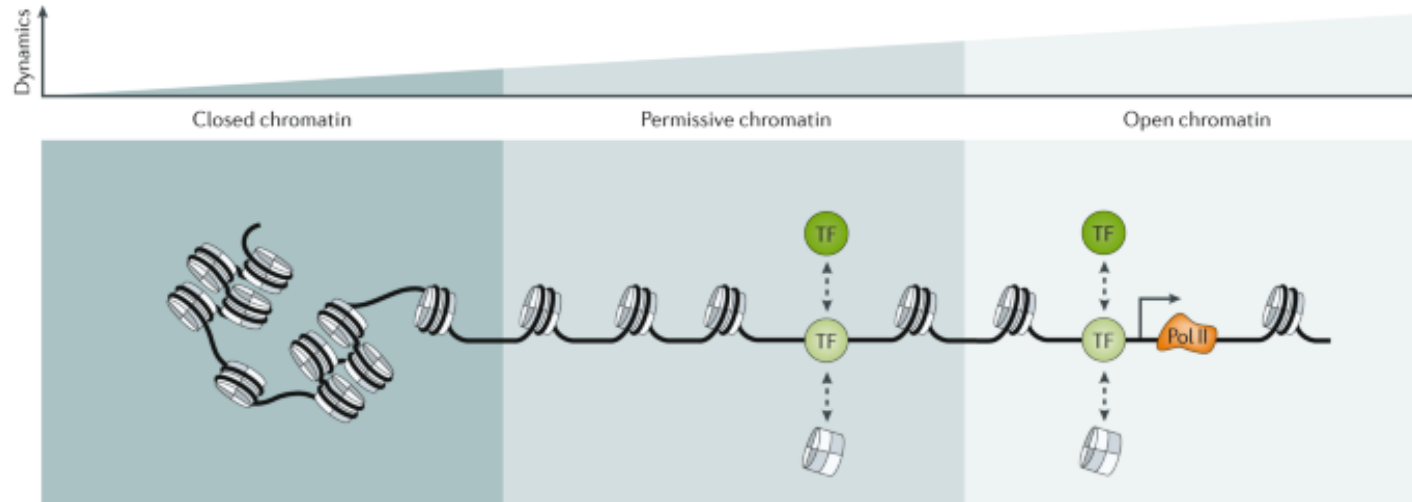


An interpretable classification of OCRs

Justin Bellavance

What's the problem?

(ATGACTAGCTACAGTGTACGA) -> does chromatin hide it?

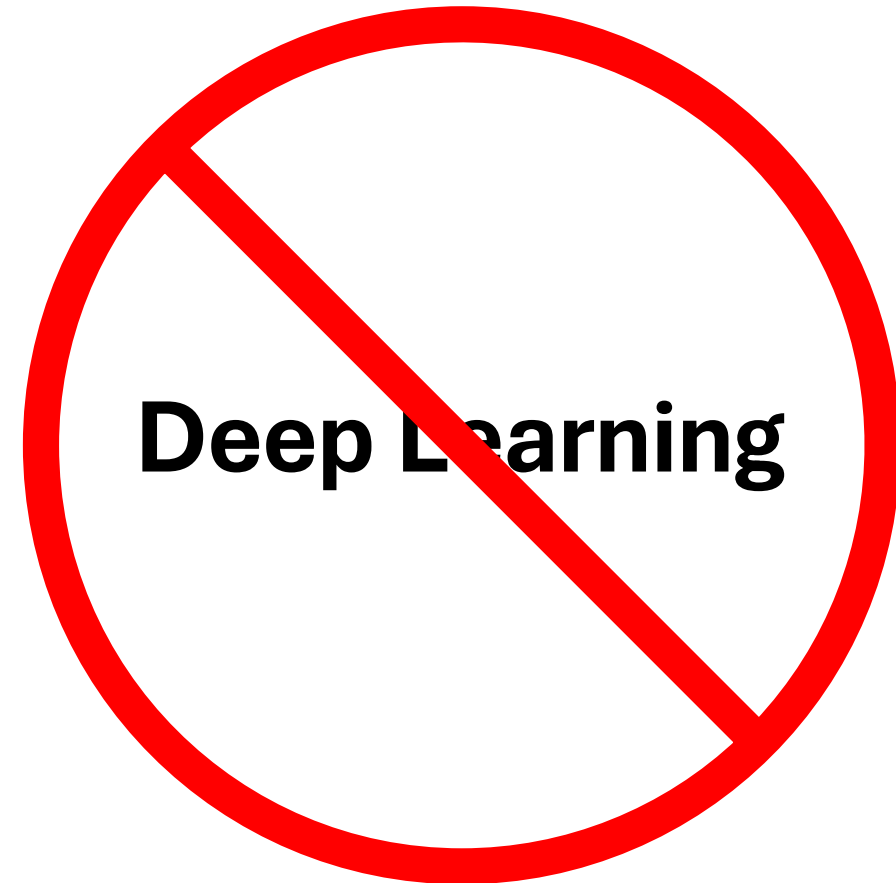


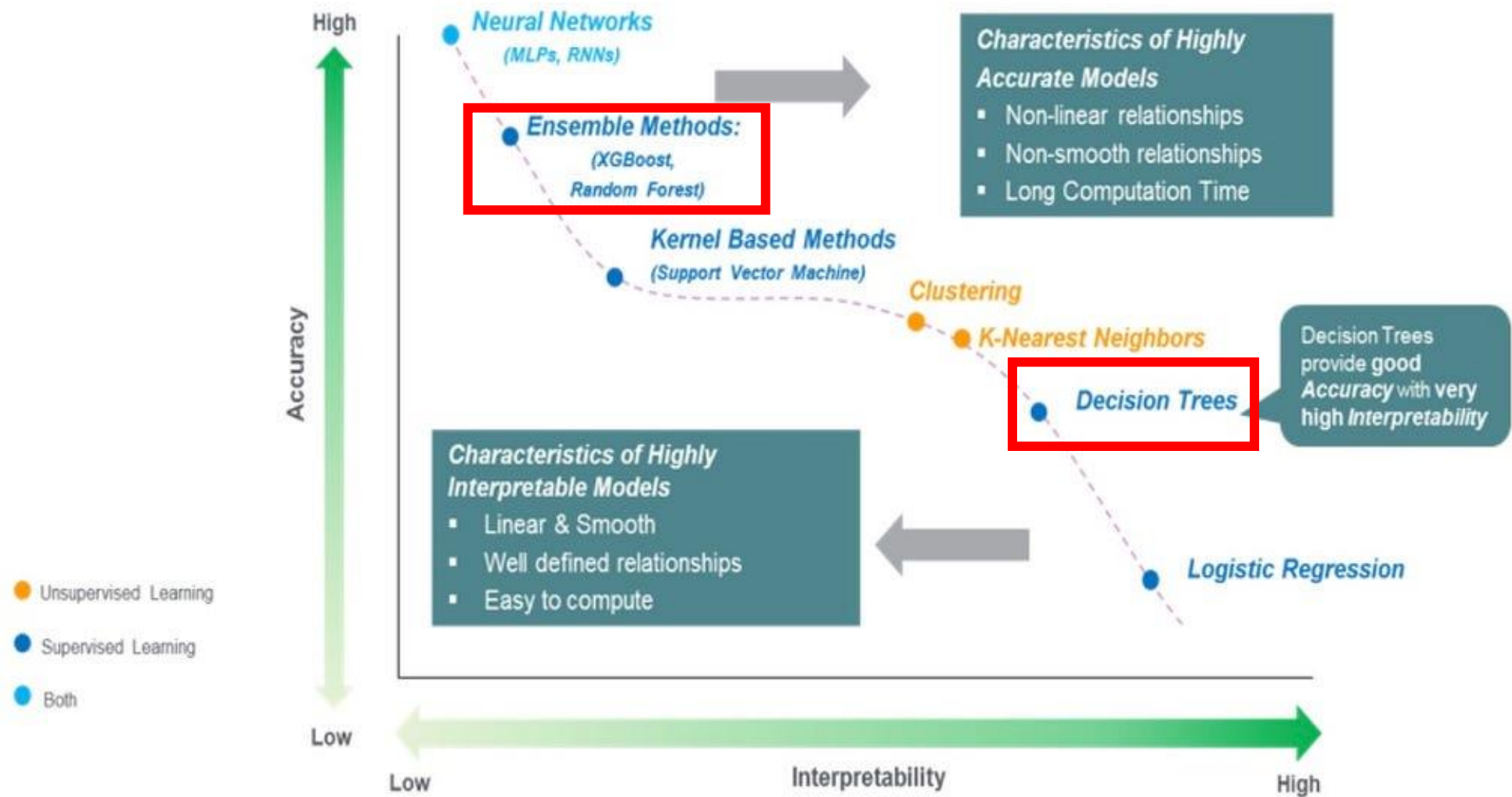
Very limited information based solely on DNA sequence

A lot of supplementary information available... but takes time to call

Deep learning can easily pick up on 'hidden' patterns without external help

- Perfect tool!





Also I have no GPU :(

Goal: Beat Deep Learning Benchmark (not SOTA)

Method	Accuracy	F1
Pytorch CNN	68	66.1
TensorFlow CNN	68.8	72

https://github.com/ML-Bioinfo-CEITEC/genomic_benchmarks/blob/main/experiments/README.md#hyenadna

Methods (run 1)

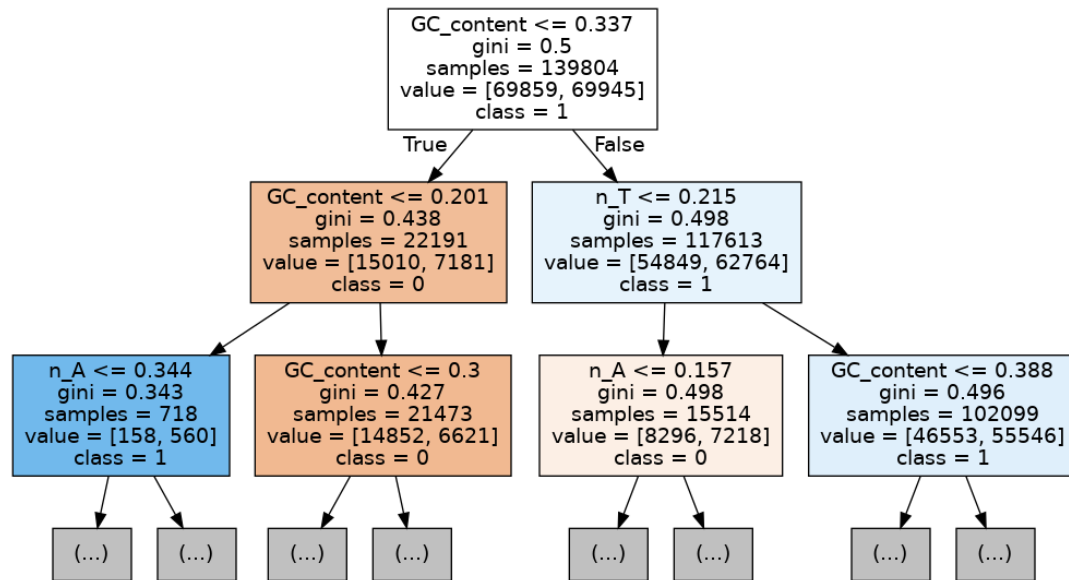
Try to think of all information I can possibly get from 4 letter alphabet DNA strands:

- A, C, T, G content
- Strand length
- CG%
- N content
- Contains start codon (ATG)
- Contains stop codon (TAG|TAA|TGA)

Use Decision tree then XGBoostClassifier (decision forest)

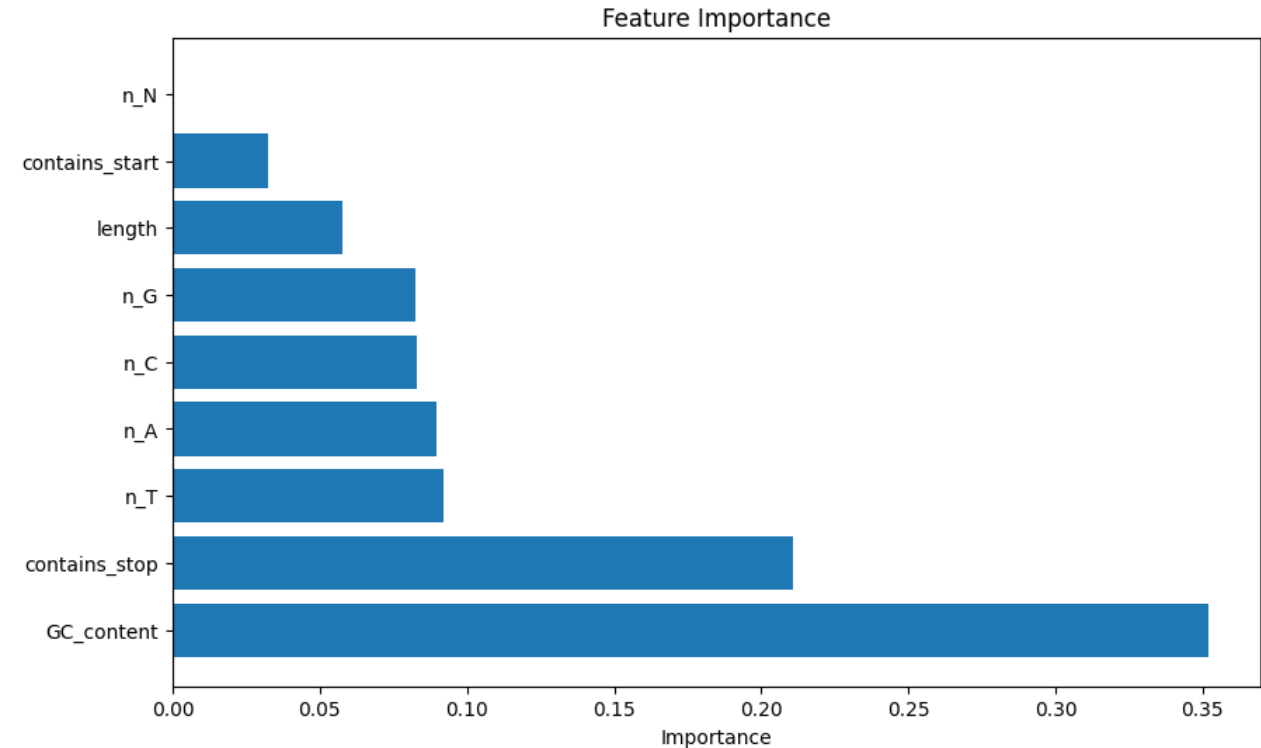
Results (Run 1)

Decision Tree



58.9% accuracy

XGBoost



60.5% accuracy

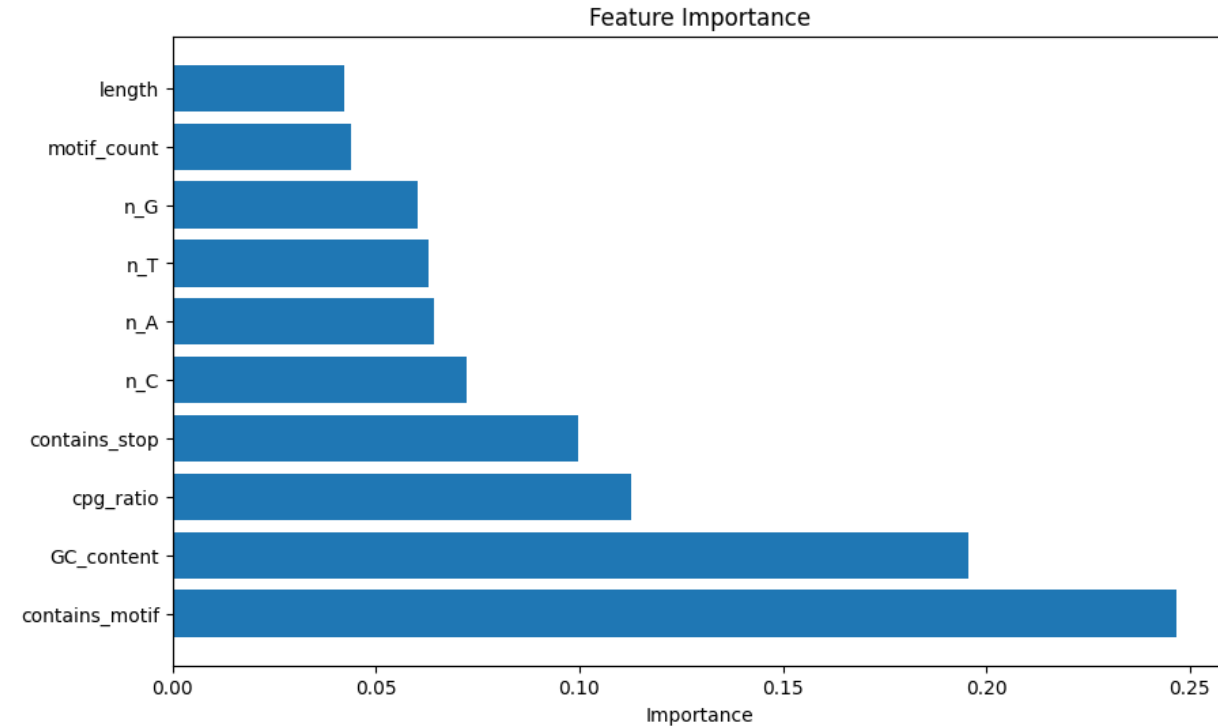
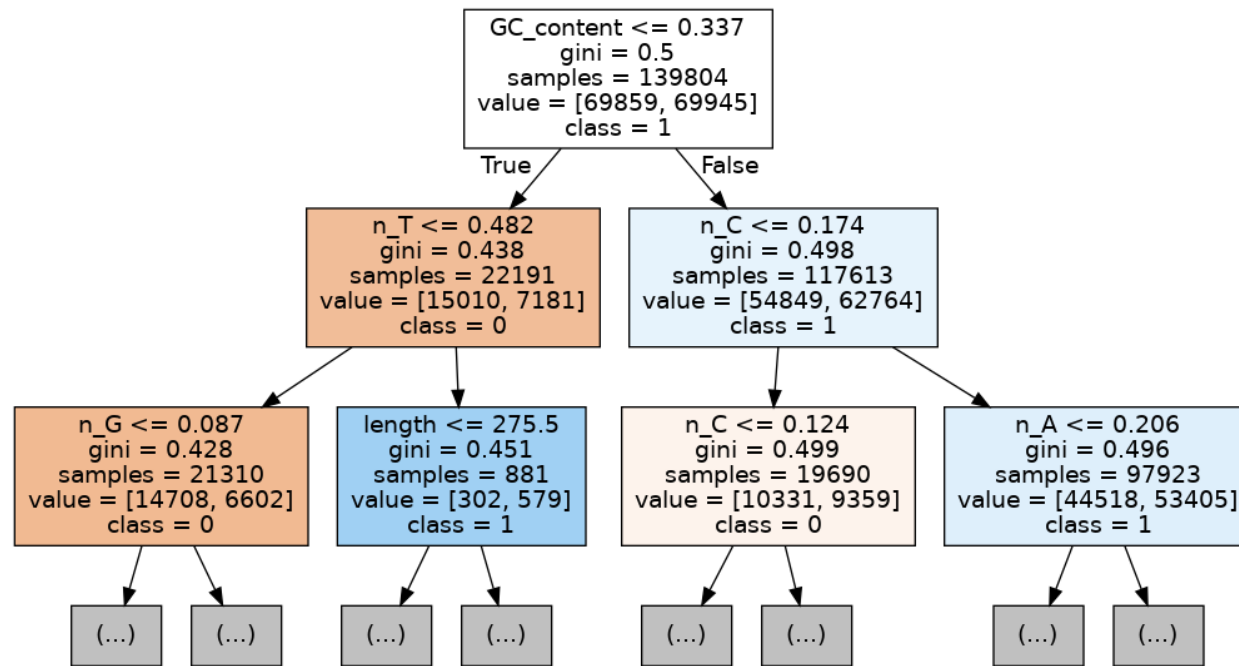
Methods (after feature engineering)

Parameters Now

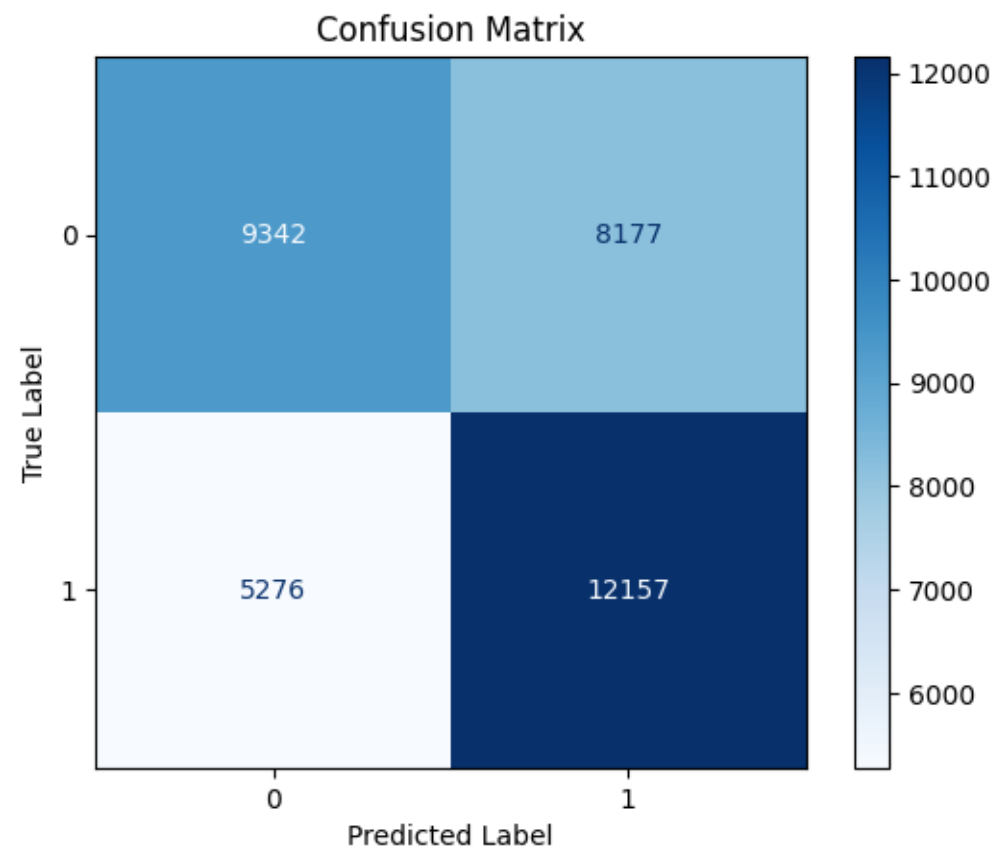
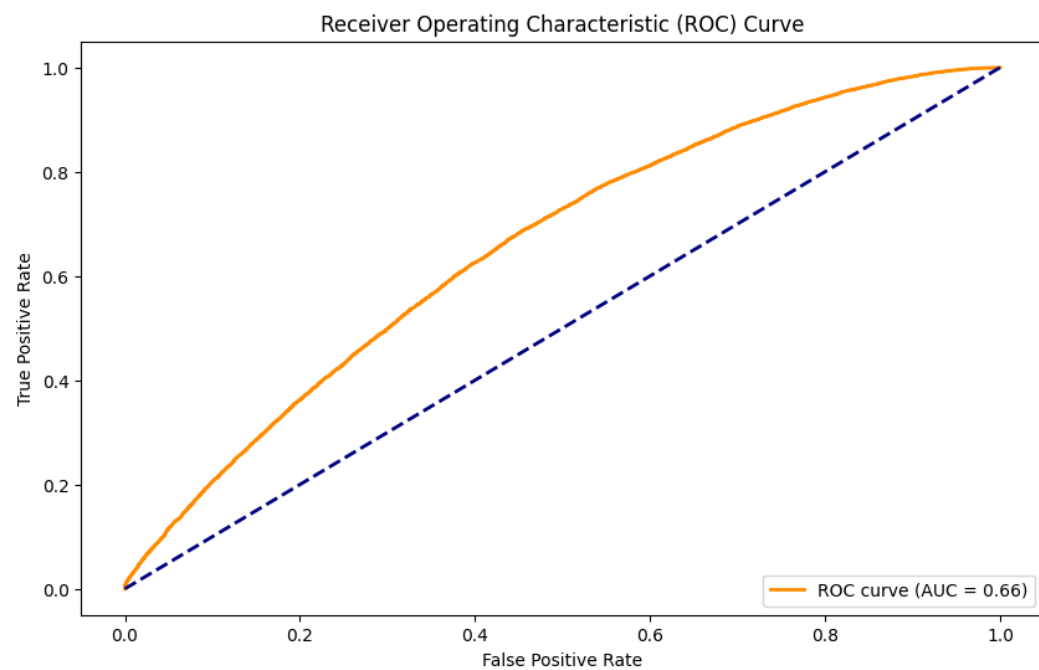
- GC content
- CpG islands
- Length
- # of TF binding motifs (from JASPAR)
 - (99% contained a motif)

BayesSearchCV for hyperparameter tuning

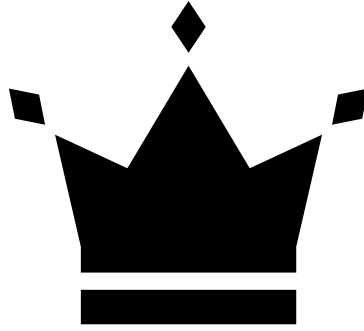
result



result



Method	Accuracy	F1
Pytorch CNN	68	66.1
TensorFlow CNN	68.8	72
XGBoostClassifier	61.5	NA



Deep Learning