

# First Project - EDA

**This is your first project during your bootcamp. You'll be working with the King County House Sales dataset. Here, the focus is on EDA though you are required to demonstrate an entire Data Science Lifecycle.**

## The data

- The dataset can be found in the file "King\_County\_House\_prices\_dataset.csv", in this folder.
- The description of the column names can be found in the column\_names.md file in this repository.
- The column names are NOT clear at times.

*In the real world we will run into similar challenges. We would then go ask our business stakeholders for more information. In this case, let us assume our business stakeholder who would give us information, left the company. Meaning we would have to identify and look up what each column names might actually mean.*

## Tasks for you

1. Create new repo, and new conda environment.
2. Through statistical analysis/EDA, above please come up with AT LEAST 3 recommendations for home sellers and/or buyers in King County.

If you use linear regression in the exploration phase remember that  $R^2$  close to 1 is good.

3. Then model this dataset with a multivariate linear regression to predict the sale price of houses as accurately as possible.
  - a. Split the dataset into a train and a test set. (use the sklearn split method  
[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.train\\_test\\_split.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html) )
  - b. Use Root Mean Squared Error (RMSE) as your metric of success and try to minimize this score on your test data.

# First Project - EDA

## The Deliverables

1. A well **documented Jupyter Notebook** containing any code you've written for this project and comments explaining it. This work will need to be pushed to your GitHub repository in order to submit your project (latest upload: 15.10.2020 11:00). Do not push all the analysis.. just the analysis that is relevant..
2. A python script for training the model, printing out the model statistics and saving the model.
3. An **organized README.md** file in the GitHub repository that describes the contents of the repository. This file should be the source of information for navigating through the repository.
4. A **short Keynote/PowerPoint/Google Slides/Jupyter slides presentation** giving a high-level overview of your methodology and recommendations for non-technical stakeholders. The duration of the presentation should be **10 minutes**, then the discussion will continue for 5 minutes. Also put your slides (delivered as a PDF export) on Github to get a well-rounded project.