# Assignment : Classification

- I have followed an LLM based approach.
- I have played with various BERT based models and tried to generate inference using zero shot, many shot & model finetuning.
- Models used → DistilBert, BioBert, BioClinicalBert & bert-based-uncased
- There are a total of 40 classes.
- Another approach which I thought of was using Sentence Transformers to generate embedding for the transcription, we can then use the train set. Now for each sentence in the test set we will again generate the embedding and use KNN method to determine the medical_speciality for the given test transcription.
- Classical ML based algos can also be applied but given the explosion of LLMs and their vast understanding of the world, I have tried this approach.
- Statistical approaches such as co-occurrence matrices can be used. We could further apply SVD to generate dense representations. But this approach however wil not be feasible due to large vocab size.
- We can also merge the label set i.e. medical_speciality, but it will rely on the medical domain knowledge.

Data Exploration:

Consult - History and Phy.

Cardiovascular / Pulmonary

Orthopedic

Radiology

General Medicine

Gastroenterology

Neurology

SOAP / Chart / Progress Notes

Urology

Obstetrics / Gynecology

Discharge Summary

ENT - Otolaryngology

Neurosurgery

Hematology - Oncology

Ophthalmology

Nephrology

Emergency Room Reports

Pediatrics - Neonatal

Pain Management

Psychiatry / Psychology

Office Notes

Podiatry

Dermatology

Cosmetic / Plastic Surgery

Dentistry

Letters

Physical Medicine - Rehab

Sleep Medicine

Endocrinology

Bariatrics

IME-QME-Work Comp etc.

Chiropractic

Rheumatology

Speech - Language

Autopsy

Diets and Nutritions

Lab Medicine - Pathology

Hospice - Palliative Care

Surgery

21.9%
10.4%
7.5%
7.1%
5.5%
5.2%
4.5%
4.5%
3.3%
3.1%
3.1%
2.2%
1.9%
1.9%
1.8%
1.7%
1.6%
1.5%
1.4%
1.2%
1.1%
1.0%
0.9%
0.6%
0.5%
0.5%
0.4%
0.4%
0.3%
0.3%
0.2%
0.2%
0.1%
0.1%
0.1%
0.1%
0.1%
0.1%