General Tips

 Stop wasting time explicitly importing all the data science libraries you need to use in your dev environment. Instead, use Pyforest to lazily import them for you only when you need them:

pip install pyforest

Tired of writing for loops to join lists? Use the zip function instead:

Initial data	a = ("John", "Charles", "Mike") b = ("For", "Against", "For")		
Zip function	x = zip(a, b)		
Result	(('John', For), ('Charles', 'Against'), ('Mike', 'For'))		

Pandas Tips

Shortcut your initial data analysis with the pandas' profile function, which generates a detailed report of your data (missing values, variables, counts, etc) in just one line of code. For example:

df = pd.read_csv('somedata.csv')

df.profile_report()

 Speed up pandas operations with pandarallel. For example, instead of df.progress_apply() use df.parallel_apply() to run the process in parallel.

Need to unstack a table? Use pandas, which can convert one level of an index into the columns of your data frame. For example:

Initial data (trimmed)	state AK AZ CA	<pre>email_provider aol.com hotmail.cm cox.net kitty.com deleo.com aol.com aol.com cox.net nikolozakes.org parvis.com gmail.com cox.net aol.com</pre>			
Unstack code	<pre>clients.groupby('state')['email_provider'].value_counts ().unstack().fillna(0)</pre>				
Result (trimmed)	email_provider state ak ar az Ca Ca Ca Ct Ct fl ga	angalich.com 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.	ankeny.org 0.0 0.0 1.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0	aol.com 2.0 2.0 9.0 0.0 1.0 0.0 4.0 1.0	



from ActiveState

Visualization Tips

Tired of rendering static plots in Jupyter with %matplotlib ? Try importing %matplotlib notebook which gives you interactive plots you can resize and zoom in on.

• Want to more easily spot patterns in tabulated data? Create a heatmap using Seaborn's gradient capabilities. For example:

hm = sns.light.palette('green', as_cmap = True)
style = df.style.background_gradient(cmap = hm)

	A	В	С	D	E
0	-1.264051	1.52791	-0.970711	0.47056	-0.100697
1	0.303793	-1.72596	1.5851	0.13297	-1.10686
2	1.57823	0.10798	-0.764048	-0.775189	1.38385
3	0.760385	-0.285647	0.538367	-2.0839	0.937782

Jupyter Notebook Tips

Need to debug your code in Jupyter notebooks? Just type %debug to launch an interactive debugger that takes you back to the point where the exception happened. Press q to exit.

Computational costs matter. To check the running time of a block of code in a Jupyter notebook preface the code block with %%time

• Working with Python in a Jupyter Notebook, but wish you had access to R? Now you can run both of them together by typing pip install rpy2.

Scikit-Learn Tips

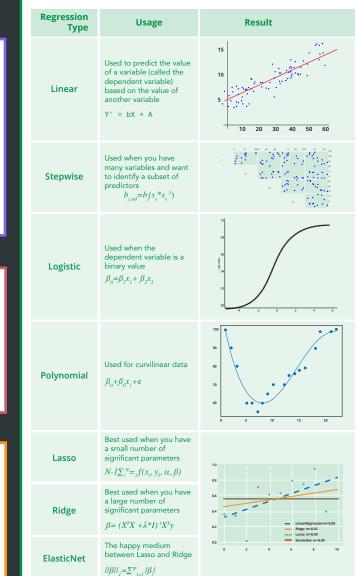
Always use the <u>stratify</u> parameter to ensure test and train sets are split into equal proportions for better prediction and reproducibility of results. For example:

test_x, train_x, test_y, train_y = tain_test_split (x, y, random_state = 59, stratify = y)

Missing values in your dataset? Don't settle for univariate methods to create the missing values when scikit-learn's multivariate input based on k-Nearest Neighbor can offer better accuracy:
 impute = KNNImputer (n neighbors=2)

Predictive Analysis Guide

And finally, a simple rule guide for when to apply which regression technique when doing predictive analysis:



ActivePython comes pre-compiled with the most popular machine learning packages. Download it from activestate.com

ActiveState