Spring 2023 Data C100/C200 Midterm Reference Sheet

Pandas

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

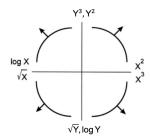
Function	Description		
df[col]	Returns the column labeled col from df as a Series.		
df[[col1, col2]]	Returns a DataFrame containing the columns labeled col1 and col2.		
<pre>s.loc[rows] / df.loc[rows, cols]</pre>	Returns a Series/DataFrame with rows (and columns) selected by their index values.		
<pre>s.iloc[rows] / df.iloc[rows, cols]</pre>	Returns a Series/DataFrame with rows (and columns) selected by their positions.		
<pre>s.isnull() / df.isnull()</pre>	Returns boolean Series/DataFrame identifying missing values		
s.fillna(value) / df.fillna(value)	Returns a Series/DataFrame where missing values are replaced by value		
df.drop(labels, axis)	Returns a DataFrame without the rows or columns named labels along axis (either 0 or 1)		
<pre>df.rename(index=None, columns=None)</pre>	Returns a DataFrame with renamed columns from a dictionary index and/or columns		
df.sort_values(by, ascending=True)	Returns a DataFrame where rows are sorted by the values in columns by		
s.sort_values(ascending=True)	Returns a sorted Series.		
s.unique()	Returns a NumPy array of the unique values		
s.value_counts()	Returns the number of times each unique value appears in a Series		
<pre>pd.merge(left, right, how='inner', on='a')</pre>	Returns a DataFrame joining DataFrames left and right on the column labeled a; the join is of type inner		
<pre>left.merge(right, left_on=col1, right_on=col2)</pre>	Returns a DataFrame joining DataFrames left and right on columns labeled col1 and col2.		
<pre>df.pivot_table(index, columns, values=None, aggfunc='mean')</pre>	Returns a DataFrame pivot table where columns are unique values from columns (column name or list), and rows are unique values from index (column name or list); cells are collected values using aggfunc. If values is not provided, cells are collected for each remaining column with multi-level column indexing.		
df.set_index(col)	Returns a DataFrame that uses the values in the column labeled col as the row index.		
<pre>df.reset_index() Let grouped = df.groupby(by) where by can</pre>	Returns a DataFrame that has row index 0, 1, etc., and adds the current index as a column. be a column label or a list of labels.		
Function	Description		
grouped.count()	Return a Series containing the size of each group, excluding missing values		
grouped.size()	Return a Series containing size of each group, including missing values		
<pre>grouped.mean()/grouped.min()/grouped.max()</pre>	Return a Series/DataFrame containing mean/min/max of each group for each column, excluding missing values		
<pre>grouped.filter(f) grouped.agg(f)</pre>	Filters or aggregates using the given function f		
Function	Description		
s.str.len()	Returns a Series containing length of each string		
<pre>s.str.lower()/s.str.upper()</pre>	Returns a Series containing lowercase/uppercase version of each string		
s.str.replace(pat, repl)	Returns a Series after replacing occurences of substrings matching regular expression pat with string repl		
s.str.contains(pat)	Returns a boolean Series indicating whether a substring matching the regular expression pat is contained in each string		
<pre>s.str.extract(pat)</pre>	Returns a Series of the first subsequence of each string that matches the regular expression pat. If pat contains one group, then only the substring matching the group is extracted		

Visualization

Matplotlib: x and y are sequences of values.

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

Tukey-Mosteller Bulge Diagram.



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

Function	Description
<pre>sns.countplot(data, x)</pre>	Create a barplot of value counts of variable x from data
<pre>sns.histplot(data, x, kde=False) sns.displot(x, data, rug=True, kde=True)</pre>	Creates a histogram of x from data; optionally overlay a kernel density estimator. displot is similar but can optionally overlay a rug plot.
<pre>sns.boxplot(data, x=None, y) sns.violinplot(data, x=None, y)</pre>	Create a boxplot of y, optionally factoring by categorical x, from data. violinplot is similar but also draws a kernel density estimator of y.
<pre>sns.scatterplot(data, x, y)</pre>	Create a scatterplot of x versus y from data
<pre>sns.lmplot(x, y, data, fit_reg=True)</pre>	Create a scatterplot of \boldsymbol{x} versus \boldsymbol{y} from data, and by default overlay a least-squares regression line
sns.jointplot(x, y, data, kind)	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
	Matches any character except \n		*	Matches preceding character/group zero or more times
\\	Escapes metacharacters		?	Matches preceding character/group zero or one times
I	Matches expression on either side of expression; has lowest priority of any operator		+	Matches preceding character/group one or more times
\d, \w, \s	Predefined character group of digits alphanumerics (a-z, A-Z, 0-9, and underscore), or whitespace, respecti	. , , ,	^, \$	Matches the beginning and end of the line, respectively
\D, \W, \S	Inverse sets of \d, \w, \s, respectivel	ly	()	Capturing group used to create a sub- expression
{m}	Matches preceding character/group exactly m times		[]	Character class used to match any of the specified characters or range (e.g. [abcde] is equivalent to [a-e])
{m, n}	Matches preceding character/group least m times and at most n times if e or n are omitted, set lower/upper bout to 0 and ∞ , respectively	ither m	[^]	Invert character class; e.g. [^a-c] matches all characters except a, b, c
Function	De	Description		
re.match(pattern, string)		Returns a match if zero or more characters at beginning of string matches pattern, else None		
re.search(pattern, string)		Returns a match if zero or more characters anywhere in string matches pattern , else None		
re.findall(pattern, string		Returns a list of all non-overlapping matches of pattern in string (if none, returns empty list)		

Function Description

re.sub(pattern, repl, string)

Returns string after replacing all occurrences of pattern with repl

Modified lecture example for a single capturing group:

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

Modeling

Concept	Formula	Concept	Formula
Variance, σ_x^2	$\frac{1}{n}\sum_{i=1}^n(x_i-\bar{x})^2$	Correlation r	$r=rac{1}{n}\sum_{i=1}^{n}rac{x_{i}-ar{x}}{\sigma_{x}}rac{y_{i}-ar{y}}{\sigma_{y}}$
L_1 loss	$L_1(y,\hat{y}) = \mid y - \hat{y} \mid$	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	$\hat{ heta}_0 = ar{y} - \hat{ heta}_1 ar{x} \qquad \hat{ heta}_1 = r$
mpirical risk with loss ${\cal L}$	$\frac{1}{n}$		

$$R(heta) = rac{1}{n} \sum_{i=1}^n L(y_i, \hat{y_i})$$

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	$R(heta) = rac{1}{n} \ \ \mathbb{Y} - \mathbb{X} heta\ \ _2^2$	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 $\ R^2$ determination)	$z = rac{ ext{variance of fitted values}}{ ext{variance of } y}$