# Creating a Machine Learning Model to Find out What Drives Compensation
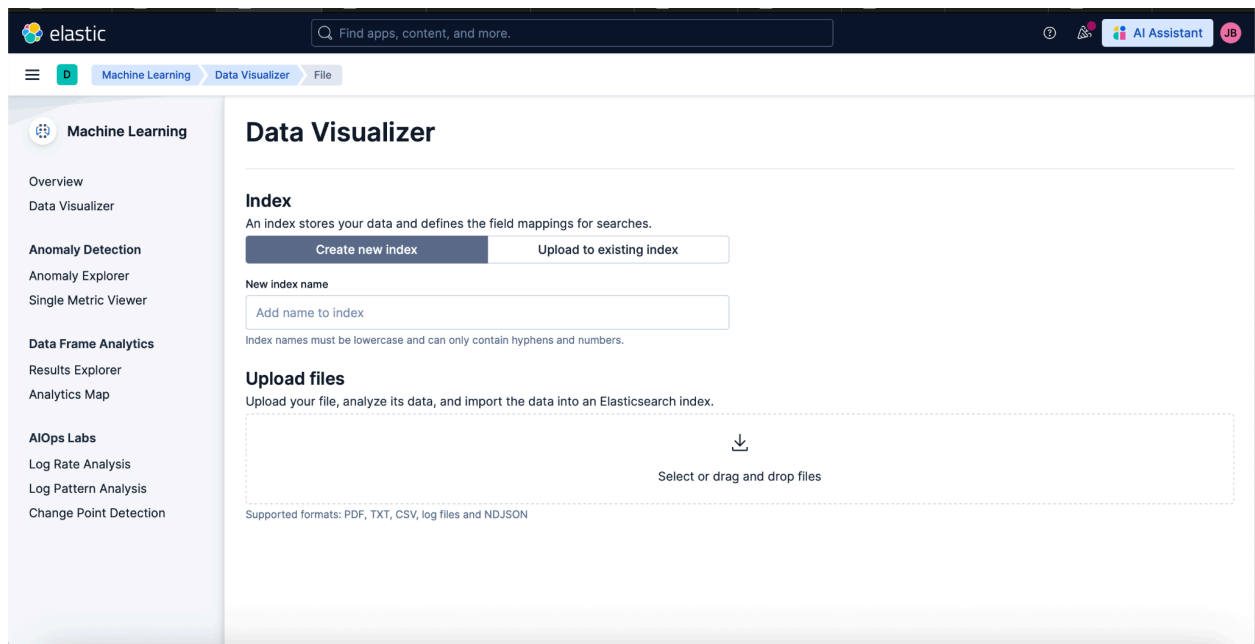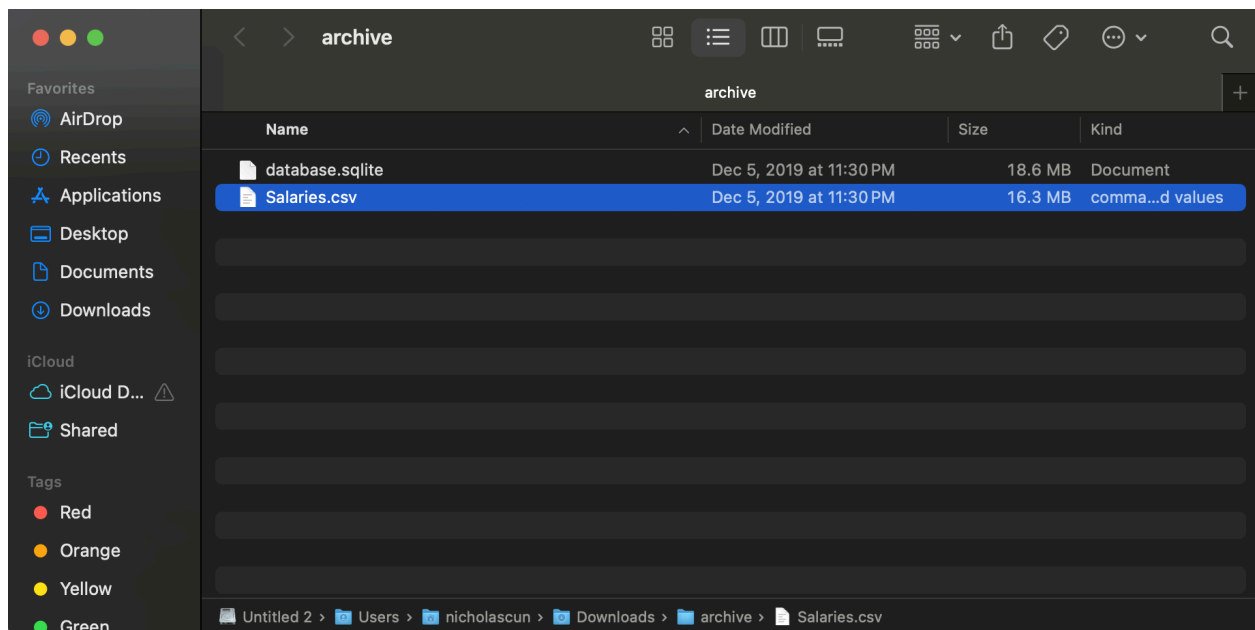
**Step 1:** First, go to this link and download the salaries dataset
https://www.kaggle.com/datasets/kaggle/sf-salaries

**Step 2:** Then go to machine learning / data visualizer to upload the dataset that was downloaded. I recommend typing this in the search bar to pull up the path immediately without having to search and scroll around for it.

Make sure to choose the .csv file



**Step 3:** Then go to Dev Tools and run these two functions in order to include job title in the ML model. Having the original JobTitle in it's text form will produce an error in the ML model, so we need to make a keyword version of it to use instead.

```
PUT sf_salaries_fixed
{
 "mappings": {
   "properties": {
     "JobTitle": {
       "type": "text",
       "fields": {
         "keyword": {
           "type": "keyword"
         }
       }
     }
   }
 }
}
```

```
POST _reindex
{
  "source": {
    "index": "sf_salaries"
  },
  "dest": {
    "index": "sf_salaries_fixed"
  }
}
```

**Step 4:** Create a data frame analytics regression model. In the search bar go to Machine Learning / Data Frame Analytics Jobs. Then choose regression.



**Step 5:** Select TotalPayBenefits as the Dependant variable.

**Step 6:** Uncheck id, EmployeeName, and Benefits.keyword from the "is included" section. This will make it easier on the ML model and safer since it has less unnecessary variables to work with. Don't forget to uncheck JobTitle aswell since it will give you an error if you do not. We will be using JobTitle.keyword instead.





**Step 7:** Set the training percent to 90

**Step 8:** Start the job and wait for all the phases to finish up.

**Step 9:** Once all phases are completed and the model has stopped running, go to Data Frame analytics / Results Explorer to view the results of the ML model.





**\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*This concludes the SF Salaries Tutorial\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\***