

# Exam 2

In [ ]:	<pre>#Instructions #Note that each question is for 1 point #Exam total of 40 points possible  #Add your name and ID below #Remember to save your work #Run the setup sections and answer the questions that follow</pre>																																																																																				
In [2]:	<pre># Student Name: Xingyu Chen # Student ID: I10291925</pre>																																																																																				
In [3]:	<pre>#Loading Matplotlib and Seaborn from matplotlib import pyplot as plt import pandas as pd import seaborn as sns import numpy as np from matplotlib import pyplot as plt sns.set_style('darkgrid')</pre>																																																																																				
In [4]:	<pre>#Run the following cell. #Use the following data set for Questions 1 - 5 below  import pandas as pd Cars93_data = pd.read_csv("https://vincentarelbundock.github.io/Rdatasets/csv/MASS/Cars93.csv", index_col=0).dx Cars93_data.head()</pre>																																																																																				
Out[4]:	<table><tr><th></th><th>Manufacturer</th><th>Model</th><th>Type</th><th>Min.Price</th><th>Price</th><th>Max.Price</th><th>MPG.city</th><th>MPG.highway</th><th>AirBags</th><th>DriveTrain</th><th>...</th><th>Passengers</th><th>Length</th></tr><tr><td>1</td><td>Acura</td><td>Integra</td><td>Small</td><td>12.9</td><td>15.9</td><td>18.8</td><td>25</td><td>31</td><td>None</td><td>Front</td><td>...</td><td>5</td><td>17</td></tr><tr><td>2</td><td>Acura</td><td>Legend</td><td>Midsize</td><td>29.2</td><td>33.9</td><td>38.7</td><td>18</td><td>25</td><td>Driver &amp; Passenger</td><td>Front</td><td>...</td><td>5</td><td>19</td></tr><tr><td>3</td><td>Audi</td><td>90</td><td>Compact</td><td>25.9</td><td>29.1</td><td>32.3</td><td>20</td><td>26</td><td>Driver only</td><td>Front</td><td>...</td><td>5</td><td>18</td></tr><tr><td>4</td><td>Audi</td><td>100</td><td>Midsize</td><td>30.8</td><td>37.7</td><td>44.6</td><td>19</td><td>26</td><td>Driver &amp; Passenger</td><td>Front</td><td>...</td><td>6</td><td>19</td></tr><tr><td>5</td><td>BMW</td><td>535i</td><td>Midsize</td><td>23.7</td><td>30.0</td><td>36.2</td><td>22</td><td>30</td><td>Driver only</td><td>Rear</td><td>...</td><td>4</td><td>18</td></tr></table> <p>5 Rows x 27 Columns</p>		Manufacturer	Model	Type	Min.Price	Price	Max.Price	MPG.city	MPG.highway	AirBags	DriveTrain	...	Passengers	Length	1	Acura	Integra	Small	12.9	15.9	18.8	25	31	None	Front	...	5	17	2	Acura	Legend	Midsize	29.2	33.9	38.7	18	25	Driver & Passenger	Front	...	5	19	3	Audi	90	Compact	25.9	29.1	32.3	20	26	Driver only	Front	...	5	18	4	Audi	100	Midsize	30.8	37.7	44.6	19	26	Driver & Passenger	Front	...	6	19	5	BMW	535i	Midsize	23.7	30.0	36.2	22	30	Driver only	Rear	...	4	18
	Manufacturer	Model	Type	Min.Price	Price	Max.Price	MPG.city	MPG.highway	AirBags	DriveTrain	...	Passengers	Length																																																																								
1	Acura	Integra	Small	12.9	15.9	18.8	25	31	None	Front	...	5	17																																																																								
2	Acura	Legend	Midsize	29.2	33.9	38.7	18	25	Driver & Passenger	Front	...	5	19																																																																								
3	Audi	90	Compact	25.9	29.1	32.3	20	26	Driver only	Front	...	5	18																																																																								
4	Audi	100	Midsize	30.8	37.7	44.6	19	26	Driver & Passenger	Front	...	6	19																																																																								
5	BMW	535i	Midsize	23.7	30.0	36.2	22	30	Driver only	Rear	...	4	18																																																																								
In [5]:	<pre>#Question 1: #Run this code to display the pairplot for features Price, MPG.city and Weight sns.pairplot(Cars93_data, hue='AirBags', vars=['Price', 'MPG.city', 'Weight'])  #Write Seaborn code in the next cell that shows that a non-linear model # is probably not appropriate to model the Price-Weight relationship</pre>	Out[5]:																																																																																			
Out[5]:																																																																																					
In [95]:	<pre>#Write your answer for Question 1 below: #Write Seaborn code in the next cell that shows that a non-linear model # is probably not appropriate to model the Price-Weight relationship sns.jointplot(x=Cars93_data.Weight, y=Cars93_data.Price, kind='resid') sns.pairplot(Cars93_data, vars=['Price', 'Weight'])</pre>	Out[95]:																																																																																			
Out[95]:																																																																																					
In [96]:	<pre>#Question 2: #Explain your answer to Question 1 below:  # The pairplot suggests (to some extent) that a non-linear model would not be appropriate for the Price-Weight # The residual plot for the same relationship appears to indicate that a linear model would be appropriate # So a non-linear model is probably not appropriate to model the Price-Weight relationship # A linear model is probably appropriate to model the Price-Weight relationship</pre>	In [15]:																																																																																			
In [15]:	<pre>#Question 3: #Write Seaborn code to build a 1-row Facet Grid of boxplots for the Price feature using Cars93_data #Breakout the chart by the AirBags feature # Use .map_dataframe() #Note you should have all lines of code in the same cell  #Write your answer below: g=sns.FacetGrid(Cars93_data, col='AirBags') g.map_dataframe(sns.histplot, x='Price') g.map_dataframe(sns.boxplot, x='Price')</pre>	Out[15]:																																																																																			
Out[15]:																																																																																					
In [16]:	<pre>#Question 4 #From your answer/chart in Question 3 above, which value of AirBags has # the lowest 50th Percentile for the feature Price?  #Write your answer below: # AirBags = None</pre>	In [31]:																																																																																			
In [31]:	<pre>#Question 5 #From your answer to Question 3 above, which value of AirBags has the largest interquartile range (IQR)?  #Write your answer below: g=sns.FacetGrid(Cars93_data, col='AirBags') g.set_axis_labels('Price', 'Weight') g.map_dataframe(sns.boxplot, x='Price') g.add_legend()  # AirBags = Driver &amp; Passenger</pre>	Out[31]:																																																																																			
Out[31]:																																																																																					
In [40]:	<pre>#Question 6 #Write Seaborn code to  #Write Seaborn code to build a 1-column Facet Grid of scatterplots using Cars93_data #Breakout the chart by the AirBags feature #Add a scatterplot to the column for the variable Price # Use .map_dataframe(). Put feature Price on the horizontal axis and Weight on the vertical axis #Note you should have all lines of code in the same cell #Use hue to breakout your chart with categorical variable DriveTrain #Label the horizontal axis Price #Label the vertical axis Weight #Remember to pass hue to the sns.FacetGrid line and not the # .map_dataframe() line #Also add legend with g.add_legend  #Write your answer below: g=sns.FacetGrid(Cars93_data, row='AirBags', hue = 'DriveTrain') g.map_dataframe(sns.scatterplot, x='Price', y = 'Weight') g.set_axis_labels('Price', 'Weight') g.add_legend()</pre>	Out[40]:																																																																																			
Out[40]:																																																																																					
In [41]:	<pre>#Question 7 # Using your chart from Question 6 above, when AirBags = None, the only dot for DriveTrain = Rear # is between what range of values for variable Weight?  #Write your answer below:  # around 3550 weight, highest value for variable weight, outlier</pre>	In [42]:																																																																																			
In [42]:	<pre>#Question 8 #Using your chart from Question 6 above, What can you infer about the Weight-Price relationship # when feature DriveTrain = 4WD and feature AirBags = Driver &amp; Passenger?  #Write your answer below:  # No green dot, there is no data when feature DriveTrain = 4WD and feature AirBags = Driver &amp; Passenger # so I can not infer about the Weight-Price relationship</pre>	In [43]:																																																																																			
In [43]:	<pre>#Question 9 #Using your chart from Question 6 above, what can you infer about the Weight-Price relationship # when DriveTrain = Front and AirBags = Driver only?  #Write your answer below:  # They have positive relationship, when price increasing, the weight increasing # meanwhile when price decreasing, the weight decreasing, when DriveTrain = Front and AirBags = Driver only</pre>	In [45]:																																																																																			
In [45]:	<pre>#Question 10 # From the chart in Question 1, what can you infer when AirBags = Driver only with respect to # feature Weight and feature MPG.city?  #Write your answer below:  # when AirBags = Driver only, they are more dense for feature Weight and feature MPG.city.</pre>	In [46]:																																																																																			
In [46]:	<pre>#Run the following cell. #Use the following data set for Questions 11 - 20 below  import pandas as pd Bacteria_data = pd.read_csv("https://vincentarelbundock.github.io/Rdatasets/csv/MASS/cabbages.csv", index_col=0) Bacteria_data.head()</pre>	Out[46]:																																																																																			
Out[46]:	<table><tr><th></th><th>Cult</th><th>Date</th><th>HeadWt</th><th>VitC</th></tr><tr><td>1</td><td>c39</td><td>d16</td><td>2.5</td><td>61</td></tr><tr><td>2</td><td>c39</td><td>d16</td><td>2.2</td><td>65</td></tr><tr><td>3</td><td>c39</td><td>d16</td><td>3.1</td><td>45</td></tr><tr><td>4</td><td>c39</td><td>d16</td><td>4.3</td><td>42</td></tr><tr><td>5</td><td>c39</td><td>d16</td><td>2.5</td><td>63</td></tr></table>		Cult	Date	HeadWt	VitC	1	c39	d16	2.5	61	2	c39	d16	2.2	65	3	c39	d16	3.1	45	4	c39	d16	4.3	42	5	c39	d16	2.5	63																																																						
	Cult	Date	HeadWt	VitC																																																																																	
1	c39	d16	2.5	61																																																																																	
2	c39	d16	2.2	65																																																																																	
3	c39	d16	3.1	45																																																																																	
4	c39	d16	4.3	42																																																																																	
5	c39	d16	2.5	63																																																																																	
In [47]:	<pre>#Question 11 #Write the Seaborn code to render a pairplot on the entire dataset above # Breakout the pairplot by the feature Cult #Write your answer below  sns.pairplot(Bacteria_data, hue='Cult')</pre>	Out[47]:																																																																																			
Out[47]:																																																																																					
In [51]:	<pre>#Question 12 #From your chart in Question 12 above, write Seaborn code that would help you confirm # whether or not a polynomial model would be appropriate to model the HeadWt-VitC relationship  #Write your answer below: sns.jointplot(x=Bacteria_data.VitC, y=Bacteria_data.HeadWt,               kind='reg',               joint_kws={'ci': 'None', 'order': 2})</pre>	Out[51]:																																																																																			
Out[51]:																																																																																					
In [52]:	<pre>#Question 13 #From your answer in Question 12 above, what does your chart suggest about whether or not to use a polynomial m #Write your answer below:  # a polynomial model would not be appropriate to model the HeadWt-VitC relationship</pre>	In [53]:																																																																																			
In [53]:	<pre>#Question 14 #What does your answer to Questions 11, 12 and 13 suggest about the HeadWt-VitC relationship?  #Write your answer below:  # they have negative linear relationship, when VitC increasing, the HeadWt decreasing # meanwhile when VitC decreasing, the HeadWt increasing.</pre>	In [57]:																																																																																			
In [57]:	<pre>#Question 15 #Write Seaborn code to draw a pairplot (using kernel density plot) for the entire data frame Bacteria_data #Write your answer below: sns.pairplot(Bacteria_data, kind = 'kde')</pre>	Out[57]:																																																																																			
Out[57]:																																																																																					
In [59]:	<pre>#Question 16 #Write Seaborn code to draw a pairplot for the entire Bacteria_data data frame #Breakout the plot by feature Date  #Write your answer below: sns.pairplot(Bacteria_data, hue = 'Date')</pre>	Out[59]:																																																																																			
Out[59]:																																																																																					
In [60]:	<pre>#Question 17 #From your chart in Question 16 above, what two values of feature Date can you conclude # are mostly responsible for the deviation of the shape of the VitC feature # from the benchmark normal distribution curve?  #Write your answer below: # d21 and d16</pre>	In [ ]:																																																																																			
In [ ]:	<pre>#Question 18 #From your chart in Question 16 above, What values of feature Date can you conclude # are mostly responsible for the deviation of the shape of the HeadWt feature # from the benchmark normal distribution curve?  #Write your answer below: # d21 and d20</pre>	In [61]:																																																																																			
In [61]:	<pre>#Write Seaborn code to draw a regression plot for the HeadWt-VitC relationship #Turn off the confidence interval  #Write your answer below: sns.jointplot(x=Bacteria_data.VitC, y=Bacteria_data.HeadWt,               kind='reg',               joint_kws={'ci': 'None'})</pre>	Out[61]:																																																																																			
Out[61]:																																																																																					
In [62]:	<pre>#Question 20 #From your chart in Question 19 above, can an increase in VitC be expected # to be associated with a decrease in HeadWt on average? (YES or NO)  #Write your answer below:  #YES</pre>	In [63]:																																																																																			
In [63]:	<pre>#Run this cell and use it to answer Questions 21 - 33 below  import pandas as pd import seaborn as sns import numpy as np from matplotlib import pyplot as plt import matplotlib Insurance_data = pd.read_csv("https://vincentarelbundock.github.io/Rdatasets/csv/MASS/Insurance.csv", index_col=0) Insurance_data.head()</pre>	Out[63]:																																																																																			
Out[63]:	<table><tr><th></th><th>District</th><th>Age</th><th>Holders</th><th>Claims</th></tr><tr><td>1</td><td>1</td><td>&lt;11</td><td>&lt;25</td><td>197</td><td>38</td></tr><tr><td>2</td><td>1</td><td>&lt;11</td><td>25-29</td><td>264</td><td>35</td></tr><tr><td>3</td><td>1</td><td>&lt;11</td><td>30-35</td><td>246</td><td>20</td></tr><tr><td>4</td><td>1</td><td>&lt;11</td><td>&gt;35</td><td>1680</td><td>156</td></tr><tr><td>5</td><td>1</td><td>1-151</td><td>&lt;25</td><td>284</td><td>63</td></tr></table>		District	Age	Holders	Claims	1	1	<11	<25	197	38	2	1	<11	25-29	264	35	3	1	<11	30-35	246	20	4	1	<11	>35	1680	156	5	1	1-151	<25	284	63																																																	
	District	Age	Holders	Claims																																																																																	
1	1	<11	<25	197	38																																																																																
2	1	<11	25-29	264	35																																																																																
3	1	<11	30-35	246	20																																																																																
4	1	<11	>35	1680	156																																																																																
5	1	1-151	<25	284	63																																																																																
In [64]:	<pre>#####Variable definitions##### #District: factor: district of residence of policyholder (1 to 4); 4 is major cities. #Coupran ordered factor: group of car with levels &lt;1 litre, 1-1.5 litre, 1.5-2 litre, &gt;2 litre. #Age:an ordered factor: the age of the insured in 4 groups labelled &lt;25, 25-29, 30-35, &gt;35. #Holders:numbers of policyholders. #Claims:numbers of claims</pre>	In [65]:																																																																																			
In [65]:	<pre>#Question 21 #Write Seaborn code to plot a FacetGrid of boxplots for the number of claims (Claims) feature, conditioned # on the age of the insured (Age)? #Put the charts on one row  #Write your code below:  g=sns.FacetGrid(Insurance_data, col='Age') g.map_dataframe(sns.boxplot, x='Claims')</pre>	Out[65]:																																																																																			
Out[65]:																																																																																					
In [66]:	<pre>#Question 22 #Write Seaborn code to plot a FacetGrid of boxplots for the number of policyholders (Holders) feature, conditio # on the age of the insured (Age)? #Put the charts on one row  #Write your code below:  g=sns.FacetGrid(Insurance_data, col='Age') g.map_dataframe(sns.boxplot, x='Holders')</pre>	Out[66]:																																																																																			
Out[66]:																																																																																					
In [68]:	<pre>#Question 22 #What can you infer about the variability of observations for the Holders feature # when Age is between 0 and 34 relative to when Age greater than 35?  #Write your answer below:  # The range for the holder feature when age is between 0 and 34 is smaller than when Age greater than 35 # The IQR for the holder feature when age is between 0 and 34 is smaller than when Age greater than 35 # The Standard deviation for the holder feature when age is between 0 and 34 is smaller than when Age greater t # The Variance for the holder feature when age is between 0 and 34 is smaller than when Age greater than 35</pre>	In [69]:																																																																																			
In [69]:	<pre>#Question 23 #What is responsible for your inference in Question 22 above?  #Write your answer below: # The shape of boxplot</pre>	In [70]:																																																																																			
In [70]:	<pre>#Question 24 #What can you infer about the variability of observations for the Claims feature around their median # as you read your chart in Question 21 from left to right?  #Write your answer below: # The value of median is increasing from left to right</pre>	In [71]:																																																																																			
In [71]:	<pre>#Question 25 #Write Seaborn code to plot a FacetGrid of regression plots for the Claims-Holders relationship, # conditioned on Age #Put the charts on one row  #Write your code below:  g=sns.FacetGrid(Insurance_data, col='Age') g.map_dataframe(sns.regplot, x='Holders', y='Claims')</pre>	Out[71]:																																																																																			
Out[71]:																																																																																					
In [73]:	<pre>#Question 26 #Using your chart in Question 25 above, what can you infer about the Claims-Holders relationship overall?  #Write your answer below  # they have positive linear relationship, when holders increasing, the Claims increasing # meanwhile when holders decreasing, the Claims decreasing</pre>	In [75]:																																																																																			
In [75]:	<pre>#Question 27 #Write Seaborn code to help you prove the type of model appropriate to use for the Claims-Holders relationship #The pairplot suggests (to some extent) that a non-linear model would be appropriate for the Claims-Holders re #The residual plot for the same relationship appears to indicate that a linear model would be appropriate  sns.jointplot(x=Insurance_data.Holders, y=Insurance_data.Claims, kind='resid') sns.pairplot(Insurance_data, vars=['Holders', 'Claims'])</pre>	Out[75]:																																																																																			
Out[75]:																																																																																					
In [76]:	<pre>#Question 28 #From your answer to Question 27, is a non-linear model the appropriate model for the Claims-Holders relatio #Write your answer below # NO</pre>	In [77]:																																																																																			
In [77]:	<pre>#Question 29 #Write Seaborn code to plot a FacetGrid of kde plots for the Claims feature, # conditioned on Age #Put the charts on one row  #Write your code below:  g=sns.FacetGrid(Insurance_data, col='Age') g.map_dataframe(sns.kdeplot, x='Claims')</pre>	Out[77]:																																																																																			
Out[77]:																																																																																					
In [78]:	<pre>#Question 30 #What does your chart in Question 29 suggest about the behaviour of Claims as you read from left to right # (i.e as you consider changes in demographics indicated by the Age feature)?  #Write your answer below:  # probability of density of claims is decreasing from left to right meanwhile the median is increasing</pre>	In [80]:																																																																																			
In [80]:	<pre>#Question 31 #Write Seaborn code to generate a categorical plot of boxplots for the Holders feature #conditioned on Age  #Write your answer below  sns.catplot(x='Holders', y = 'Age', data=Insurance_data, kind='box')</pre>	Out[80]:																																																																																			
Out[80]:																																																																																					



```
#Question 32

#Assuming constant fitability of existing insurance policies across all Age groups.

#If you are the head of the insurance company sales team, and your current goal is to prevent as many
#policyholders as possible from canceling their insurance policies, which age group would
#you most likely focus marketing resources on to keep their insurance policies active based on your chart in Q

#<25
#25-29
#30-35
#>35

#Write your answer below
#>35

In [82]: #Question 33

#Justify your answer to Question 32 above.

#Write your answer below

#The age group greater than 35 has the most distribution so it is better focus on marketing on this age group

In [83]: #Run this cell and use it to answer Questions 35 - 40 below

import pandas as pd
import seaborn as sns
import numpy as np
from matplotlib import pyplot as plt
import matplotlib
School_data = pd.read_csv("https://vincentarelbundock.github.io/Rdatasets/csv/mlmRev/bdf.csv", index_col=0).dro
School_data.head()

#Write your answer below

Out[83]: schoolNR pupilNR IQ.verb IQ.perf sex Minority repeatgr aritPRET classNR aritPOST ... curmeet mixedgra percmino ar
1 1 17001 15.0 12.33333 0 N 0 14 180 24 ... 1.83333 0 60
2 1 17002 14.5 10.00000 0 Y 0 12 180 19 ... 1.83333 0 60
3 1 17003 9.5 11.00000 0 N 0 10 180 24 ... 1.83333 0 60
4 1 17004 11.0 10.00000 0 N 0 13 180 26 ... 1.83333 0 60
5 1 17005 8.0 6.66667 0 N 0 8 180 9 ... 1.83333 0 60

5 rows x 28 columns

In [44]: #####Variable Definitions#####

# schoolNR
# a factor denoting the school.

# pupilNR
# a factor denoting the pupil.

# IQ.verb
# a numeric vector of verbal IQ scores

# IQ.perf
# a numeric vector of IQ scores.

# sex
# Sex of the student.

# Minority
# a factor indicating if the student is a member of a minority group.

# repeatgr
# an ordered factor indicating if one or more grades have been repeated.

# aritPRET
# a numeric vector

# classNR
# a numeric vector

# aritPOST
# a numeric vector

# langPRET
# a numeric vector

# langPOST
# a numeric vector

# ses
# a numeric vector of socioeconomic status indicators.

# denomina
# a factor indicating if the school is a public school, a Protestant private school, a Catholic private school,

# schoolSES
# a numeric vector

# satiprin
# a numeric vector

# natitest
# a factor with levels 0 and 1

# meetings
# a numeric vector

# curmeet
# a numeric vector

# mixedgra
# a factor indicating if the class is a mixed-grade class.

# percmino
# a numeric vector

# aritdiff
# a numeric vector

# homework
# a numeric vector

# classsis
# a numeric vector

# groupsis
# a numeric vector

In [86]: #Question 34

#Write Seaborn code to generate a FacetGrid of scatter plots for the IQ.perf-IQ.verb relationship according to
#Condition on the denomina feature first
#Condition on the Minority feature second
#Condition on the sex feature last

#Add legend

#Label the horizontal axis "Verbal IQ Score"
#Label the vertical axis "Overall IQ Score"

#You should have a 2-row X 4-column matrix of scatterplots

#Write your code below

g=sns.FacetGrid(School_data, col='denomina', row = 'Minority', hue='sex')
g.map_dataframe(sns.scatterplot, x='IQ.verb', y='IQ.perf')
g.set_axis_labels('Verbal IQ Score', 'Overall IQ Score')
g.add_legend()

Out[86]: Seaborn.axisgrid.FacetGrid at 0x7fcd079b3890>



In [87]: #Question 35

#Which of these charts represent a subset of data does your chart in Question 34 suggest you could safely delete
#from your dataset with minimal impact on your regression modeling results?

#Bottom rightmost scatterplot
#Top leftmost scatterplot
#Top rightmost scatterplot
#Bottom leftmost scatterplot

#Write your answer below
#Bottom rightmost scatterplot

In [88]: #Question 36

#Write Seaborn code to generate a FacetGrid of kde plots for the IQ.perf feature according to the following:
#Condition on the denomina feature

#Add legend

#Write your code below

g=sns.FacetGrid(School_data, col='denomina')
g.map_dataframe(sns.kdeplot, x='IQ.perf')
g.add_legend()

Out[88]: Seaborn.axisgrid.FacetGrid at 0x7fcd08444fa0>



In [90]: #Question 37

#For which value of the denomina feature is your chart in Question 36 closest to the benchmark
#standard normal deviation curve?

#Write your answer below

#denomina = 4

In [91]: #Question 38

#Write Seaborn code to build a pairplot with kde plots in the diagonal using only the following variables:
#aritPRET
#classNR
#langPOST

#langPRET
#langPOST

#Write your answer below

sns.pairplot(School_data, vars=['aritPRET','classNR','aritPOST','langPRET','langPOST'], diag_kind='kde')

Out[91]: Seaborn.axisgrid.PairGrid at 0x7fcd0887d3d0>



In [92]: #Question 39

#From your chart in Question 38, which of these relationships does not appear to exist?

#langPOST-aritPRET
#langPRET-classNR
#langPRET-langPOST
#aritPRET-langPRET

#Write your answer below

#langPRET-classNR

In [93]: #Question 40

#Write Seaborn code to justify why a linear model would be appropriate for the langPOST-langPRET relationship
#Write your answer below

#The pairplot suggests (to some extent) that a non-linear model would be appropriate for the langPOST-langPRET
#The residual plot for the same relationship appears to indicate that a linear model would be appropriate

sns.jointplot(x=School_data.langPRET, y=School_data.langPOST, kind='resid')

Out[93]: Seaborn.axisgrid.JointGrid at 0x7fcd08335850>


```