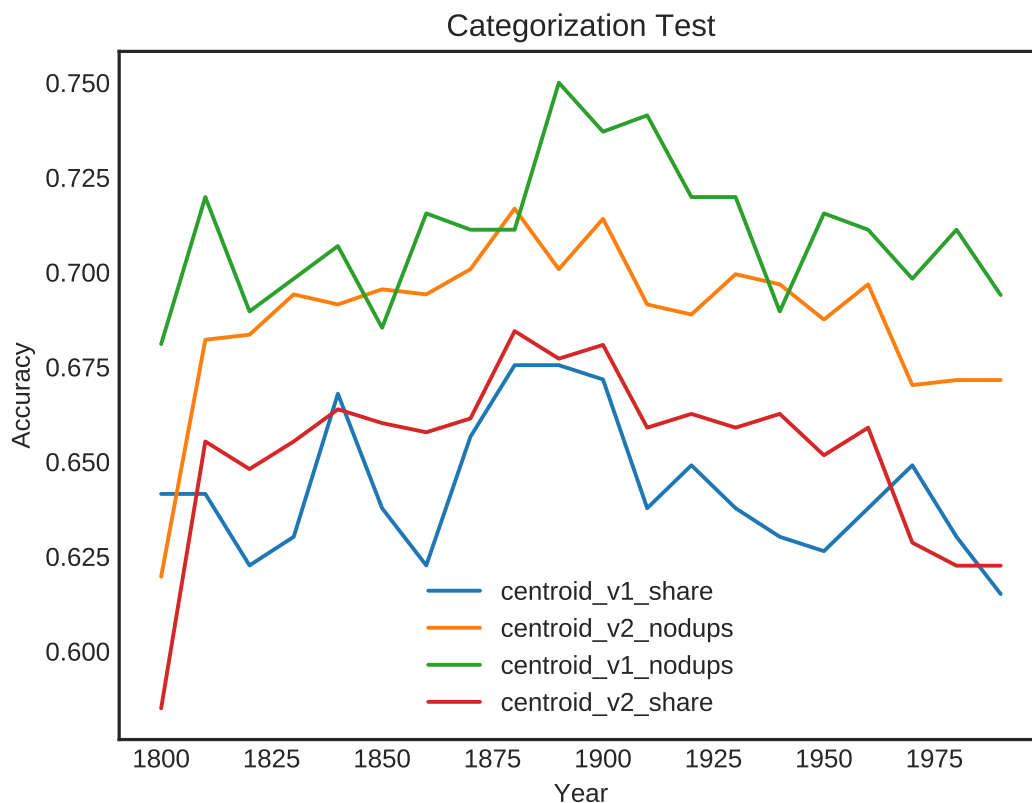Table 1: Average accuracy of all models and data schemes.
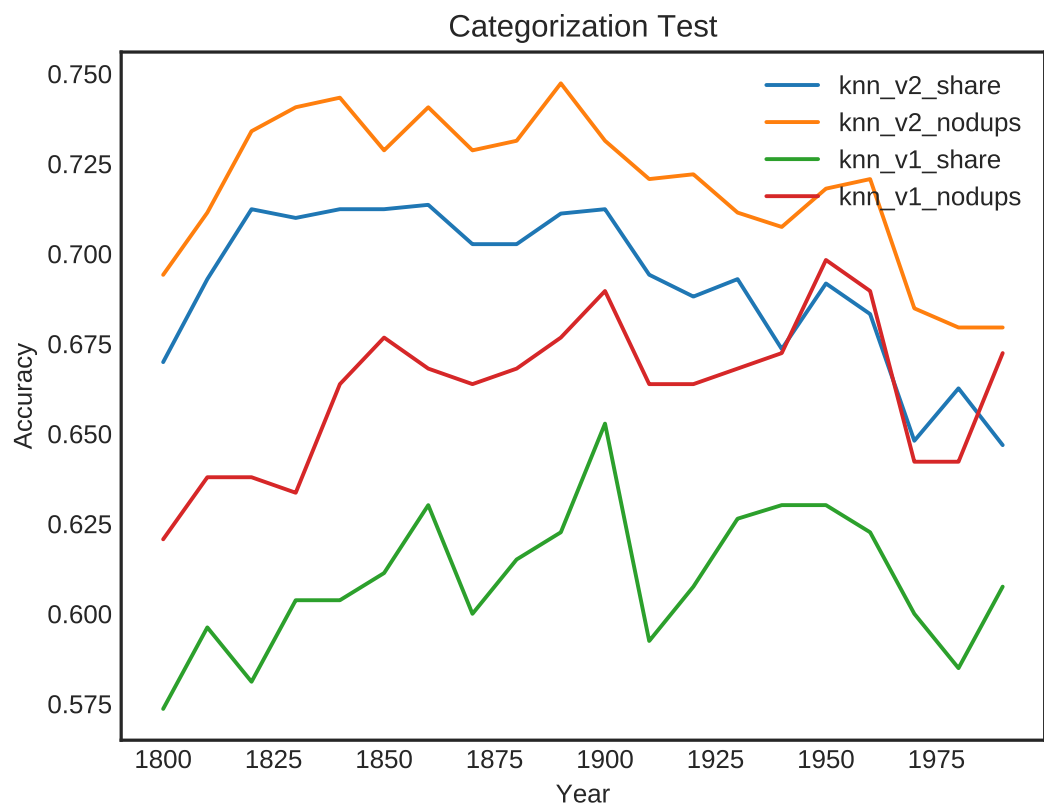
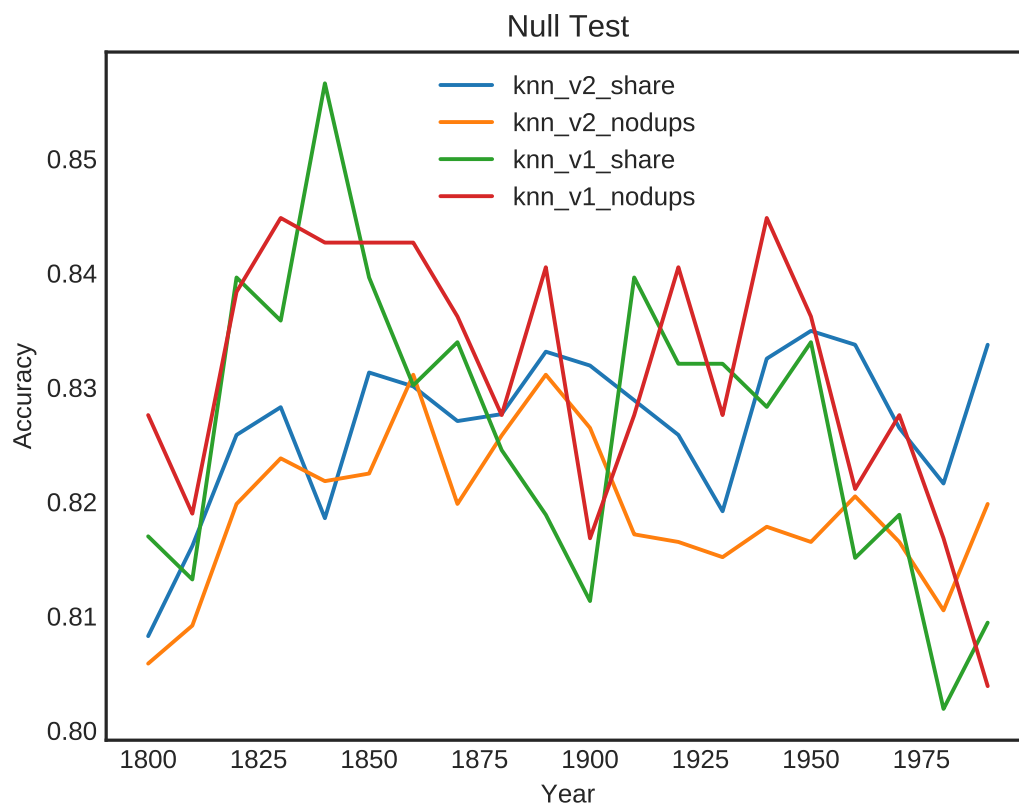| Model Name | Categorization Test | Null Test | Polarity Test |
|:---:|:---:|:---:|:---:|
| centroid˙v1˙share | 0.64 | 0.82 | 0.89 |
| centroid˙v1˙nodups | 0.71 | 0.82 | 0.89 |
| centroid˙v2˙share | 0.65 | 0.80 | 0.90 |
| centroid˙v2˙nodups | 0.69 | 0.79 | 0.90 |
| knn˙v1˙share | 0.61 | 0.83 | 0.91 |
| knn˙v1˙nodups | 0.66 | 0.83 | 0.91 |
| knn˙v2˙share | 0.69 | 0.83 | 0.91 |
| knn˙v2˙nodups | 0.72 | 0.82 | 0.91 |

The following plots compare the centroid and kNN models under variation of the MFD version (v1 vs v2) and whether duplicate seed words are kept ("share") or removed ("nodups").
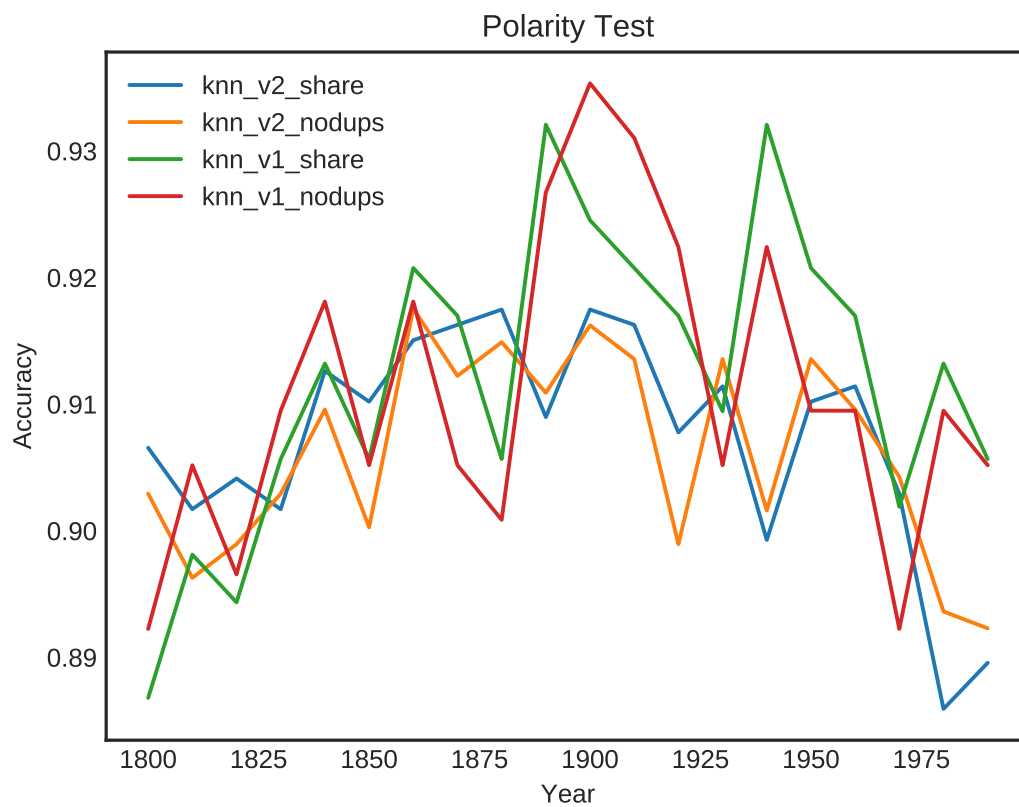
Null Test

Polarity Test

Categorization Test

Null Test

Polarity Test

confusion_matrix_categorization

| True label | +authority | +care | +fairness | +loyalty | +sanctity | -authority | -care | -fairness | -loyalty | -sanctity |
|---|---|---|---|---|---|---|---|---|---|---|
| +authority | 0.68 | 0.03 | 0.05 | 0.10 | 0.04 | 0.03 | 0.03 | 0.00 | 0.01 | 0.03 |
| +care | 0.04 | 0.65 | 0.03 | 0.09 | 0.08 | 0.01 | 0.01 | 0.01 | 0.04 | 0.05 |
| +fairness | 0.04 | 0.04 | 0.76 | 0.01 | 0.07 | 0.01 | 0.01 | 0.03 | 0.00 | 0.03 |
| +loyalty | 0.13 | 0.04 | 0.01 | 0.66 | 0.06 | 0.03 | 0.01 | 0.00 | 0.07 | 0.00 |
| +sanctity | 0.04 | 0.05 | 0.02 | 0.01 | 0.77 | 0.00 | 0.01 | 0.01 | 0.00 | 0.09 |
| -authority | 0.09 | 0.00 | 0.00 | 0.00 | 0.00 | 0.36 | 0.11 | 0.02 | 0.32 | 0.09 |
| -care | 0.00 | 0.01 | 0.01 | 0.02 | 0.00 | 0.09 | 0.68 | 0.02 | 0.04 | 0.14 |
| -fairness | 0.03 | 0.00 | 0.14 | 0.00 | 0.00 | 0.04 | 0.03 | 0.65 | 0.00 | 0.12 |
| -loyalty | 0.03 | 0.00 | 0.00 | 0.06 | 0.00 | 0.47 | 0.11 | 0.01 | 0.23 | 0.08 |
| -sanctity | 0.00 | 0.00 | 0.00 | 0.00 | 0.03 | 0.06 | 0.12 | 0.01 | 0.01 | 0.75 |

Predicted label

Figure 2: Categorization confusion matrix: Centroid, V1, 'nodups'



confusion_matrix_categorization

Figure 3: Categorization confusion matrix: Centroid, V2, 'share'



confusion_matrix_categorization

Figure 4: Categorization confusion matrix: Centroid, V2, 'nodups'



confusion_matrix_categorization

confusion_matrix_categorization

|  | +authority | +care | +fairness | +loyalty | +sanctity | -authority | -care | -fairness | -loyalty | -sanctity |
|---|---|---|---|---|---|---|---|---|---|---|
| +authority | 0.73 | 0.04 | 0.01 | 0.06 | 0.06 | 0.02 | 0.03 | 0.00 | 0.00 | 0.05 |
| +care | 0.12 | 0.63 | 0.02 | 0.03 | 0.10 | 0.00 | 0.05 | 0.01 | 0.00 | 0.05 |
| +fairness | 0.19 | 0.05 | 0.46 | 0.00 | 0.18 | 0.00 | 0.01 | 0.04 | 0.00 | 0.06 |
| +loyalty | 0.29 | 0.07 | 0.00 | 0.46 | 0.04 | 0.03 | 0.07 | 0.00 | 0.03 | 0.01 |
| +sanctity | 0.04 | 0.04 | 0.00 | 0.00 | 0.80 | 0.00 | 0.01 | 0.00 | 0.00 | 0.11 |
| -authority | 0.12 | 0.01 | 0.01 | 0.00 | 0.00 | 0.24 | 0.19 | 0.01 | 0.24 | 0.19 |
| -care | 0.02 | 0.00 | 0.00 | 0.00 | 0.00 | 0.07 | 0.74 | 0.00 | 0.02 | 0.15 |
| -fairness | 0.09 | 0.01 | 0.06 | 0.00 | 0.00 | 0.01 | 0.14 | 0.54 | 0.00 | 0.15 |
| -loyalty | 0.03 | 0.02 | 0.00 | 0.06 | 0.00 | 0.43 | 0.19 | 0.01 | 0.12 | 0.13 |
| -sanctity | 0.00 | 0.00 | 0.00 | 0.00 | 0.02 | 0.05 | 0.11 | 0.01 | 0.00 | 0.80 |

True label

Predicted label

Figure 6: Categorization confusion matrix: kNN, V1, 'nodups'



confusion_matrix_categorization

Figure 7: Categorization confusion matrix: kNN, V2, 'share'
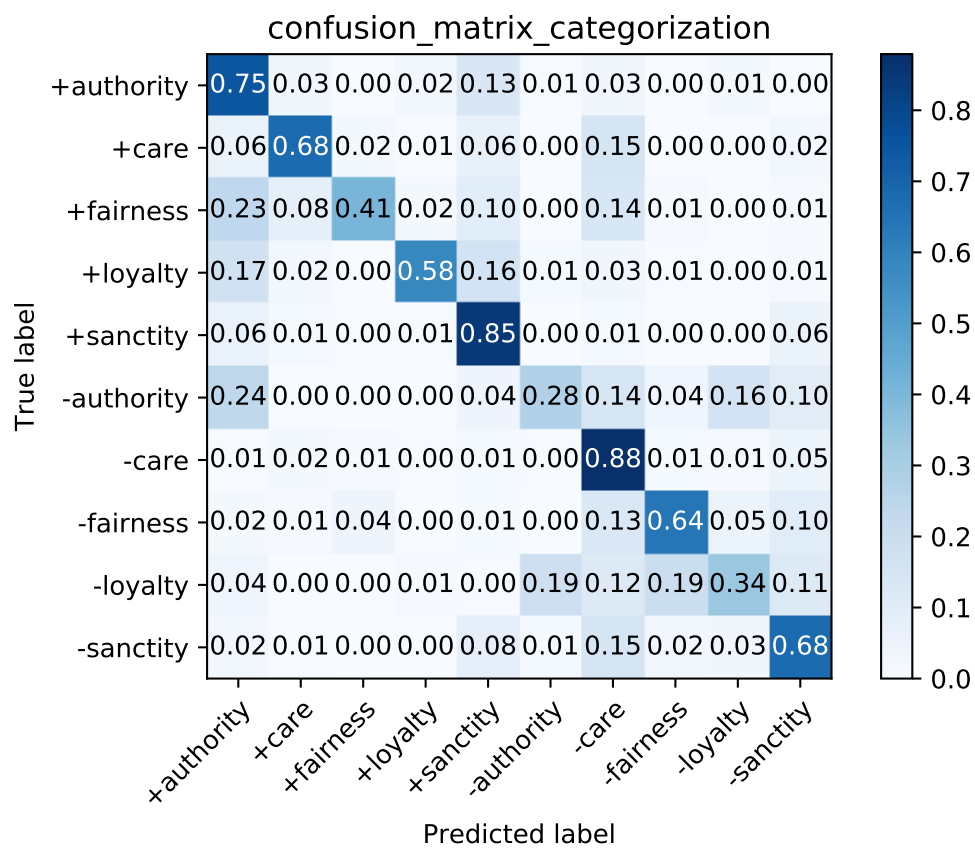
## confusion_matrix_categorization

Figure 8: Categorization confusion matrix: kNN, V2, 'nodups'



confusion_matrix_categorization