

In [115]:

```
1 #CS156 - 10.1.py
2 %matplotlib inline
3
4 """
5 See data at: https://gist.github.com/oba2311/c1c63e1f635910e87a4302fdef902f3f
6 """
7
8 import csv
9 import numpy as np
10 import pandas as pd
11 from sklearn.neighbors.kde import KernelDensity
12 import matplotlib.pyplot as plt
13 import random
14
15 random.seed(23)
16
17 data=pd.read_csv('data.csv',delim_whitespace=True)
18 data.head()
```

Out[115]:

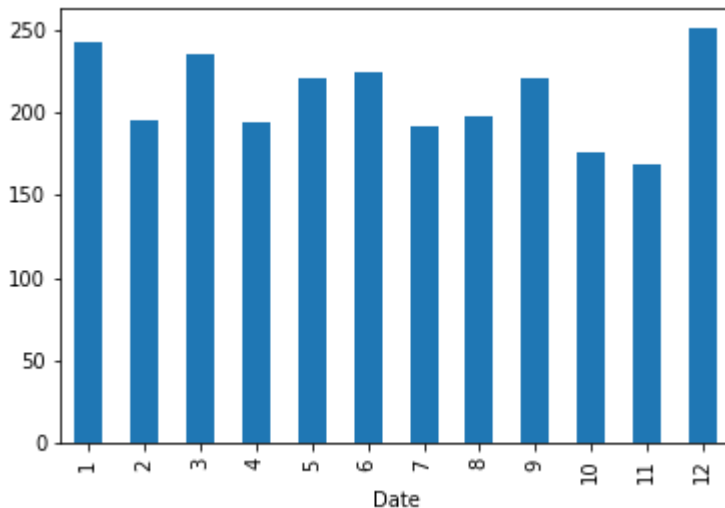
	Date	Amount
0	25.5.2016	54,241.35
1	29.5.2017	54,008.83
2	30.6.2017	54,008.82
3	05.1.2017	52,704.37
4	23.2.2017	52,704.36

In [3]:

```
1 data['Date']=pd.to_datetime(data['Date'],infer_datetime_format=True)
2 months = data.groupby(data["Date"].dt.month).count()
3 plot_object_of_months = data.groupby(data["Date"].dt.month).count()
4 plot_object_of_months= plot_object_of_months.loc[:, "Amount"]
5 plot_object_of_months.plot(kind="bar")
```

Out[3]:

<matplotlib.axes._subplots.AxesSubplot at 0x10d77bf90>

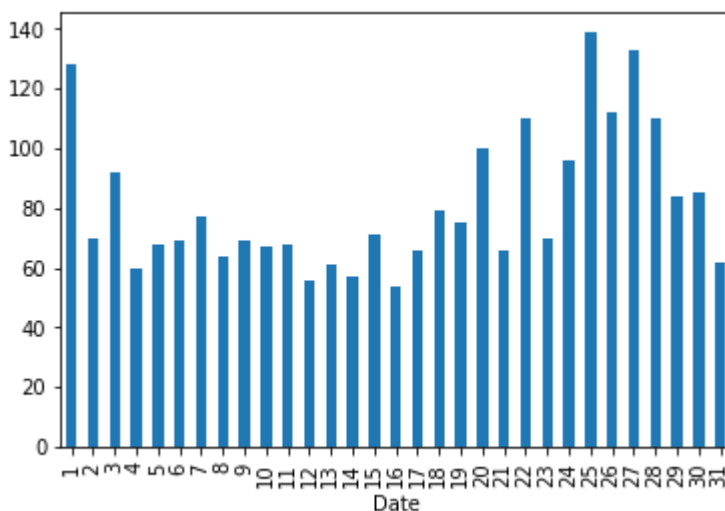


In [125]:

```
1 data['Date']=pd.to_datetime(data['Date'],infer_datetime_format=True)
2 plot_object_of_days = data.groupby(data["Date"].dt.day).count()
3 days = data.groupby(data["Date"].dt.day).count()
4 plot_object_of_days= plot_object_of_days.loc[:, "Amount"]
5
6 plot_object_of_days.plot(kind="bar")
```

Out[125]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a18d3cf50>

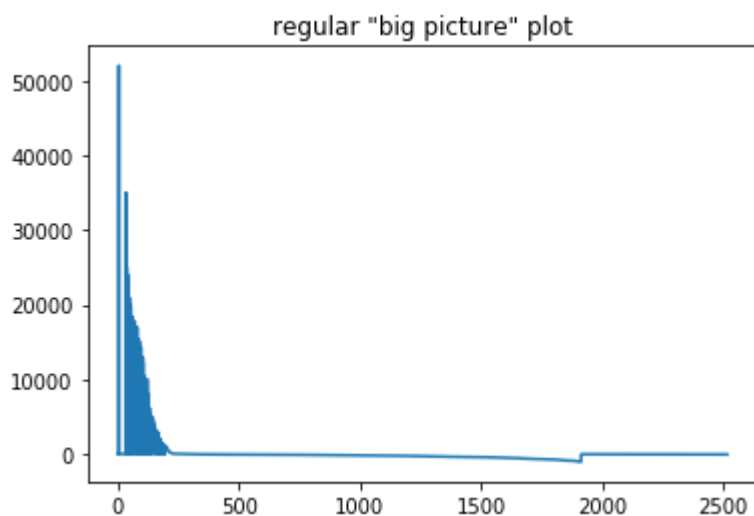


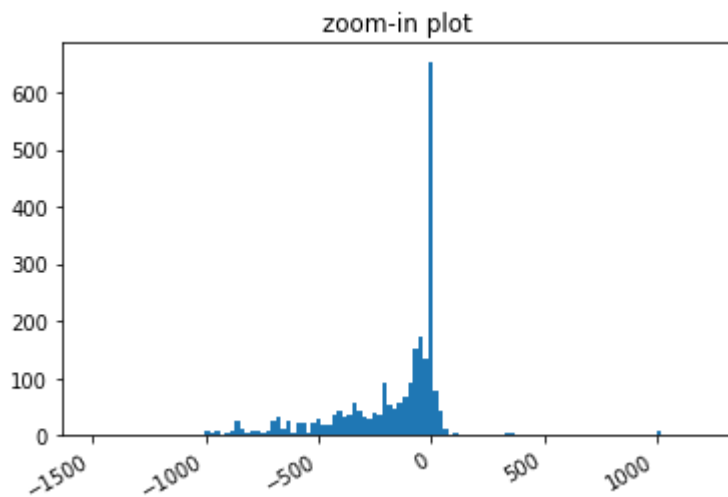
In [126]:

```
1 amount = data['Amount']
2 amount = amount.values.reshape(-1,1)
3
4 amount_as_num = []
5 for i in amount:
6     i = i[0].replace(',','.')
7     i = np.fromstring(i, dtype=np.float, sep=',')
8     amount_as_num.append(i)
9
10 data["amount_as_num"] = amount_as_num
11
12 amount_as_num = np.array(amount_as_num).reshape(-1,1)
13
14 print data["amount_as_num"].head()
15
16 #Here we zoom in to see most of the data. This can be justified by looking at t
17 plt.plot(amount_as_num)
18 plt.title("regular \"big picture\" plot")
19
20 #Just a regular plotting of the data shows that the data is somewhat weird with
21
22 fig, ax = plt.subplots()
23 #Therefore, an option is to zoom in and try to better understand the area under
24 # we are taking the risk of generating data that will not fake accurately outlie
25 plt.hist(amount_as_num,bins = 120,range=(-1500,1200),align='mid')
26 plt.title("zoom-in plot")
27
28 fig.autofmt_xdate() # make space for and rotate the x-axis tick labels
29 plt.show()
```

```
0    [54.241]
1    [54.008]
2    [54.008]
3    [52.704]
4    [52.704]
```

Name: amount_as_num, dtype: object





In [127]:

```
1 kde_month = KernelDensity(kernel='gaussian', bandwidth=0.2).fit(months)
2 kde_month.score_samples(months)
```

Out[127]:

```
array([-1.10390789, -1.10390789, -1.10390789, -1.10390789, -0.4107607
1,
      -1.10390789, -1.10390789, -1.10390789, -0.41076071, -1.1039078
9,
      -1.10390789, -1.10390789])
```

In [128]:

```
1 kde_month.sample(10)
```

Out[128]:

```
array([[ 176.37641122,  175.8224081 ],
      [ 190.95049999,  191.2712638 ],
      [ 223.91574546,  223.70898105],
      [ 221.04182044,  220.79627387],
      [ 250.99192953,  250.98516124],
      [ 175.95329179,  175.96377093],
      [ 243.03881323,  242.9664534 ],
      [ 250.91757258,  250.89227902],
      [ 175.67315668,  176.20988767],
      [ 224.07053017,  223.8826399 ]])
```

In [129]:

```
1 kde_day = KernelDensity(kernel='gaussian', bandwidth=0.2).fit(days)
2 kde_day.score_samples(days)
```

Out[129]:

```
array([-2.05298845, -1.35984127, -2.05298845, -2.05298845, -1.3598412
7,
      -1.35984127, -2.05298845, -2.05298845, -1.35984127, -2.0529884
5,
      -1.35984127, -2.05298845, -2.05298845, -2.05298845, -2.0529884
5,
      -2.05298845, -1.35984127, -2.05298845, -2.05298845, -2.0529884
5,
      -1.35984127, -1.35984127, -1.35984127, -2.05298845, -2.0529884
5,
      -2.05298845, -2.05298845, -1.35984127, -2.05298845, -2.0529884
5,
      -2.05298845])
```

In [130]:

```
1 kde_day.sample(10)
```

Out[130]:

```
array([[ 67.77015889,  67.54265004],
       [ 65.988236 ,  65.70493764],
       [ 75.06384941,  75.2135438 ],
       [110.02243013, 109.90202193],
       [ 56.65548486,  57.17725381],
       [ 68.54562583,  69.55292289],
       [ 55.68722254,  55.81888602],
       [ 60.01100802,  60.0279363 ],
       [ 92.16851497,  92.2393542 ],
       [128.04089359, 127.93803894]])
```

In [131]:

```
1 kde_amount = KernelDensity(kernel='gaussian', bandwidth=0.2).fit(amount_as_num)
2 kde_amount.score_samples(amount_as_num)
3 kde_amount.sample(10)
```

Out[131]:

```
array([[ -58.77152887],
       [ -47.14064439],
       [-260.07169609],
       [-164.8674097 ],
       [  -7.85387687],
       [ 38.01915293],
       [  -1.1661075 ],
       [  -5.73407325],
       [-27.20152874],
       [-846.83686649]])
```

Sampling from these density models, create a fictitious month of personal transactions:

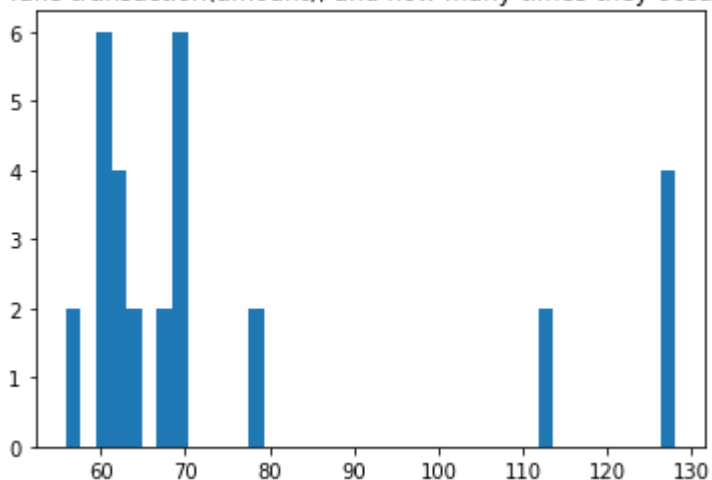
In [132]:

```
1 Month = kde_day.sample(15)
2 Month=Month.reshape(-1,1)
3
4 plt.hist(Month, bins = 40)
5 plt.title("fake transaction(amount), and how many times they occur:")
6 print "here's the data:", Month
```

here's the data: [[60.90855642]

```
[ 61.26446341]
[ 127.97035796]
[ 127.85768522]
[ 60.10264321]
[ 60.01020326]
[ 79.24471743]
[ 78.74509348]
[ 66.98870391]
[ 67.08969987]
[ 111.84698394]
[ 112.20338867]
[ 127.96245835]
[ 128.071591 ]
[ 62.16235396]
[ 62.0289875 ]
[ 55.85694027]
[ 55.86262854]
[ 62.01164561]
[ 61.60400723]
[ 64.35600784]
[ 64.36502073]
[ 70.0263983 ]
[ 70.06341142]
[ 69.87134987]
[ 70.1280842 ]
[ 59.8175005 ]
[ 60.23727729]
[ 68.69229151]
[ 69.2972583 ]]
```

fake transaction(amount), and how many times they occur:



In [149]:

```
1
2 print "total number of transactions:" ,np.mean(kde_month.sample(10))
3 print ""
4 notice that by averaging high numbers of generated transactions, we are getting
5 central limit theorem. That said, by doing so we are still interpolating in a sense
6 ""
```

total number of transactions: 213.734830778

notice that by averaging high numbers of generated transactions, we are getting at a more likely expected number, based on the central limit theorem. That said, by doing so we are still interpolating in a sense, since the data is still based only on the data we trained on. So, the risk of overfitting is real.

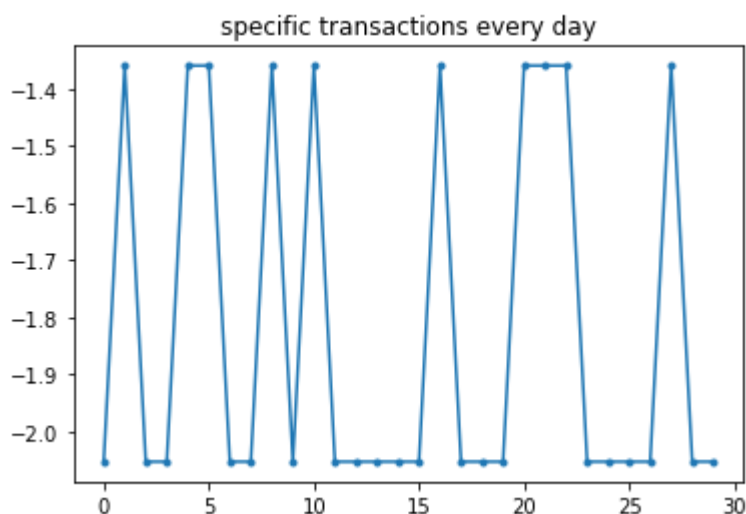
In [154]:

```
1 days_data = np.array(zip(kde_day.score_samples(days),range(1,31)))[1:,0]
2 print days_data
3 plt.plot(days_data,marker='.')
4 plt.title("specific transactions every day")
```

```
[ -2.05298845 -1.35984127 -2.05298845 -2.05298845 -1.35984127 -1.359841
27 -2.05298845 -2.05298845 -1.35984127 -2.05298845 -1.35984127 -2.052988
45 -2.05298845 -2.05298845 -2.05298845 -2.05298845 -1.35984127 -2.052988
45 -2.05298845 -2.05298845 -1.35984127 -1.35984127 -1.35984127 -2.052988
45 -2.05298845 -2.05298845 -2.05298845 -1.35984127 -2.05298845 -2.052988
45]
```

Out[154]:

Text(0.5,1,u'specific transactions every day')



Discussion regarding parameters: The parameter used in KDE is bandwidth h , which is the denominator in the kernel, influencing sensitivity. "Mathematically, a kernel is a positive function $K(x;h)$ which is controlled by the bandwidth parameter h . Given this kernel form, the density estimate at a point y within a group of points $x_i; i = 1 \dots N$ is given by:"

$$\rho_K(y) = \sum_{i=1}^N K((y - x_i)/h)$$

Therefore, we see that it influences the bias-variance tradoff: small bandwidth will lead to an unsmooth density distribution, with higher variance (and therefore bigger range that is more likely to include potential outcome), and higher h will lead to higher bias and smooth density distribution.

Part Two - fiction poetry:

In [153]:

```
1  #week10_assignment_topic_extraction.py
2
3  from collections import OrderedDict
4  import pandas as pd
5  from nltk.corpus import stopwords
6  import string
7  import re
8  import gensim
9  import numpy as np
10 from nltk.corpus import stopwords
11 from time import time
12
13 print "===== Data Prep: ====="
14
15
16 headings = ['CHAPTER I', 'CHAPTER II',
17             'CHAPTER III', 'CHAPTER IV', 'CHAPTER V',
18             'CHAPTER VI', 'CHAPTER VII', 'CHAPTER VIII',
19             'CHAPTER IX', 'CHAPTER X', 'CHAPTER XI',
20             'CHAPTER XII', 'PREFACE', 'NOTE']
21
22 terminate_code = "THE END"
23
24 path_to_book = "ALICE.txt"
25
26 def extract_chapters(path_to_book):
27
28     with open(path_to_book, 'r') as book:
29         t = book.readlines()
30
31         # We're going to put all the
32         # sections/chapters in a dictionary
33         chapter_dict = OrderedDict()
34
35         # Initialize empty sections
36         chapter_text = ''
37         chapter_name = None
38
39         for line in t:
40
41             if not chapter_name and any([heading in line for heading in headings]):
42                 chapter_name = line.replace('\r\n', '')
43
44             if chapter_name:
45
46                 # Populate section string with line
47                 if not any([heading in line for heading in headings]):
48                     chapter_text += line
49
50                 # If a heading line, we've hit a new
51                 # chapter/section; throw the text string into
52                 # the dictionary and start a new section
53                 if any([heading in line for heading in headings]):
54                     chapter_dict[chapter_name] = [chapter_text]
55                     chapter_name = line.replace('\r\n', '')
56                     chapter_text = ''
57
58                 # If we hit the end of the book,
59                 # throw everything into last chapter
```

```

60         if terminate_code in line:
61             chapter_dict[chapter_name] = [chapter_text]
62
63         # Make dataframe to store chapters
64         df = pd.DataFrame.from_dict(chapter_dict, orient='index')
65         df.columns = ['raw_text']
66
67         return df
68
69     '''Clean up the text!'''
70
71     df = extract_chapters(path_to_book)
72
73     # Establish stopwords
74     sw = set(stopwords.words('english'))
75
76     # Remove punctuation & make lowercase
77     df['no_punctuation'] = df.raw_text.map(
78         lambda x: x.translate(None, string.punctuation).lower())
79
80     # Remove stopwords
81     df['no_stopwords'] = df.no_punctuation.map(
82         lambda x: [word for word in x.split() if word not in sw])
83
84     df['no_stopwords'] = df.no_stopwords.map(
85         lambda x: " ".join(x))
86
87     print "===== Top Words Extraction : =====
88
89     n_components = 10
90     n_samples = 1000
91     n_features = 1000
92     n_top_words = 20
93
94     def print_top_words(model, feature_names, n_top_words):
95         for topic_idx, topic in enumerate(model.components_):
96             message = "Topic #%d: " % topic_idx
97             message += " ".join([feature_names[i]
98                                 for i in topic.argsort()[::-n_top_words - 1:-1]])
99             print(message)
100         print()
101
102     print "===== TF - for LDA : =====
103
104     from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
105
106     # Use tf (raw term count) features for LDA.
107     print("Extracting tf features for LDA...")
108     tf_vectorizer = CountVectorizer(max_df=0.95, min_df=2,
109                                   max_features=n_features,
110                                   stop_words='english')
111     t0 = time()
112     from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
113     tf = tf_vectorizer.fit_transform(df["no_stopwords"])
114     print("done in %0.3fs." % (time() - t0))
115
116     print "===== LDA - Topic Extraction - Pe
117
118     from sklearn.decomposition import LatentDirichletAllocation
119
120     print("Fitting LDA models with tf features, "

```

```

121         "n_samples=%d and n_features=%d..."
122         % (n_samples, n_features))
123     lda = LatentDirichletAllocation(n_components=n_components, max_iter=5,
124                                     learning_method='online',
125                                     learning_offset=50.,
126                                     random_state=0)
127     t0 = time()
128     lda.fit(tf)
129     print("done in %0.3fs." % (time() - t0))
130
131     print("\nTopics in LDA model: These are the top 10 topics the model extracts fo
132     tf_feature_names = tf_vectorizer.get_feature_names()
133     print_top_words(lda, tf_feature_names, n_top_words)
134
135     print """Put together: \n
136     a hatter a king and a dormouse march with a hare a mouse and a queen.
137     The hatter didn stand the summer jury so he startled twinkling anxiously and sa
138     caterpillar came and the father felt bit mouth size took girls good inches.
139     Sharp question effected mushroom, the mouse, the rabbit ran to bring gloves soc
140     queen and duchess quite, but the cat say came to the turtle. The rabbit used la
141     The turtle gryphon mock soup beautiful(y) voice and won.
142     mouse came nad the cat quite. the caterpillar tell a thing and let voice of the
143     queen and king('s) voice came going, and soldiers anxiously quite white tone wi
144     A bottle wrote players small angrily eaglet sight unfolded march lives believe
145     king and the gryphon voice made the jury, queen,and the course and the mouse l
146     """
147
148     print "===== LDA - Topic Extraction, per
149
150     for index, row in enumerate(df["no_stopwords"]):
151         lda = LatentDirichletAllocation(n_components=n_components, max_iter=5,
152                                         learning_method='online',
153                                         learning_offset=50.,
154                                         random_state=0)
155         print("\nTopics in LDA model: These are the top 10 topics the model extract
156         tf = tf_vectorizer.fit_transform(row.split())
157         lda.fit(tf)
158         tf_feature_names = tf_vectorizer.get_feature_names()
159         print_top_words(lda, tf_feature_names, n_top_words)
160
161     print "In this part, I worked with Skye Herish, and she provided a very handy wa

```

```

=====
=====
===== TF - for LDA : =====
=====
Extracting tf features for LDA...
done in 0.041s.
===== LDA - Topic Extraction -
Per Book : =====
=====
Fitting LDA models with tf features, n_samples=1000 and n_features=100
0...
done in 0.851s.

```

Topics in LDA model: These are the top 10 topics the model extracts for the whole book:

Topic #0: hatter king dormouse march hare mouse queen court ll rabbit

white thing great quite say jury added voice course looking

Topic #1: sleep hatter shut didn stand summer jury sleepy march whiske

The code works in two different options, to examine different possibilities: running LDA over the whole book, and running over each chapter. The results of the first option are summarized into a short, interesting plot.

In []:

1	
---	--