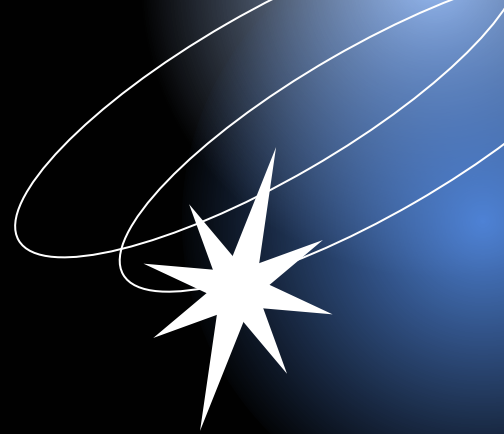


DNA Sequence Analysis

Jordan Dube, Luke Unterman, Blue
Arevalo, Kiersten Kofi Adomfrimpong



Motivation

- Generated multiple datasets from Pfeature for our model
- Noticed that average MCC values were still low as a result of an overrepresentation of nonDRNA data
- Used SMOTE oversampling operator as an attempt to correct underrepresentation of minority classes
- Refined model to reduce overfitting

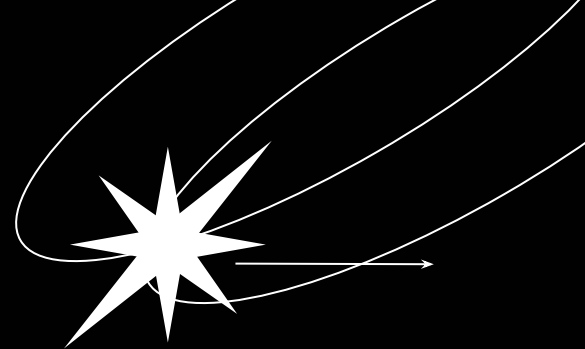
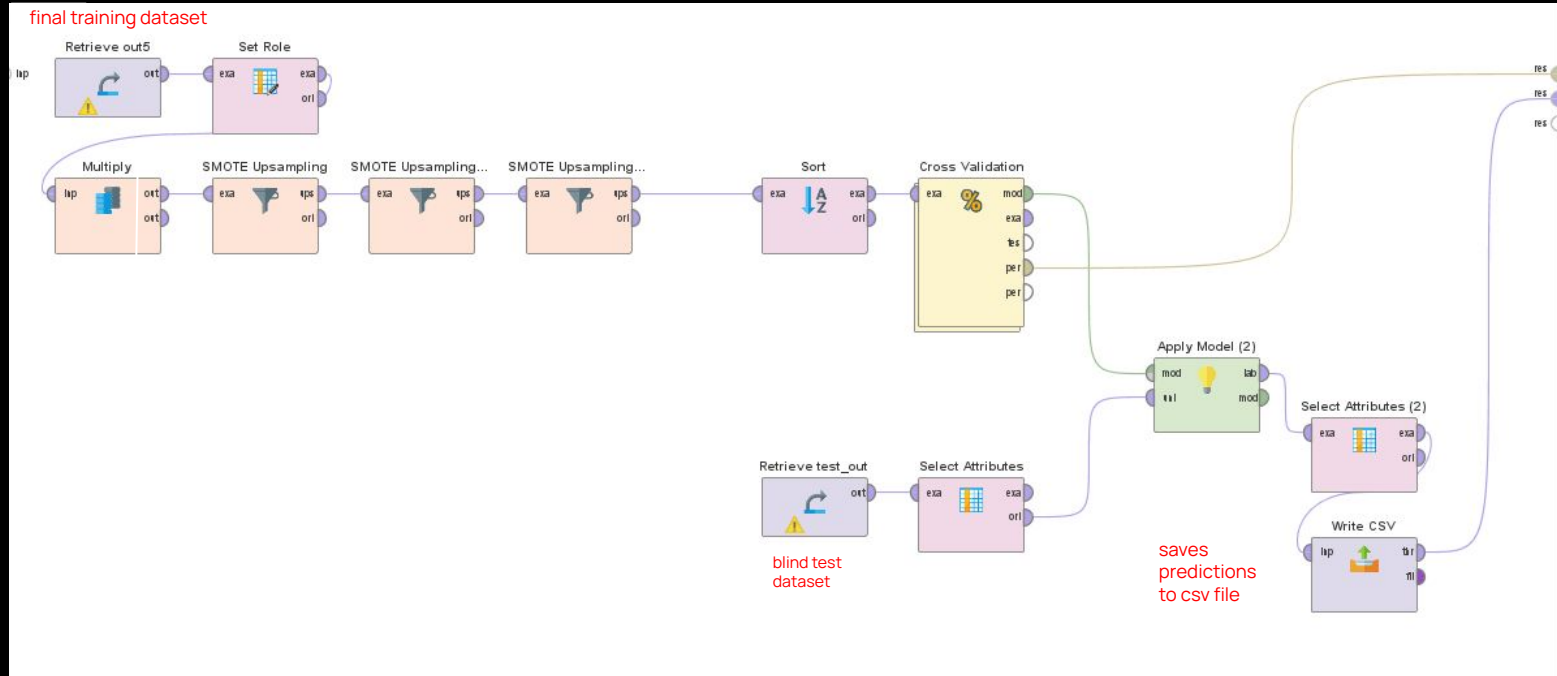


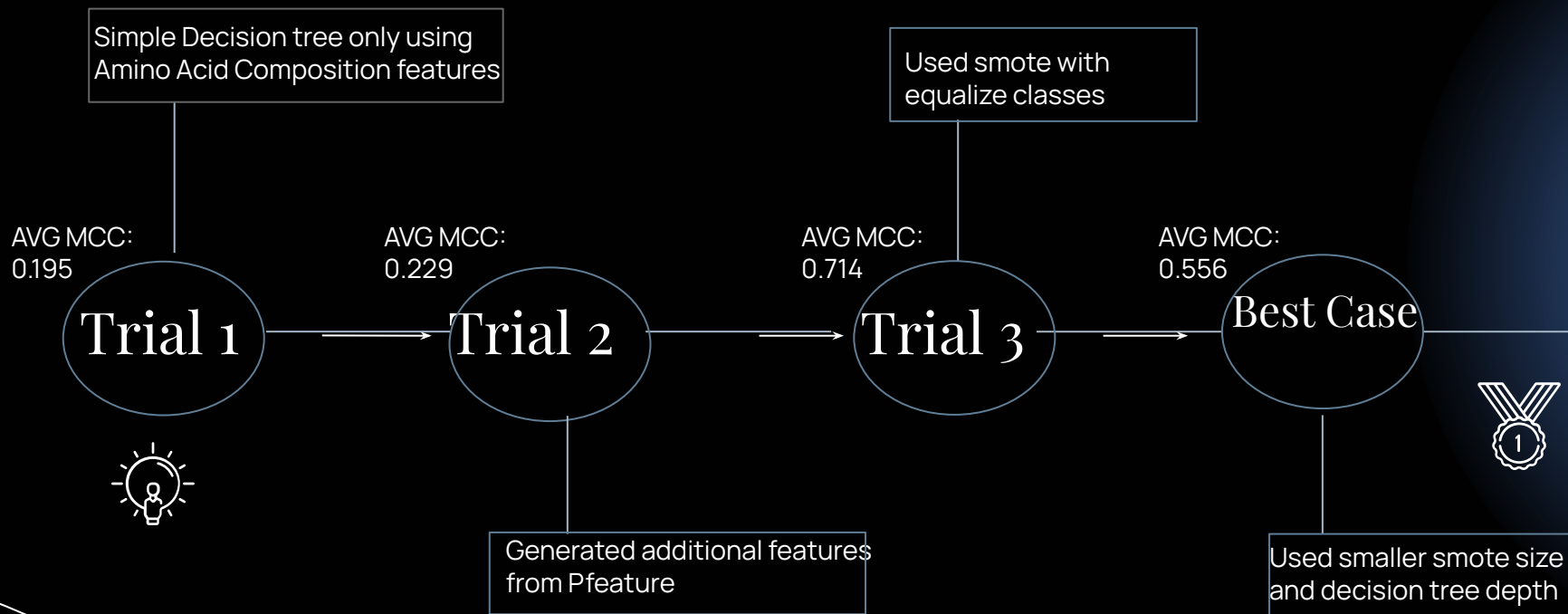
Diagram of Best Model

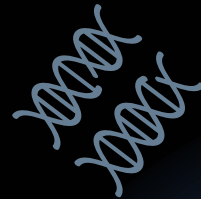


Final Parameters

| SMOTE Upsampling | Decision Tree |
|------------------------------------------------------|------------------------------------------------------------------|
| number of neighbours <input type="text" value="5"/> | criterion <input type="text" value="information_g..."/> |
| <input checked="" type="checkbox"/> normalize | maximal depth <input type="text" value="10"/> |
| <input type="checkbox"/> equalize classes | <input checked="" type="checkbox"/> apply pruning |
| upsampling size <input type="text" value="1000"/> | confidence <input type="text" value="0.1"/> |
| <input type="checkbox"/> auto detect minority class | <input checked="" type="checkbox"/> apply prepruning |
| minority class <input type="text" value="DNA"/> | minimal gain <input type="text" value="0.01"/> |
| <input type="checkbox"/> round integers | minimal leaf size <input type="text" value="2"/> |
| nominal change rate <input type="text" value="0.5"/> | minimal size for split <input type="text" value="4"/> |
| <input type="checkbox"/> use local random seed | number of prepruning alternatives <input type="text" value="4"/> |

Changes in iterations





Results

| Class | Absolute Count (test set, training set) | Representation in Dataset (test set, training set) |
|---------|--------------------------------------------|-------------------------------------------------------|
| nonDRNA | 7558, 7859 | 85.9%, 89.4% |
| RNA | 594, 523 | 6.8%, 5.9% |
| DNA | 542, 391 | 6.2%, 4.4% |
| DRNA | 100, 22 | 1.1%, 0.3% |

| Outcome | Quality measure | Baseline result | Design 1 | Design 2 | Design 3 | Best Design |
|--------------------------|--------------------|-----------------|---------------|---------------|--------------|--------------|
| DNA | <i>Sensitivity</i> | 6.9% | 16.4% | 17.6% | 93.9% | 57.3% |
| | <i>Specificity</i> | 99.3% | 96.4% | 96.4% | 96.0% | 97.0% |
| | <i>Accuracy</i> | 95.2% | 92.8% | 92.9% | 95.9% | 95.3% |
| | <i>MCC</i> | 0.132 | 0.131 | 0.143 | 0.681 | 0.496 |
| RNA | <i>Sensitivity</i> | 39.6% | 34.6% | 42.6% | 94.5% | 71.1% |
| | <i>Specificity</i> | 98.9% | 96.9% | 97.1% | 97.7% | 98.4% |
| | <i>Accuracy</i> | 95.3% | 93.2% | 93.9% | 97.5% | 96.8% |
| | <i>MCC</i> | 0.501 | 0.343 | 0.421 | 0.814 | 0.706 |
| DRNA | <i>Sensitivity</i> | 4.5% | 0.0% | 0.0% | 81.8% | 59.1% |
| | <i>Specificity</i> | 100% | 99.8% | 99.9% | 99.8% | 99.6% |
| | <i>Accuracy</i> | 99.7% | 99.6% | 99.6% | 99.7% | 99.5% |
| | <i>MCC</i> | 0.122 | -0.002 | -0.002 | 0.613 | 0.398 |
| nonDRNA | <i>Sensitivity</i> | 98.6% | 93.8% | 94.0% | 93.3% | 95.2% |
| | <i>Specificity</i> | 29.8% | 35.4% | 40.0% | 96.2% | 69.9% |
| | <i>Accuracy</i> | 91.3% | 87.6% | 88.2% | 93.6% | 92.5% |
| | <i>MCC</i> | 0.428 | 0.309 | 0.354 | 0.748 | 0.625 |
| <i>averageMCC</i> | | 0.296 | 0.195 | 0.229 | 0.714 | 0.556 |
| <i>accuracy4labels</i> | | 90.8% | 86.6% | 87.3% | 93.3% | 92.0% |

Conclusion – Part 1

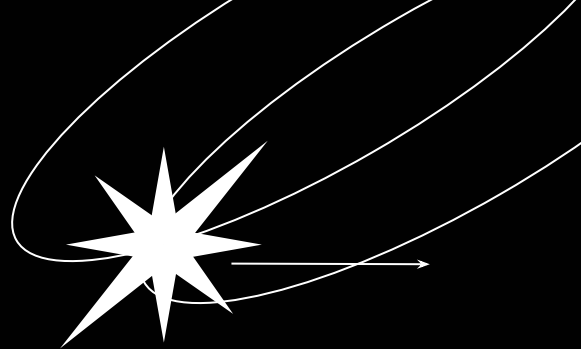
Advantages

- Uses SMOTE to oversample minority classes
- Uses a multitude of features generated from Pfeature
- Decision Tree parameters have been fine tuned

Disadvantages

- model uses the somewhat arbitrary feature selection from Pfeature
- a small number of classification models have been explored
- Possible overfitting

Experience – Concl. Pt 2



- Ultimately pleased with results from project
- Adjusted for overfitting, underrepresentation, and noise in our model
- Results steadily improved as a result of generating more features and using SMOTE
- Confident that we came close to desired results without facing oversampling
- Project was positive experience and taught us to process, generate, and analyze real-world data
- Helped prepare us for jobs in Data Science

Q/A

Any questions?

