

Data Analysis with Python

Cheat Sheet: Data Wrangling

Package/Method	Description	Code Example
Replace missing data with frequency	Replace the missing values of the data set attribute with the mode common occurring entry in the column.	<pre>1. 1 2. 2 3. mode_frequency = df[attribute_name].value_counts().idxmax() 4. df[attribute_name].replace(np.nan, mode_frequency, inplace=True) 5. df[attribute_name].replace(np.nan, mode_frequency, inplace=True)</pre> <div>Copy</div>
Replace missing data with mean	Replace the missing values of the data set attribute with the mean of all the entries in the column.	<pre>1. 1 2. 2 3. 3 4. mean_frequency = df[attribute_name].mean() 5. df[attribute_name].replace(np.nan, mean_frequency, inplace=True)</pre> <div>Copy</div>
Fix the data types	Fix the data types of the columns in the dataframe.	<pre>1. 1 2. 2 3. 3 4. df[attribute_name].astype('float64') 5. df[attribute_name].astype('int64') 6. df[attribute_name].astype('object')</pre> <div>Copy</div>
Data Normalization	Normalize the data in a column such that the values are restricted between 0 and 1.	<pre>1. 1 2. 2 3. 3 4. df[attribute_name] = df[attribute_name].min() 5. df[attribute_name] = df[attribute_name].max()</pre> <div>Copy</div>
Binning	Create bins of data for better analysis and visualization.	<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> <div>Copy</div>
Change column name	Change the label name of a dataframe column.	<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> <div>Copy</div>
Indicator Variables	Create indicator variables for categorical data.	<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> <div>Copy</div>

