x_data, y_data, cv=5) 12. print(score)</pre> <div>Copy</div>
Grid Search	Use Grid Search to find the correct alpha value for which the Ridge regression model gives the best performance. It further uses cross-validation to create a more refined model.	<pre>1. 1 2. 2 3. 3 4. 4 5. 5 6. 6 7. 7 8. 8 9. 9 10. 10 11. 11 12. 12 13. 13 14. 14 15. 15 16. 16 17. 17 18. 18 19. 19 20. 20 21. 21 22. 22 23. 23 24. 24 25. 25 26. 26 27. 27 28. 28 29. 29 30. 30 31. 31 32. 32 33. 33 34. 34 35. 35 36. 36 37. 37 38. 38 39. 39 40. 40 41. 41 42. 42 43. 43 44. 44 45. 45 46. 46 47. 47 48. 48 49. 49 50. 50 51. 51 52. 52 53. 53 54. 54 55. 55 56. 56 57. 57 58. 58 59. 59 60. 60 61. 61 62. 62 63. 63 64. 64 65. 65 66. 66 67. 67 68. 68 69. 69 70. 70 71. 71 72. 72 73. 73 74. 74 75. 75 76. 76 77. 77 78. 78 79. 79 80. 80 81. 81 82. 82 83. 83 84. 84 85. 85 86. 86 87. 87 88. 88 89. 89 90. 90 91. 91 92. 92 93. 93 94. 94 95. 95 96. 96 97. 97 98. 98 99. 99 100. 100</pre> <pre>1. from sklearn.linear_model import Ridge 2. from sklearn.preprocessing import PolynomialFeatures 3. from sklearn.model_selection import GridSearchCV 4. x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=0.3, random_state=0) 5. # Create a polynomial regression model 6. model = Ridge(alpha=1.0) 7. # Create a grid of parameters to search 8. param_grid = {'alpha': [0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000, 100000, 1000000]} 9. # Create a GridSearchCV object 10. grid_search = GridSearchCV(model, param_grid, cv=5, scoring='r2') 11. # Fit the model 12. grid_search.fit(x_train, y_train) 13. # Get the best model 14. best_model = grid_search.best_estimator_ 15. # Predict on the test set 16. y_pred = best_model.predict(x_test) 17. # Calculate the cross-validation score 18. score = cross_val_score(best_model, x_data, y_data, cv=5) 19. print(score)</pre> <div>Copy</div>

