

Fall 2023 Data C100/C200 Midterm Reference Sheet

Pandas

Suppose `df` is a `DataFrame`; `s` is a `Series`. `import pandas as pd`

Function	Description
<code>df[col]</code>	Returns the column labeled <code>col</code> from <code>df</code> as a <code>Series</code> .
<code>df[[col1, col2]]</code>	Returns a <code>DataFrame</code> containing the columns labeled <code>col1</code> and <code>col2</code> .
<code>s.loc[rows] / df.loc[rows, cols]</code>	Returns a <code>Series/DataFrame</code> with rows (and columns) selected by their index values.
<code>s.iloc[rows] / df.iloc[rows, cols]</code>	Returns a <code>Series/DataFrame</code> with rows (and columns) selected by their positions.
<code>s.isnull() / df.isnull()</code>	Returns boolean <code>Series/DataFrame</code> identifying missing values
<code>s.fillna(value) / df.fillna(value)</code>	Returns a <code>Series/DataFrame</code> where missing values are replaced by <code>value</code>
<code>s.isin(values) / df.isin(values)</code>	Returns a <code>Series/DataFrame</code> of booleans indicating if each element is in <code>values</code> .
<code>df.drop(labels, axis)</code>	Returns a <code>DataFrame</code> without the rows or columns named <code>labels</code> along <code>axis</code> (either 0 or 1)
<code>df.rename(index=None, columns=None)</code>	Returns a <code>DataFrame</code> with renamed columns from a dictionary <code>index</code> and/or <code>columns</code>
<code>df.sort_values(by, ascending=True)</code>	Returns a <code>DataFrame</code> where rows are sorted by the values in columns <code>by</code>
<code>s.sort_values(ascending=True)</code>	Returns a sorted <code>Series</code> .
<code>s.unique()</code>	Returns a <code>NumPy</code> array of the unique values
<code>s.value_counts()</code>	Returns the number of times each unique value appears in a <code>Series</code>
<code>pd.merge(left, right, how='inner', on='a')</code>	Returns a <code>DataFrame</code> joining <code>DataFrames</code> <code>left</code> and <code>right</code> on the column labeled <code>a</code> ; the join is of type <code>inner</code>
<code>left.merge(right, left_on=col1, right_on=col2)</code>	Returns a <code>DataFrame</code> joining <code>DataFrames</code> <code>left</code> and <code>right</code> on columns labeled <code>col1</code> and <code>col2</code> .
<code>df.pivot_table(index, columns, values=None, aggfunc='mean', fill_value=None)</code>	Returns a <code>DataFrame</code> pivot table where columns are unique values from <code>columns</code> (column name or list), and rows are unique values from <code>index</code> (column name or list); cells are collected <code>values</code> using <code>aggfunc</code> . If <code>values</code> is not provided, cells are collected for each remaining column with multi-level column indexing. If a <code>fill_value</code> is provided, any <code>NaN</code> values will be replaced with that <code>fill_value</code> .
<code>df.set_index(col)</code>	Returns a <code>DataFrame</code> that uses the values in the column labeled <code>col</code> as the row index.
<code>df.reset_index()</code>	Returns a <code>DataFrame</code> that has row index 0, 1, etc., and adds the current index as a column.

Let `grouped = df.groupby(by)` where `by` can be a column label or a list of labels.

Function	Description
<code>grouped.count()</code>	Return a <code>DataFrame</code> containing the size of each group, excluding missing values
<code>grouped.size()</code>	Return a <code>Series</code> containing size of each group, including missing values
<code>grouped.mean()/grouped.min()/grouped.max()</code>	Return a <code>Series/DataFrame</code> containing mean/min/max of each group for each column, excluding missing values
<code>grouped.filter(f)</code> <code>grouped.agg(f)</code>	Filters or aggregates using the given function <code>f</code>

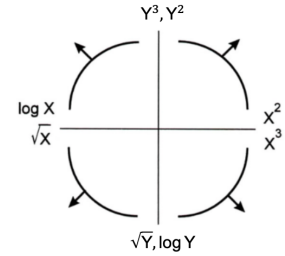
Function	Description
<code>s.str.len()</code>	Returns a <code>Series</code> containing length of each string
<code>s.str[a:b]</code>	Returns a <code>Series</code> where each element is a slice of the corresponding string indexed from <code>a</code> (inclusive, optional) to <code>b</code> (non-inclusive, optional)
<code>s.str.lower()/s.str.upper()</code>	Returns a <code>Series</code> of lowercase/uppercase versions of each string
<code>s.str.replace(pat, repl)</code>	Returns a <code>Series</code> that replaces occurrences of substrings matching the regex <code>pat</code> with string <code>repl</code>
<code>s.str.contains(pat)</code>	Returns a boolean <code>Series</code> indicating if a substring matching the regex <code>pat</code> is contained in each string
<code>s.str.extract(pat)</code>	Returns a <code>Series</code> of the first subsequence of each string that matches the regex <code>pat</code> . If <code>pat</code> contains one group, then only the substring matching the group is extracted

Visualization

Matplotlib: x and y are sequences of values. `import matplotlib.pyplot as plt`

Tukey-Mosteller Bulge Diagram.

Function	Description
<code>plt.plot(x, y)</code>	Creates a line plot of x against y
<code>plt.scatter(x, y)</code>	Creates a scatter plot of x against y
<code>plt.hist(x, bins=None)</code>	Creates a histogram of x; bins can be an integer or a sequence
<code>plt.bar(x, height)</code>	Creates a bar plot of categories x and corresponding heights height



Seaborn: x and y are column names in a DataFrame data. `import seaborn as sns`

Function	Description
<code>sns.countplot(data, x)</code>	Create a barplot of value counts of variable x from data
<code>sns.histplot(data, x, stat='count', kde=False)</code> <code>sns.displot(data, x, kind='hist', rug=False, kde=False)</code>	Creates a histogram of x from data, where bin statistics stat is one of 'count', 'frequency', 'probability', 'percent', and 'density'; optionally overlay a kernel density estimator. <code>displot</code> is similar but can optionally overlay a rug plot and/or a KDE plot
<code>sns.boxplot(data, x=None, y)</code> <code>sns.violinplot(data, x=None, y)</code>	Create a boxplot of y, optionally factoring by categorical x, from data. <code>violinplot</code> is similar but also draws a kernel density estimator of y
<code>sns.rugplot(data, x)</code>	Adds a rug plot on the x-axis of variable x from data
<code>sns.scatterplot(data, x, y)</code>	Create a scatterplot of x versus y from data
<code>sns.lmplot(x, y, data, fit_reg=True)</code>	Create a scatterplot of x versus y from data, and by default overlay a least-squares regression line
<code>sns.jointplot(x, y, data, kind)</code>	Combine a bivariate scatterplot of x versus y from data, with univariate density plots of each variable overlaid on the axes; kind determines the visualization type for the distribution plot, can be scatter, kde or hist

Regular Expressions

Operator	Description	Operator	Description
<code>.</code>	Matches any character except <code>\n</code>	<code>*</code>	Matches preceding character/group zero or more times
<code>\</code>	Escapes metacharacters	<code>?</code>	Matches preceding character/group zero or one times
<code> </code>	Matches expression on either side of expression; has lowest priority of any operator	<code>+</code>	Matches preceding character/group one or more times
<code>\d, \w, \s</code>	Predefined character group of digits (0-9), alphanumerics (a-z, A-Z, 0-9, and underscore), or whitespace, respectively	<code>^, \$</code>	Matches the beginning and end of the line, respectively
<code>\D, \W, \S</code>	Inverse sets of <code>\d, \w, \s</code> , respectively	<code>()</code>	Capturing group used to create a sub-expression
<code>{m}</code>	Matches preceding character/group exactly m times	<code>[]</code>	Character class used to match any of the specified characters or range (e.g. <code>[abcde]</code> is equivalent to <code>[a-e]</code>)
<code>{m, n}</code>	Matches preceding character/group at least m times and at most n times. If either m or n are omitted, set lower/upper bounds to 0 and ∞ , respectively	<code>[^]</code>	Invert character class; e.g. <code>[^a-c]</code> matches all characters except a, b, c

Modified lecture example for capture groups:

```
import re
lines = '169.237.46.168 -- [26/Jan/2014:10:47:58 -0800] "GET ... HTTP/1.1"'
re.findall(r'[\d+\/(\w+)\]\d+:\d+:\d+:\d+ .+\]', line) # returns ['Jan']
```

Function	Description
<code>re.match(pattern, string)</code>	Returns a match if zero or more characters at beginning of string matches pattern, else None
<code>re.search(pattern, string)</code>	Returns a match if zero or more characters anywhere in string matches pattern, else None
<code>re.findall(pattern, string)</code>	Returns a list of all non-overlapping matches of pattern in string (if none, returns empty list)
<code>re.sub(pattern, repl, string)</code>	Returns string that replaces all occurrences of pattern with repl

Modeling

Concept	Formula	Concept	Formula
Variance, σ_x^2	$\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$	Correlation r	$r = \frac{1}{n} \sum_{i=1}^n \frac{x_i - \bar{x}}{\sigma_x} \frac{y_i - \bar{y}}{\sigma_y}$
L_1 loss	$L_1(y, \hat{y}) = y - \hat{y} $	Linear regression estimate of y	$\hat{y} = \theta_0 + \theta_1 x$
L_2 loss	$L_2(y, \hat{y}) = (y - \hat{y})^2$	Least squares linear regression	$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 \bar{x} \quad \hat{\theta}_1 = r \frac{\sigma_y}{\sigma_x}$
Empirical risk with loss L	$R(\theta) = \frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i)$		

Ordinary Least Squares

Multiple Linear Regression Model: $\hat{\mathbb{Y}} = \mathbb{X}\theta$ with design matrix \mathbb{X} , response vector \mathbb{Y} , and predicted vector $\hat{\mathbb{Y}}$. If there are p features plus a bias/intercept, then the vector of parameters $\theta = [\theta_0, \theta_1, \dots, \theta_p]^T \in \mathbb{R}^{p+1}$. The vector of estimates $\hat{\theta}$ is obtained from fitting the model to the sample (\mathbb{X}, \mathbb{Y}) .

Concept	Formula	Concept	Formula
Mean squared error	$R(\theta) = \frac{1}{n} \ \mathbb{Y} - \mathbb{X}\theta\ _2^2$	Normal equation	$\mathbb{X}^T \mathbb{X} \hat{\theta} = \mathbb{X}^T \mathbb{Y}$
		Least squares estimate, if \mathbb{X} is full rank	$\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$
Residual vector, e	$e = \mathbb{Y} - \hat{\mathbb{Y}}$	Multiple R^2 (coefficient of determination)	$R^2 = \frac{\text{variance of fitted values}}{\text{variance of } y}$