7 202	state = df_kaggi ("After filtered state.head(10) # Is dataset, we have filtered dataset Ate Northeastern_state 11- 01 vern 21- 01 pennsylva	le.loc[df_kaggle[ed dataset to nor## prints first 1 ave 3112 records t to northeastern ates state_id median	"rhode island" "Northeastern_state: theastern region, we o rows of filtered of and 39 columns. I region, we have 54	e have", df_nestate.s dataframe 9 records and 39	setts"] n_states)] # sele shape[0], " recor columns.	ecting all rows from rds and ", df_nestate	given options of sta e.shape[1], "column ting_count_mm active_list 0.0702	ns.")	n_days_on_market me 46.0 41.0	edian_square_feet_yy aver -0.0562 -0.0564	rage_listing_price aver 5.383490e+05 4.012947e+05	rage_listing_pric - -
13 202 07- 26 202 07- 28 202 07- 32 202 07- 41 202 07- 50 202 07-	massachus massac	setts ma aine me ticut ct shire nh	419900.0 649000.0 339900.0 489900.0 424900.0 445000.0	0.0344 -0.0159 0.0034 -0.0092 0.0130 -0.0089 -0.0300	-0.0107 0.0469 0.1368 0.2295 0.0827 0.0023	1417 9209 3536 7764 2321 25084 57328	0.1400 0.0205 0.1802 0.0774 0.1105 0.0599	-0.2811 -0.2615 -0.3571 -0.5490 -0.2664 -0.1966 -0.1335	30.0 31.0 33.0 36.0 28.0 73.0	-0.0962 -0.0911 -0.0070 0.0021 -0.0486 -1.0000 -0.1042	7.778582e+05 1.182808e+06 5.076741e+05 1.143219e+06 5.736333e+05 6.511915e+05 1.215842e+06	
62 202 06- 10 rows #df_n	× 39 columns estate.describe oving the hashtags(#	aine me ().T #describes t ##) or uncommenting	338750.0 The statistics of the	-0.0007 e dataframe, transpos nat each column of data o	0.1701 se for better vie	2996 ewing	0.2249	-0.4651	28.0 max and quartiles of each	-0.0047	5.062248e+05	-(
Box One way	yter Notebook analy plot y to plot boxplot is u	yzed the average hou sing seaborn library.	using price of these state Vertical box plot is group	oed by categorical variable	ce: Kaggle) through t	box plots and bar chart ra	ces. a dataframe (our data) into	o seaborn's boxplot.				
box_pplt.yplt.xplt.xplt.xplt.t	lot = sns.boxplot label('Price in label('Price in label('State') ticks(rotation itle("Avg. hous: 0.5, 1.0, 'Avg. label("Avg. hous: 2.5, 1.0, 'Avg. label("Avg. hous: 2.6, housing pri	ct (x = 'Northeast millions') = 45) ing price of Northeast housing price of ce of Northeast thuse the connecticut manner State ex plot that represent	heastern states", for Northeastern states tern states tern states tern states states Angle Representation of the states of	ontsize = 15, fontwe	ight = 'bold')	g price value (somewhere	•		e (between 1.0 - 1.2) is m value present at the top o			
df_ra #crea #pass df_ra df_ra Northea	years. ce = df_sort[['te spreadsheet sed in values arg ce=df_race.pivotce astern_states conrect Date 2016-07-01 7.5983 2016-08-01 7.5065 2016-09-01 7.5657 2016-10-01 7.6577 2016-11-01 7.7227 2021-03-01 9.4362 2021-04-01 9.1653 2021-05-01 8.9846 2021-06-01 1.1684 2021-07-01 1.1432 × 9 columns	Date', 'Northeaste style pivot table gument. Columns at table (values = necticut maine 300e+05 331254.6000 384e+05 332513.9006 340e+05 328398.3357 320e+05 328159.4010	massachusetts new harmassachusetts new harmassachus	_listing_price']] e statistical summary cording to the feature ice',index=['Date'],6 mpshire new jersey 14.6333	y of feature, incres. columns = 'Northe new york pennsylv 362448e+05 272366.6 734834e+05 270407.8 347256e+05 271526.0 247852e+05 271621.3 744485e+05 270294.0 295947e+06 432894.3 315247e+06 438361.3 297325e+06 426755.0 269033e+06 413517.8	dex is the featured eastern_states') rania rhode island ver 6627 539143.5216 358644 8052 535634.2990 358844 0676 524995.7656 360144 7676 517423.0132 363224 0491 522212.3088 369524 0	3.3840 3.4227 3.8716 5.2592 1.7021 3.2998 5.3699 2.2729 3.4730		chart. Since it is self-expla , filename = None)	natory chart and is visual	ly appealing, it receive	ed huge audienc
Refer ba	els now be observing th		manually populated in E	xcel. Majority of data wer _Test.xlsx) to understanc	re collected from US	-	w, Wikipedia. 21 features	of the real estates lik	te median listing price, Ho	me Value Index, housing	units, Crime rate per 1	100,000 populat
df_re df_re df_re 0 Co	g = pd.read_csv(g.fillna(df_reg) g Crime rate per 100,000 population nnecticut 181.6	ZHVI (Zillow Home Value Index) 247706 1491786 244109 1524992	e=True) # data wrang	0.0 11930.0 8.0 17249.0			1240.0 1828.0 2336.0	dian_square_feet ave 2016.0 2083.5 1888.0	erage_listing_price total_lis 847461.2371 724776.1123 732724.4979	sting_count pending_ration 12088.0 0.013244 17340.0 0.005276 20747.0 0.003385	54	tness_score su 64.806595 31.261890 59.321497
4 Col 94 95 96 97	nnecticut 227.7 Vermont 102.6 Vermont 114.9 Vermont 142.3 Vermont 135.2	2429211512528227102328991227340327315227048327315	74304 33990 75923 33990 60708 33945 54842 33945 55582 33945 51862 33945	0.0 21450.0 0.0 14294.0 0.0 14294.0 0.0 14294.0 0.0 14294.0	2644.0 3252.0 2644.0 2644.0 2644.0 2644.0	116.0 92.0 92.0 92.0 92.0 92.0 92.0	2216.0 2320.0 1828.0 1828.0 1828.0 1828.0 1828.0	1579.0 1976.0 1677.0 1677.0 1677.0 1677.0	715891.4913 768576.4007 522433.7284 522433.7284 522433.7284 522433.7284 522433.7284	20147.0 0.012666 25152.0 0.172587 16375.0 0.205254 16375.0 0.205254 16375.0 0.205254 16375.0 0.205254	450 387 823 1295 801	51.870640 73.208624 56.436271 14.838301 31.388713 19.942930 16.930881
Corre The follo by a cer plt.f cm =	elation bwing chart shows he tain amount the other igure(figsize=(1) df_reg.corr(methor)	now related are each er changes on an ave	variables with each othe erage by a certain amoui	ached with the Honors Proceeds. The Red shows positive cores of the contract	relation and blue sho		The darker the shade of co	olor, the stronger the	relation between these va	ariables. If the two variabl	es are related it mean	s that when one
mask sns.h plt.g plt.s	= np.triu(np.one eatmap(cm, annot cf().set_size_ir how() rate per 100,000 popul IVI (Zillow Home Value Ir Total housing I Median Inc median_listing_ active_listing_c price_increased_c	es_like(df_reg.co t= True , cmap = 'R nches(20, 10)	orr())) cdBu_r',linewidth=0.9 6 60.082 1 0.13 0.96 6 0.074 0.88 0.94	5, square= True , mask =	mask)	- 0.8 - 0.6 - 0.4						
median_	median_square median_square average_listing_o total_listing_o pending_ nielsen_hh_ hotness_s supply_s demand_s median_days_on_ma	count - 0.34 0.21 0.62 0.05 ratio -0.0150.0260.0510.1 rank - 0.28 -0.36 -0.15 -0.1 score -0.0050.24 -0.15 0.24 score -0.13 0.370.00210.2 score -0.150.017 -0.26 0.14 arket - 0.06 -0.25 0.22 -0.3	1 0.59 0.43 0.46 0.38 0.45 0.0 250.023 0.53 0.64 0.53 0.57 0.3 1 0.85 0.4 0.36 0.31 0.31 0.4 64 0.13 0.97 0.92 0.89 0.93 0.8 7 0.13 0.11 0.0750.029 0.12 0.3 2 0.2 0.0250.0470.0150.0310.1 4 0.34 0.016.00250.0260.020.0 7 0.33 0.068 0.11 0.0560.0810.0 4 0.24 0.1 0.12 0.11 0.120.0 3 0.16 0.030.0560.0540.040.0 8 0.13 0.0610.0750.0720.0630.0	24 0.64 8 0.48-0.022 16 0.34-0.47 0.45 13 0.29 0.03 0.13 0.032 13-0.0770.051-0.160.027 0.2 18 0.1 -0.18 0.27 0.056.00220.6 14 0.26 0.064 0.16-0.0580.24 0.0 160.0810.092-0.15-0.0270.0680.3	34 0.83 0.44 37 -0.68 -0.68 -0.47 3 -0.67 -0.63 -0.49 0.98 uarket	-0.2 -0.0 0.2 0.4 0.6						
matrix, s Multi Logistic To unde from from	class logist regression is only sometime the performant sklearn.model_sesklearn.linear_refine the features	ic regression coefficient ic regression uitable for predicting ance of the model, the election import to model import Logical and target (X and target (X and target)	of ZHVI and Crime rate. (Softmax reg probability on a binary rate dataset is divided into the crain_test_split esticRegression	Similarly, the correlation of ression) ange from 0 to 1. As my ta	coefficient of pending	g ratio with Total housing o	units is negative 0.051.	i.e. 9 states) , I am us	et, gives us numerical val sing Softmax Regression t model testing. This way, w	o handle multiple states.		41 at the top of
<pre>pd.un X = d #crea y = d y_num y = y</pre>	<pre>ique(df_reg["Sta f_reg.drop(["Sta ting target vect f_reg["States"] eric_dict={'Cont 'New ' 'Vermo' .apply(lambda x)</pre>	ates"]) # target ates"], axis = 1) tor and convertin mecticut':1,'Main Jersey': 5, 'New ont': 9} cross y_numeric_dict[class .astype(float) g to numeric codes e':2,'Massachusetts York': 6, 'Pennsylv x])	':3, 'New Hampshire' vania': 7, 'Rhode Is	land':8,	ata. K-nearest neighbors e	estimates the likelihood of	f data point based on	what group the data point	s is nearest to.		
from from # Set # gen knn_m	sklearn.model_se sklearn.metrics up the model we eral rule of the odel = KNeighbor	umb for n_neighborsClassifier(n_ne	Fold fusion_matrix ers : square root of ighbors=10)	samples (eg:79 for s		et to measure the accurac	cy of our model.					
In the most to under to under to under to under to under the tour tour tour tour tour tour tour tour	atrix, consider: 'Cordel shows that Massestand how the house of the stand how	nnecticut':1,'Maine':2, sachusetts (labeled a sing price of whole N' CCUTACY simport make_cla election import celection import R ession(multi_clas valuation procedu iedKFold(n_splits and collect the _score(model, X, erformance : %.3f (%.3f)' % (0.069) nodel achieved a mea understand that since estate model as it doe ON (AIC) compares the le by R as the result i	'Massachusetts':3, 'New is 3) 's data was more con y state can be averaged. "ssification ross_val_score repeatedStratifiedKF6 res='multinomial', so re reliable. "an classification accuracy (n_scores.mean(), n_man classification accuracy response there is a huge variation response rot give idea on variation reliable.	nsistent and has similar has been state and has similar has been and similar has been and similar has been similar	nousing patterns thro zed inconsistent data) state dataset. This is s within the same geo e rather gives an idea	an impressive accuracy of ographical region (eg: New on what states was most	hand, since there is a hig it well, hence 1 out of 4. on the test set. w York has huge disparity favored by our data. This	in real house price in results in biasednes	g price of New York, rangi n regards to its location as as as not all zip codes were forms in Python, which do	s stated in our visualizatio e taken into consideration	n). The Softmax Regr ns but an average or m	ession was not a
<pre>y = d #add X2 = est = #atte est2 print</pre>	f_reg["average_1 constant to pred sm.add_constant(sm.OLS(y, X2.as mpt to fit regre = est.fit() (est2.summary()) cariable: average_1	listing_price"] dictor variables (X) stype(float)) ession model OLS Regr erage_listing_pri	Lce R-squared: DLS Adj. R-squared Tes F-statistic:	======================================	.955 .943 2.51							
Time: No. Ob Df Res Df Mod Covari ====== const Crime ZHVI (Total Median median active	ance Type: rate per 100,000 Zillow Home Value housing Units Income Listing_price	nonrobu ====================================	24 Log-Likelihood 99 AIC: 78 BIC: 20 ust coef std er 193.0974 93.67 -0.5119 0.14 -0.0080 0.00 -0.4015 0.64 1.8584 0.18 2.018e+04 5117.45	-11 2 2 2 T t P> 4 -2.707 0.1 2.061 0.5 -3.531 0.4 -2.135 0.9 -0.619 0.8 9.892 0.9 9.3.944 0.	.76.5 .395. .449. 	0.975] -5.15e+04 379.583 -0.223 -0.001 0.891 2.232 3.04e+04						
price_ price_ pendir mediar total_ pendir nielse hotnes supply demand mediar		per_square_foot - - - -	75.7659 17.84 -616.7901 144.39 -114.4701 13.75 2.018e+04 5117.87 785.4590 281.63 79.8258 30.30 2.017e+04 5117.90 1.049e+04 1.8e+0 7.3074 18.04 6.436e+09 1.38e+1 3.218e+09 6.91e+0 3.218e+09 6.91e+0 -891.8093 561.03 760.5292 523.90 ===================================	9 -4.271 0. 3 -8.323 0. 9 3.942 0. 5 2.789 0. 4 2.634 0. 5 -3.942 0. 4 -0.581 0. 6 0.405 0. 0 0.466 0. 9 -0.466 0. 9 -0.466 0. 8 -1.590 0.		111.285 -329.314 -87.090 3.04e+04 1346.151 140.157 -9983.878 2.54e+04 43.235 3.39e+10 1.05e+10 1.05e+10 225.131 1803.551						
Skew: Kurtos ====== Notes: [1] St [2] Th strong	andard Errors and second secon	0.441 4.915 ssume that the convalue is 5.45e-1 ity problems or to Variables Step Df Develored		0.00010 1.59e+1 ====================================	3 :=							
6 7 - 8 9 10 11 12 R predict	- Me ime.rate.per.100.000 - ho ZHVIZillow.Home.V - pe - active_li - median_day - niel - price_incr eted that 11 data fea	dian.Income 1 630 .population 1 1420 tness_score 1 21290 alue.Index. 1 59933 nding_ratio 1 95646 sting_count 1 127036 s_on_market 1 145685 sen_hh_rank 1 253892 eased_count 1 349805 atures were considered	993.6 20 95078964 918.0 21 95291873 329.2 22 95891207 327.4 23 96847670 868.9 24 98118039 925.8 25 99574898 553.3 26 102113823 158.9 27 105611875 ed valuable attributes to commercial entributes and commercial entributes are	83 914.5859 77 912.5927 95 910.6933 24 908.9755 51 907.4221 20 906.0085 46 904.6718 99 903.8048 58 903.3205 Dur average real estate promeaningless variables to	our model. That leav				ng_count, new_listing_cou , the intercept tells the ave			
We have Rate, In First, we compon	e 22 sets of variable frastructures can be start by standardizents of pca with (.co	grouped into locatio	e to fit each variables with n features to reduce dim mon practice to normalize	ensionality of data.					ated data together. For eg			
<pre>df_PC from pca = pca.f #Acce pca.c</pre>	A = pd.DataFrame sklearn.decompos PCA(n_component it(df_PCA) ssing the components_	e(df_PCA) sition import PCA ts = 4) # projec	orm(df_norm) # norma. It from 21 to 4 dimen	nsions								
; array(5.98240902e-0 3.69726953e-0 2.92397809e-0 1.93136928e-0 -4.94753779e-0 -2.42335583e-0 [-6.63891596e-0 -2.23033592e-0 6.04975419e-0 1.72169962e-0 -2.09758546e-0 2.11954285e-0 -3.02956713e-0 [2.56561183e-0 9.40170173e-0 -8.45402242e-0 9.70052594e-0 3.60801468e-0 2.89667234e-0 -1.31304684e-0 [-3.09681672e-0	02, 1.03087678e- 01, 3.52691427e- 01, -1.23612317e- 01, 3.64561577e- 02, 3.54905839e- 02, -4.39818693e- 03, -2.25715490e- 01, -2.67646721e- 02, 6.57173421e- 02, -2.30651615e- 01, -3.99434118e- 01, 3.61754026e- 01, 2.44194226e- 01, 2.44194226e- 02, 3.89139106e- 02, -6.22582772e- 02, 4.34617139e- 01, -1.94213278e- 01, -1.75234188e- 01, 2.79717452e- 01, -2.24221570e-	2.06303792e-01 3.62521104e-01 01, 3.59405929e-01 01, -2.23163508e-01 01, -2.47936053e-05 02, 7.97936196e-02 02, -4.65495100e-02 01, 7.51042544e-02 02, 7.85795628e-02 01, -4.06485988e-02 01, -4.06485988e-02 01, -3.71355081e-01 01, 3.44685251e-01 01, 7.36918653e-02 01, -7.36918653e-02 01, -7.36918653e-02 01, 2.55150921e-01 01, 2.55150921e-01 01, 1.64383466e-01 01, 3.00221741e-01 01, -2.04831140e-01 04, 4.08315394e-02									
colum df_pc df_pc	1.25107550e-0 3.24258546e-0 8.67800088e-0 5.27522305e-0 3.24966735e-0 y the transformation of the second of the s	02, 9.13367413e- 01, 1.19297172e- 02, 1.42877180e- 01, 5.52300368e- 01, 7.27496210e- ation and convert (x) for x in range (pca.transform(d	1.11944411e-02 01, 7.31741830e-02 01, 4.59761790e-01 02, -2.07633972e-01 04, -5.16749618e-02 the result into a late (1, pca.n_components) e(1, pca.n_components)	DataFrame.								
95 -1.9 96 -1.9 97 -1.4 98 -1.4 We will i	('Explained var	0.630199 0.604719 0.429860 1.715031 0.365540 0.571030 0.351001 0.696570 ance_ratio. Highest valiation per princi	pal component: {}'.	vill be the first principal co	_variance_ratio_)							
SCre	ents to understand Plot lot(pca.explaine label('number of label('cumulativ	oonent 1 holds 32% o how our samples are ed_variance_)	of information, principal control of information, principal control of the distributed among nine of the distributed among nin	·	_	v principal component 3 a	nd 8% by principal compo	onent 4. About 25% o	f data information was los	t during the analysis. Nov	v let's visualize the sai	mples along 4 p
cumulative explained variance												
3D F	numb is just an evidence Plot of Princi gure helps us obser	that most of the feat ipal Compor ve classes (listed in land) other states. Hov	legend) and how they are		·	onging to the same class a	are close to each other, ar	nd the points or imag	es that are very different s	semantically are further a	way from each other. I	New York's hous
is sprea		oca.transform(df_	PCA), x=0, y=1,z =2	color=df_reg['State	es']) A 2							color Connect Maine Massach New Hal New Jer New Yor Pennsylv Rhode Is
is sprea	et of variable	es on each d	components		2 0 2 2 4 6 2 1		6 0					
impor fig = fig.s		components_, p='YlGnBu', cklabels=["PCA"+ cklabels=list(df_	estr(x) for x in rand PCA.columns), on": "vertical"})	ge(1,pca.n_component:	s_+1)],							
importing = fig.s Effect ax = try: excep	ytid xtid cbai		-0.4 -0.2 -0.0 0.2									