

Kings County Housing Sales

Final Project

Presented By:
Ahmed Raza Amaan

1. Provide a Table with descriptive statistics of all the numerical variables. What can you say about the variability in price and size of houses (sqft_living, sqft_lot, sqft_above, sqft_basement) in King County?

	price	sqft_living	sqft_lot	sqft_above	sqft_basement	
Mean	540088.14	Mean	2079.90	Mean	15106.97	Mean
Standard Error	2497.23	Standard Error	6.25	Standard Error	281.75	Standard Error
Median	450000.00	Median	1910.00	Median	7618.00	Median
Mode	450000.00	Mode	1300.00	Mode	5000.00	Mode
Standard Deviation	367127.20	Standard Deviation	918.44	Standard Deviation	41420.51	Standard Deviation
Sample Variance	134782378397.25	Sample Variance	843533.68	Sample Variance	1715658774.18	Sample Variance
Kurtosis	34.59	Kurtosis	5.24	Kurtosis	285.08	Kurtosis
Skewness	4.02	Skewness	1.47	Skewness	13.06	Skewness
Range	7625000.00	Range	13250.00	Range	1650839.00	Range
Minimum	75000.00	Minimum	290.00	Minimum	520.00	Minimum
Maximum	7700000.00	Maximum	13540.00	Maximum	1651359.00	Maximum
Sum	11672925008.00	Sum	44952873.00	Sum	326506890.00	Sum
count	21613.00	count	21613.00	count	38652488.00	count
					6300385.00	
					21613.00	
					21613.00	

Price:

- Prices fluctuate widely, from \$75,000 to \$7,700,000, with a median of \$450,000.
- The data has a heavy-tailed distribution shown by its very high kurtosis of 34.59 and its sharp rightward skewness score of 4.02.

Sqft Living:

- Living spaces vary in size from 290 square feet to 13,540 square feet, with a median of 1,910 square feet.
- The data has a rather heavy-tailed distribution, as shown by the skewness value of 1.47 and the kurtosis value of 5.24.

Sqft Lot:

- Lot areas differ greatly in size, from 520 square feet to 1,651,359 square feet, with a median of 7,618 square feet.

□ With a skewness score of 13.06 and a very high kurtosis of 285.08, the data is strongly skewed to the right, suggesting an extremely heavy-tailed distribution.

Sqft above:

- Above-ground living spaces vary in size from 290 to 9,410 square feet, with a typical of 1,560 square feet.

- The data has a little heavy-tailed distribution indicated by a kurtosis value of 3.40 and a skewness value of 1.45, which is somewhat skewed to the right.

Sqft Basement:

- Basement spaces vary in size from 0 to 4,820 square feet, with a 0 square foot median.

- A heavy-tailed distribution is indicated by the data's strong rightward skewness value of 1.58 and kurtosis value of 2.72.

2. Develop a regression model to predict the price of houses in King County. What are the variables affecting the price? Be mindful of multicollinearity.

To solve this issue, we will use multiple regression analysis and conduct iterations. We will also exclude factors that are not statistically significant for determining the price. Specifically, we will remove variables with a p-value less than 0.05. Given that the dataset contains 19 variables, it is important to note that Excel can only accommodate 16 variables for regression analysis at a time. Therefore, we will first do the regression analysis using 16 variables. After identifying and eliminating the variables that do not significantly contribute to the analysis, we will then introduce more variables.

We have taken a dummy variable for the Year.

□ Run regression taking Price on the Y-axis and the rest of variables on the X-axis

	bedrooms	bathrooms	sqft_living	sqft_lot	floors	grade	sqft_above	sqft_basement	yr_built	renovate	zipcode	lat	long	sqft_living15	sqft_lot15	waterfront	View 0	View 1	View 2	View 3	Condition 2	Condition 3	Condition 4	price
bedrooms	1										t													
bathrooms	0.515884	1																						
sqft_living	0.579677	-0.754661																						
sqft_lot	0.507001	0.672285																						
floors	0.174242	0.500681	0.351949	-0.005201																				
grade	0.356967	0.664982	0.762794	0.113621	0.438183	1																		
sqft_above	0.4778	0.683342	0.826932	0.183512	0.523885	0.759232	1																	
sqft_basement	0.303093	0.283777	0.435943	0.015286	-0.245705	0.168392	-0.051941																	
yr_built	0.154178	0.506013	0.318048	0.053038	0.489319	0.446963	0.423898	-0.133124	1															
yr_renovate	0.018841	0.050739	0.055361	0.007644	0.006338	0.014414	0.023285	0.071324	-0.224874	1														
tipcode	-0.132668	-0.203886	-0.19943	-0.125974	-0.059212	-0.184862	-0.26112	0.074845	-0.346869	0.064357	1													
lat	-0.008931	0.024578	0.052529	-0.085681	0.049612	0.114084	-0.000816	0.110538	-0.148122	0.029398	0.267048	1												
long	0.291688	0.060851	0.048465	0.225602	0.195402	0.125412	0.195402	0.125412	0.195402	0.125412	0.195402	0.125412	1											
sqft_living15	0.391638	0.060851	0.075984	0.071851	-0.012989	0.113248	0.18405	0.017276	0.070958	0.007804	0.147221	0.086419	0.254531	0.183192	0.183192	1								
sqft_lot15	0.023244	0.087175	0.181209	0.718591	-0.012989	0.113248	0.18405	0.017276	0.070958	0.007804	0.147221	0.086419	0.254531	0.183192	0.183192	1								
waterfront	-0.00582	0.067344	0.065511	-0.067847	-0.022721	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	1							
View 0	-0.080106	-0.177445	-0.239907	-0.067847	-0.012989	-0.237327	0.013325	-0.27531	0.092021	0.092123	0.092205	0.030385	0.030385	0.030417	0.04191	0.058643	0.030703	1						
View 1	0.0222	0.038054	0.065511	-0.005287	0.021818	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	0.023698	1						
View 2	0.04506	0.087204	0.135285	0.037278	0.0059751	0.121919	0.077861	0.121919	0.121919	0.121904	0.044613	0.032559	0.052004	0.050905	-0.039607	0.139827	0.031523	0.019111	0.054541	0.031372	0.053391	0.047002	0.0269057	
View 3	0.050431	0.112296	0.158885	0.073871	-0.023721	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	0.048944	
Condition	-0.036346	-0.045491	-0.035069	0.006332	-0.023377	-0.058891	-0.028994	-0.01851	-0.050101	-0.01671	0.010148	0.046269	-0.012356	-0.010003	-0.03862	0.011115	0.004648	-0.008051	-0.005796	1				
Condition	-0.051395	-0.074741	-0.065324	0.037671	-0.055953	-0.087709	-0.058925	-0.023076	-0.067277	-0.008571	0.023615	-0.026265	0.014842	-0.023417	0.02224	-0.01789	0.019073	-0.006953	-0.01177	-0.007062	-0.003339	1		
Condition	0.004873	0.190548	0.102413	0.01432	0.018049	0.196592	0.194555	0.151498	0.391719	0.092968	0.042297	0.105384	0.113784	0.012694	0.021606	0.027743	-0.019334	0.017467	-0.024109	-0.050717	-0.121841	1		
Condition	-0.008931	-0.166347	-0.083794	0.013157	-0.257795	-0.139973	-0.142486	0.097212	-0.257454	-0.054813	-0.060003	-0.074781	-0.023625	0.008713	0.020401	0.008713	0.015184	0.007624	0.011354	-0.022258	-0.035471	-0.812331	1	
price	0.30883	0.525138	0.703039	0.089661	0.674348	0.605569	0.323816	0.054012	0.126434	-0.053203	0.307003	0.023447	0.266369	-0.35912	0.092607	0.140418	0.152088	-0.020885	-0.051917	-0.071311	-0.030715	1		

Regression 1:

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
intercept	6091985.51							
year dummy	26736.17	3125725.73	1.95	0.051	-34667.70	12218638.71	-34667.70	12218638.71
bedrooms	-38884.03	3150.23	8.49	0.000	20561.50	32910.84	20561.50	32910.84
bathrooms	46557.67	2027.00	-19.18	0.000	-42857.10	-34910.97	-42857.10	-34910.97
sqft_living	166.77	3489.21	13.34	0.000	39718.56	53396.79	39718.56	53396.79
sqft_lot	-28807.48	4.69	-7.00	0.000	-0.33	-0.18	-0.33	-0.18
floors	578300.86	3830.67	7.52	0.000	21299.09	36315.88	21299.09	36315.88
waterfront	42893.91	18627.46	31.05	0.000	541789.66	614812.05	541789.66	614812.05
view	20211.36	2284.78	18.77	0.000	38415.57	47372.26	38415.57	47372.26
condition	120698.57	2523.27	8.01	0.000	15265.57	25157.15	15265.57	25157.15
grade		2251.72	53.60	0.000	116285.03	125112.10	116285.03	125112.10
sqft_above	0.00	4.58	-1.69	0.092	-16.70	1.26	-16.70	1.26
sqft_basement	-3597.30	0.00	65535.00	#NUM!	0.00	0.00	0.00	0.00
yr_built	10.44	74.52	-48.28	#NUM!	-3743			

sqft_lot	-0.00523	0.05120	-0.10212	0.91866	-0.10559	0.09513	-0.10559	0.09513
floors	27510.75015	3778.36534	7.28112	0.00000	20104.87511	34916.62519	20104.87511	34916.62519
waterfront	579670.81318	18606.11951	31.15485	0.00000	543201.44519	616140.18116	543201.44519	616140.18116
view	42887.67636	2269.82393	18.89472	0.00000	38438.65386	47336.69885	38438.65386	47336.69885
condition	20604.23397	2496.03805	8.25478	0.00000	15711.81510	25496.65284	15711.81510	25496.65284
grade	119983.48738	2245.68755	53.42840	0.00000	115581.77398	124385.20078	115581.77398	124385.20078
sqft_above	-0.52942	4.55709	-1.45912	0.15013	-15.42245	2.36361	-15.42245	2.36361
sqft_basement	0.00000	0.00000	65535.00000	#NUM!	0.00000	0.00000	0.00000	0.00000
yr_built	3572.51728	70.87307	-50.40726	#NUM!	3711.43372	3433.60083	-3711.43372	-3433.60083
yr_renovated	10.97154	3.90861	2.80702	0.00500	3.31037	18.63272	3.31037	18.63272
sqft_living15	24.97683	3.59452	6.94859	0.00000	17.93132	32.02235	17.93132	32.02235
sqft_lot15	-0.54800	0.07823	-7.00467	0.00000	-0.70134	-0.39466	-0.70134	-0.39466

- The variable "sqft_lot" will be eliminated first because of its highest p-value.

Regression 3

	Coefficients 6184509.4599 26649.6554	Standard Error 138191.6053	t Stat 44.7531 8.4697	P-value 0.0000 0.0000	Lower 95% 5913643.7110 20482.3337	Upper 95% 6455375.2087 32816.9771	Lower 95.0% 5913643.7110	Upper 95.0%
intercept	-39506.9075	3146.4742	-19.5308	0.0000	-43471.7412	-35542.0739	20482.3337	6455375.2087
year dummy	45851.1413	2022.7981	13.1511	0.0000	39017.3597	52684.9229	-43471.7412	32816.9771
bedrooms	167.3886	3486.4919	35.9230	0.0000	158.2553	176.5218	39017.3597	-35542.0739
bathrooms	27524.9503	4.6596	7.2900	0.0000	20124.2621	34925.6384	158.2553	52684.9229
sqft_living	579704.5573	3775.7191	31.1623	0.0000	543241.7767	616167.3378	176.5218	
floors	42880.6848	18602.7587	18.9007	0.0000	38433.7886	47327.5809	20124.2621	34925.6384
waterfront	20607.1803	2268.7391	8.2567	0.0000	15715.2004	25499.1602	543241.7767	616167.3378
view							38433.7886	
condition	119982.9150	2495.8141	53.4295	0.0000	115581.3162	124384.5139	47327.5809	
grade		2245.6291					15715.2004	25499.1602
sqft_above	0.0000		65535.00		0.0000	0.0000	115581.3162	124384.5139
sqft_basement	-3572.4247	0.0000	00		-3711.3265	-3433.5228	15715.2004	2.3267
yr_built	10.9747	70.8656	-50.4112	#NUM!	3.3139	18.6355	0.0000	0.0000
	24.9961		2.8080	#NUM!	17.9604	32.0317	-3711.3265	
yr_renovated	-0.5536	3.9084	6.9637		-0.6629	-0.4443	3.3139	-3433.5228
sqft_living15		3.5895	-9.9251	0.0050			17.9604	18.6355
sqft_tot15		0.0558		0.0000			-0.6629	32.0317
				0.0000				-0.4443

- We are removing the variable “sqft_above” first because it has the highest p-value.

- Regression 4:

	Coefficients	Standard Error	t Stat	Pvalue	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
--	--------------	----------------	--------	--------	-----------	-----------	-------------	-------------

intercept	6184509.4599	138191.6053	44.7531	0.0000	5913643.7110	6455375.2087	5913643.7110	6455375.2087
year dummy	26649.6554	3146.4742	8.4697	0.0000	20482.3337	32816.9771	20482.3337	32816.9771
bedrooms	-39506.9075	2022.7981	-	0.0000	-43471.7412	-35542.0739	-43471.7412	-35542.0739
	45851.1413	3486.4919	19.5308	0.0000	39017.3597	52684.9229	39017.3597	
bathrooms	160.8348	3.8840	13.1511	0.0000	153.2219	168.4477	153.2219	52684.9229
sqft_living	27524.9503	3775.7191	41.4097	0.0000	20124.2621	34925.6384	20124.2621	168.4477
floors	579704.5573	18602.7587	7.2900	0.0000	543241.7767	616167.3378	543241.7767	34925.6384
waterfront	42880.6848	2268.7391	31.1623	0.0000	38433.7886	47327.5809	38433.7886	616167.3378
view	20607.1803	2495.8141	18.9007	0.0000	15715.2004	25499.1602	15715.2004	47327.5809
condition	119982.9150	2245.6291	8.2567	0.0000	115581.3162	124384.5139	115581.3162	25499.1602
grade			53.4295					124384.5139
sqft_basement	-3572.4247	70.8656	1.4465	0.0000	-3711.3265	-3433.5228	-3711.3265	15.4343
yr_built	10.9747	3.9084	-	0.0050	3.3139	18.6355	3.3139	-3433.5228
	24.9961	3.5895	50.4112	0.0000	17.9604	32.0317	17.9604	
yr_renovated	-0.5536	0.0558	2.8080	0.0000	-0.6629	-0.4443	-0.6629	18.6355
sqft_living15			6.9637					32.0317
sqft_lot15			-9.9251					-0.4443

- Now, we are removing the variable “sqft_basement” first as it has the highest p-value

• Regression 5:

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
intercept	6205923.2142	137399.8938	45.1669	0.0000	20445.5176	32780.0683	5936609.2790	6475237.1493
year dummy	26612.7929	3146.4505	8.4580	0.0000	-43469.3138	-35539.4489	20445.5176	32780.0683
bedrooms	-39504.3814	2022.8485	-19.5291	0.0000	39939.5062	53421.3679	-43469.3138	-35539.4489
bathrooms	46680.4370	3439.1209	13.5734	0.0000			39939.5062	53421.3679
sqft_living	163.0538	3.5683	45.6950	0.0000	156.0596	170.0479	156.0596	170.0479
floors	25242.6274	3430.3400	7.3586	0.0000	18518.9077	31966.3470	18518.9077	31966.3470
waterfront	578740.0249	18591.2750	31.1297	0.0000	542299.7534	615180.2963	542299.7534	615180.2963
view	43477.3392	2230.9854	19.4880	0.0000	39104.4432	47850.2352	39104.4432	47850.2352
condition	20819.5345	2491.5560	8.3560	0.0000	15935.9008	25703.1681	15935.9008	25703.1681
grade	119668.4649	2235.1393	53.5396	0.0000	115287.4268	124049.5030	115287.4268	124049.5030
yr_built	-3581.9591	70.5602	-50.7646	0.0000	-3720.2624	-3443.6559	-3720.2624	-3443.6559
yr_renovated	10.9248	3.9083	2.7952	0.0052	3.2641	18.5854	3.2641	18.5854
sqft_living15	23.9884	3.5213	6.8123	0.0000	17.0863	30.8905	17.0863	30.8905
sqft_lot15	-0.5621	0.0555	-10.1342	0.0000	-0.6708	-0.4534	-0.6708	-0.4534

- All the above variables are significant as the p values are 0.05.

- R square on Regression 1

- R square on Regression 5

Regression Statistics		Regression Statistics	
Multiple R	0.808760912	Multiple R	0.80922472
R Square	0.654094212	R Square	0.654844647
Adjusted R Square	0.653807664	Adjusted R Square	0.654636905
Standard Error	215996.263	Standard Error	215751.8466
Observations	21613	Observations	21613

Additionally, it is evident that the R square value has risen, indicating that the model has improved in accuracy after eliminating the variables that were not statistically significant.

3. Test the following hypotheses and provide your conclusion (bonus)

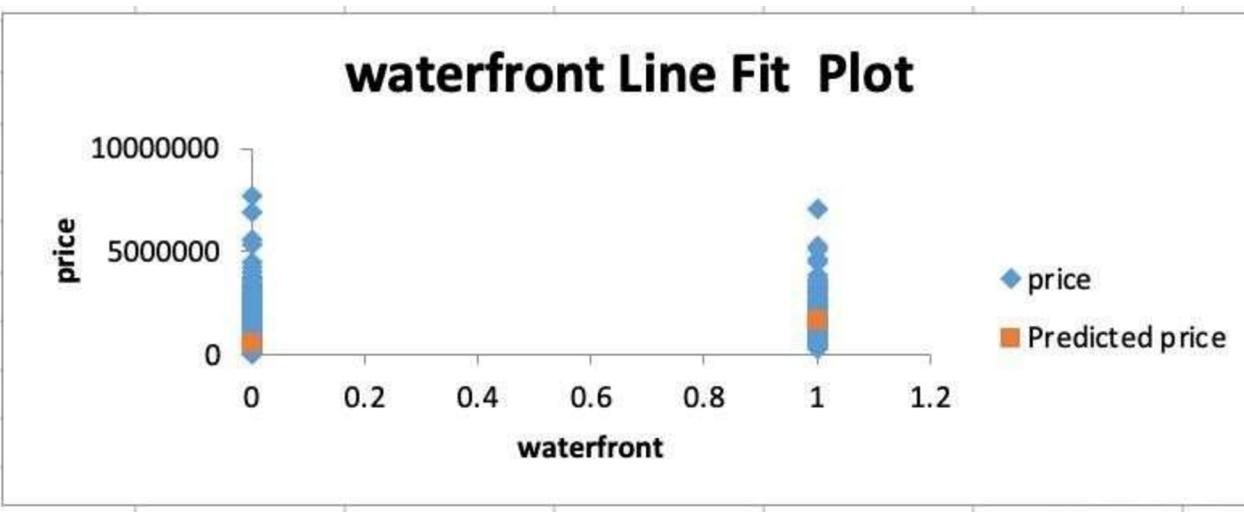
- a. Average price of houses with waterfront are higher than those without a waterfront.

For this, running regression on Price and waterfront, taking price on the Y axis and waterfront on the X axis.

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	531563.600	2416.194	220.000	0.000	526827.681	536299.519	526827.681	536299.519
waterfront	1130312.425	27822.465	40.626	0.000	1075778.342	1184846.508	1075778.342	1184846.508

By the above regression, it is evident that the price of residences increases by 1130312.425 units when one unit of waterfront is added.

To get a more comprehensive clarification, see the graph provided below.



As we can see from the graph above, “0” denotes the houses with no waterfront and “1” with houses with waterfront.

The predicted price for houses with a waterfront is higher than the houses without a waterfront.

b. Older houses have lower prices. (Create the “age” variable with respect to 2014 and 2015 using yr_built data)

For above question, A new variable, named "age," was formed by subtracting the year from the year_built, resulting in the determination of the house's "Age."

price	Year	yr_built	Age
221900	2014	1955	59
538000	2014	1951	63
180000	2015	1933	82
604000	2014	1965	49
510000	2015	1987	28
1.23E+06	2014	2001	13
257500	2014	1995	19
291850	2015	1963	52
229500	2015	1960	55

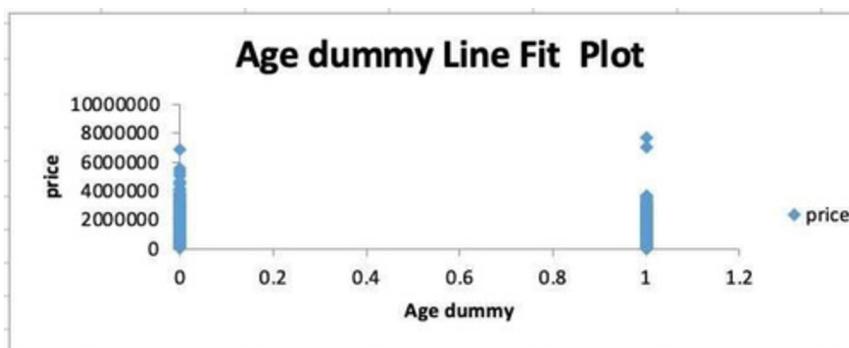
Subsequently, we computed the Mean of the column labeled "Age," yielding an average value of 43. Subsequently, we generated a dummy variable in which buildings with an age beyond 43 were assigned a value of 1 (representing old), whereas dwellings with an age below 43 were assigned a value of 0 (representing new).

Upon running regression, we got the following results:

	Coefficients	Standard Error	t Stat	P-value	Lower 95%	Upper 95%	Lower 95.0%	Upper 95.0%
Intercept	571124.4609	3426.3261	166.6871	0.0000	564408.6090	577840.3128	564408.6090	577840.3128
Age dummy	-65615.5694	4981.9233	-13.1707	0.0000	-75380.5066	-55850.6322	-75380.5066	-55850.6322

The data above indicates that for each additional year of the "Age dummy" variable, the price of the property decreases by \$65615.5694.

For more clarification, we can see the graph below:



The prices of old houses (represented by dummy variable 1) are much lower than the prices of new houses (represented by dummy variable 0).