# Abstract

The Historical Population Register (HPR) is a project to build the longitudinal life history of individuals by integrating the historical records of the people in Norway since the 19th century. This study attempted to improve the linking rate between the 1875-1900 censuses in HPR, which is currently low, using machine learning approaches. To this end, I developed a machine learning model for linking that is suitable for the Norwegian census and tested different feature sets and match selection options. I compared the results in terms of performance and match size, and also examined their representativeness to the entire population. The study results showed that the linking rate of HPR can be significantly improved by machine learning approaches while maintaining high accuracy. In addition, this study presented a reference for future use by demonstrating how the performance varies depending on the feature set and match selection. On the other hand, this study also revealed that linked data generally do not represent the population of the census, and the characteristics and degree of bias vary depending on the linking method, suggesting that caution is needed when using linked data for research.

**Keywords**: historical record linkage, Norwegian census, Historical Population Register (HPR), machine learning

**Record linkage between the 1865-1875 censuses**

*\* Based on population in 1875 (over 10 years old)*

**Record linkage between the 1875-1900 censuses**

*\* Based on population in 1900 (over 25 years old)*

**Record linkage between the 1900-1910 censuses**

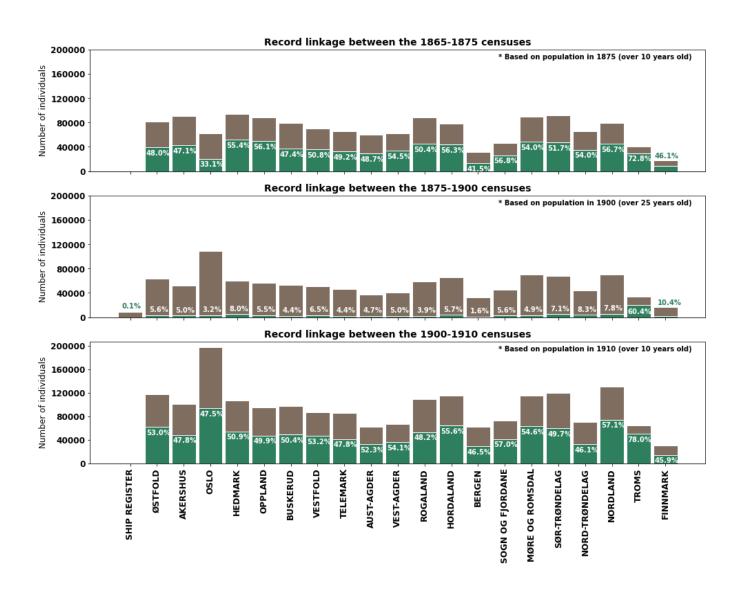*\* Based on population in 1910 (over 10 years old)*

Figure.1.1. Linking rates between historical censuses in the current HPR, by county.[1] The green area in each bar represents the proportion of people whose records are linked in that county.

---

[1] It includes the two largest cities of the time, Oslo (its name at the time was Kristiania) and Bergen.

Table 4.1. Performance measurement by algorithm

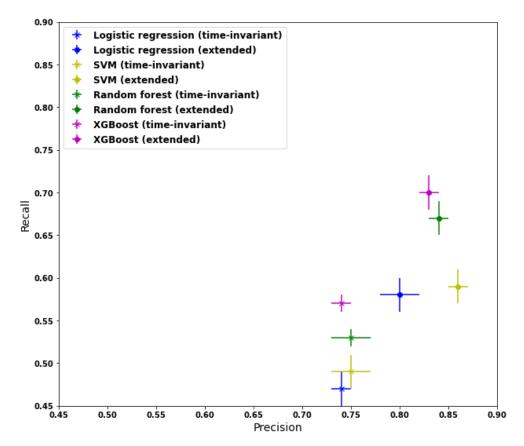| Algorithms | Features | Precision | | Recall | | F1-score | |
|---|---|---|---|---|---|---|---|
| | | mean | std | mean | std | mean | std |
| Logistic regression | time-invariant | 0.74 | 0.01 | 0.47 | 0.02 | 0.57 | 0.02 |
| | extended | 0.80 | 0.02 | 0.58 | 0.02 | 0.67 | 0.02 |
| SVM | time-invariant | 0.75 | 0.02 | 0.49 | 0.02 | 0.59 | 0.02 |
| | extended | 0.86 | 0.01 | 0.59 | 0.02 | 0.70 | 0.01 |
| Random forest | time-invariant | 0.75 | 0.02 | 0.53 | 0.01 | 0.62 | 0.01 |
| | extended | 0.84 | 0.01 | 0.67 | 0.02 | 0.74 | 0.01 |
| XGBoost | time-invariant | 0.74 | 0.01 | 0.57 | 0.01 | 0.64 | 0.01 |
| | extended | 0.83 | 0.01 | 0.70 | 0.02 | 0.76 | 0.01 |



Figure 4.1. Performance measurement by algorithm. Visualization of Table 4.1. Each point and its x-axis and y-axis lines represent the mean and standard deviation in 5-folds cross-validation.

Table 4.5. Model performance according to match selection parameters for the test set (Time-invariant feature based model). The highest F1-scores across the table and in unique matches are shown in color.

| Absolute cutoff | | Relative cutoff | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1.1 | 1.2 | 1.3 | 1.4 | 1.5 | 1.6 | 1.7 | 1.8 | 1.9 | 2 | Unique |
| 0.1 | Precision | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.51 | 0.52 |
| | Recall | 0.82 | 0.82 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.80 | 0.80 | 0.80 | 0.80 | 0.77 |
| | F1-score | **0.63** | **0.63** | **0.63** | **0.63** | **0.63** | **0.63** | **0.63** | **0.63** | **0.63** | **0.63** | **0.63** | **0.62** |
| 0.2 | Precision | 0.60 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 |
| | Recall | 0.75 | 0.75 | 0.75 | 0.75 | 0.74 | 0.74 | 0.74 | 0.74 | 0.74 | 0.74 | 0.74 | 0.73 |
| | F1-score | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.66** |
| 0.3 | Precision | 0.66 | 0.66 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 |
| | Recall | 0.69 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 | 0.67 |
| | F1-score | **0.68** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** | **0.67** |
| 0.4 | Precision | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.73 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 |
| | Recall | 0.62 | 0.62 | 0.62 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 |
| | F1-score | **0.67** | **0.66** | **0.66** | **0.66** | **0.66** | **0.66** | **0.66** | **0.66** | **0.66** | **0.66** | **0.66** | **0.66** |
| 0.5 | Precision | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 |
| | Recall | 0.56 | 0.56 | 0.56 | 0.56 | 0.56 | 0.56 | 0.56 | 0.55 | 0.55 | 0.55 | 0.55 | 0.55 |
| | F1-score | **0.65** | **0.65** | **0.65** | **0.65** | **0.65** | **0.65** | **0.65** | **0.64** | **0.64** | **0.64** | **0.64** | **0.64** |
| 0.6 | Precision | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 |
| | Recall | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| | F1-score | **0.62** | **0.61** | **0.61** | **0.61** | **0.61** | **0.61** | **0.61** | **0.61** | **0.61** | **0.61** | **0.61** | **0.61** |
| 0.7 | Precision | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 |
| | Recall | 0.42 | 0.42 | 0.42 | 0.41 | 0.41 | 0.41 | 0.41 | 0.41 | 0.41 | 0.41 | 0.41 | 0.41 |
| | F1-score | **0.56** | **0.55** | **0.55** | **0.55** | **0.55** | **0.55** | **0.55** | **0.55** | **0.55** | **0.55** | **0.55** | **0.55** |
| 0.8 | Precision | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 |
| | Recall | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 | 0.29 |
| | F1-score | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** | **0.44** |
| 0.9 | Precision | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 |
| | Recall | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 | 0.13 |
| | F1-score | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** | **0.23** |

Table 4.6. Model performance according to match selection parameters for the test set (Extended feature based model). The highest F1-scores across the table and in unique matches are shown in color.

| Absolute cutoff | | Relative cutoff | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1.1 | 1.2 | 1.3 | 1.4 | 1.5 | 1.6 | 1.7 | 1.8 | 1.9 | 2 | Unique |
| 0.1 | Precision | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 |
| | Recall | 0.87 | 0.87 | 0.87 | 0.87 | 0.87 | 0.87 | 0.86 | 0.86 | 0.86 | 0.85 | 0.85 | 0.82 |
| | F1-score | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | 0.71 |
| 0.2 | Precision | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 |
| | Recall | 0.83 | 0.83 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.81 | 0.81 | 0.81 | 0.80 |
| | F1-score | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** | **0.76** | **0.76** | 0.76 |
| 0.3 | Precision | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 |
| | Recall | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.78 | 0.77 | 0.77 | 0.77 | 0.77 | 0.76 |
| | F1-score | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.77** |
| 0.4 | Precision | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 |
| | Recall | 0.74 | 0.74 | 0.74 | 0.73 | 0.73 | 0.73 | 0.73 | 0.73 | 0.73 | 0.73 | 0.73 | 0.72 |
| | F1-score | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.78** | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** |
| 0.5 | Precision | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 | 0.85 |
| | Recall | 0.70 | 0.70 | 0.70 | 0.70 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 |
| | F1-score | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** | **0.77** | **0.76** | **0.76** | **0.76** | **0.76** | **0.76** | 0.76 |
| 0.6 | Precision | 0.88 | 0.89 | 0.89 | 0.89 | 0.89 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 |
| | Recall | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 | 0.65 |
| | F1-score | **0.75** | **0.75** | **0.75** | **0.75** | **0.75** | **0.75** | **0.75** | **0.75** | **0.75** | **0.75** | **0.75** | 0.75 |
| 0.7 | Precision | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 | 0.91 |
| | Recall | 0.60 | 0.60 | 0.60 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 |
| | F1-score | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | **0.72** | 0.72 |
| 0.8 | Precision | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 | 0.93 |
| | Recall | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 |
| | F1-score | **0.68** | **0.68** | **0.68** | **0.68** | **0.68** | **0.68** | **0.68** | **0.68** | **0.68** | **0.68** | **0.68** | 0.68 |
| 0.9 | Precision | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 |
| | Recall | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 | 0.43 |
| | F1-score | **0.59** | **0.59** | **0.59** | **0.59** | **0.59** | **0.59** | **0.59** | **0.59** | **0.59** | **0.59** | **0.59** | 0.59 |

Table 4.9. Model performance according to match selection parameters for the full test set provided by the NHDC (Time-invariant feature based model). The highest F1-scores across the table and in unique matches are shown in color.

| Absolute cutoff | | Relative cutoff | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1.1 | 1.2 | 1.3 | 1.4 | 1.5 | 1.6 | 1.7 | 1.8 | 1.9 | 2 | Unique |
| | Precision | 0.93 | 0.94 | 0.95 | 0.95 | 0.95 | 0.95 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.97 |
| 0.1 | Recall | 0.78 | 0.76 | 0.75 | 0.75 | 0.74 | 0.74 | 0.73 | 0.73 | 0.72 | 0.72 | 0.72 | 0.61 |
| | F1-score | 0.84 | 0.84 | 0.84 | 0.84 | 0.83 | 0.83 | 0.83 | 0.83 | 0.83 | 0.82 | 0.82 | 0.74 |
| | Precision | 0.94 | 0.96 | 0.96 | 0.96 | 0.96 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| 0.2 | Recall | 0.74 | 0.72 | 0.71 | 0.71 | 0.70 | 0.70 | 0.70 | 0.69 | 0.69 | 0.69 | 0.68 | 0.64 |
| | F1-score | 0.82 | 0.82 | 0.82 | 0.82 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.80 | 0.80 | 0.77 |
| | Precision | 0.95 | 0.96 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.3 | Recall | 0.70 | 0.68 | 0.68 | 0.67 | 0.67 | 0.66 | 0.66 | 0.66 | 0.65 | 0.65 | 0.65 | 0.63 |
| | F1-score | 0.80 | 0.80 | 0.80 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.78 | 0.78 | 0.78 | 0.77 |
| | Precision | 0.95 | 0.97 | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.4 | Recall | 0.65 | 0.64 | 0.64 | 0.63 | 0.63 | 0.62 | 0.62 | 0.62 | 0.62 | 0.61 | 0.61 | 0.61 |
| | F1-score | 0.78 | 0.77 | 0.77 | 0.77 | 0.76 | 0.76 | 0.76 | 0.76 | 0.76 | 0.75 | 0.75 | 0.75 |
| | Precision | 0.96 | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.5 | Recall | 0.60 | 0.59 | 0.58 | 0.58 | 0.58 | 0.57 | 0.57 | 0.57 | 0.57 | 0.57 | 0.57 | 0.57 |
| | F1-score | 0.74 | 0.73 | 0.73 | 0.73 | 0.73 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 | 0.72 |
| | Precision | 0.96 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.6 | Recall | 0.54 | 0.53 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 | 0.52 |
| | F1-score | 0.69 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 | 0.68 |
| | Precision | 0.96 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.7 | Recall | 0.45 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 |
| | F1-score | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 | 0.61 |
| | Precision | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.8 | Recall | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 | 0.32 |
| | F1-score | 0.49 | 0.48 | 0.48 | 0.48 | 0.48 | 0.48 | 0.48 | 0.48 | 0.48 | 0.48 | 0.48 | 0.48 |
| | Precision | 0.97 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.9 | Recall | 0.15 | 0.14 | 0.14 | 0.14 | 0.14 | 0.14 | 0.14 | 0.14 | 0.14 | 0.14 | 0.14 | 0.14 |
| | F1-score | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 | 0.25 |

Table 4.10. Model performance according to match selection parameters for the full test set provided by the NHDC (Extended feature based model). The highest F1-scores across the table and in unique matches are shown in color.

| Absolute cutoff | | Relative cutoff | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1.1 | 1.2 | 1.3 | 1.4 | 1.5 | 1.6 | 1.7 | 1.8 | 1.9 | 2 | Unique |
| | Precision | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.99 |
| 0.1 | Recall | 0.90 | 0.89 | 0.88 | 0.87 | 0.86 | 0.85 | 0.84 | 0.84 | 0.83 | 0.82 | 0.82 | 0.61 |
| | F1-score | 0.94 | 0.93 | 0.93 | 0.92 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.90 | 0.89 | 0.76 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.2 | Recall | 0.90 | 0.88 | 0.87 | 0.86 | 0.85 | 0.84 | 0.84 | 0.83 | 0.82 | 0.81 | 0.81 | 0.69 |
| | F1-score | 0.94 | 0.93 | 0.93 | 0.92 | 0.91 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.82 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.3 | Recall | 0.89 | 0.88 | 0.87 | 0.85 | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.81 | 0.80 | 0.73 |
| | F1-score | 0.94 | 0.93 | 0.92 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.89 | 0.84 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.4 | Recall | 0.88 | 0.87 | 0.86 | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.80 | 0.80 | 0.79 | 0.76 |
| | F1-score | 0.93 | 0.93 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.88 | 0.88 | 0.86 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| 0.5 | Recall | 0.87 | 0.86 | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.80 | 0.79 | 0.79 | 0.78 | 0.78 |
| | F1-score | 0.93 | 0.92 | 0.92 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.88 | 0.88 | 0.88 | 0.88 |
| | Precision | 0.99 | 0.99 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| 0.6 | Recall | 0.87 | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 |
| | F1-score | 0.93 | 0.92 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 |
| | Precision | 0.99 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| 0.7 | Recall | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 |
| | F1-score | 0.92 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 |
| | Precision | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| 0.8 | Recall | 0.84 | 0.82 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 |
| | F1-score | 0.91 | 0.90 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 |
| | Precision | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| 0.9 | Recall | 0.80 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 |
| | F1-score | 0.89 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 |

Table 4.11. Model performance according to match selection parameters for the sub test set provided by the NHDC (Time-invariant feature based model). The highest F1-scores across the table and in unique matches are shown in color.

| Absolute cutoff | | Relative cutoff | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1.1 | 1.2 | 1.3 | 1.4 | 1.5 | 1.6 | 1.7 | 1.8 | 1.9 | 2 | Unique |
| | Precision | 0.92 | 0.94 | 0.95 | 0.95 | 0.95 | 0.95 | 0.96 | 0.96 | 0.96 | 0.96 | 0.96 | 0.97 |
| 0.1 | Recall | 0.78 | 0.77 | 0.76 | 0.75 | 0.75 | 0.75 | 0.74 | 0.74 | 0.74 | 0.73 | 0.73 | 0.62 |
| | F1-score | 0.85 | 0.85 | 0.84 | 0.84 | 0.84 | 0.84 | 0.84 | 0.83 | 0.83 | 0.83 | 0.83 | 0.76 |
| | Precision | 0.94 | 0.95 | 0.96 | 0.96 | 0.96 | 0.96 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| 0.2 | Recall | 0.75 | 0.73 | 0.73 | 0.72 | 0.72 | 0.71 | 0.71 | 0.71 | 0.70 | 0.70 | 0.70 | 0.66 |
| | F1-score | 0.83 | 0.83 | 0.83 | 0.82 | 0.82 | 0.82 | 0.82 | 0.82 | 0.81 | 0.81 | 0.81 | 0.78 |
| | Precision | 0.94 | 0.96 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 | 0.97 |
| 0.3 | Recall | 0.71 | 0.69 | 0.69 | 0.68 | 0.68 | 0.68 | 0.67 | 0.67 | 0.67 | 0.66 | 0.66 | 0.65 |
| | F1-score | 0.81 | 0.81 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.79 | 0.79 | 0.79 | 0.79 | 0.78 |
| | Precision | 0.95 | 0.97 | 0.97 | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.4 | Recall | 0.66 | 0.65 | 0.65 | 0.64 | 0.64 | 0.64 | 0.63 | 0.63 | 0.63 | 0.63 | 0.63 | 0.62 |
| | F1-score | 0.78 | 0.78 | 0.78 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.77 | 0.76 | 0.76 | 0.76 |
| | Precision | 0.95 | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.5 | Recall | 0.61 | 0.60 | 0.60 | 0.59 | 0.59 | 0.59 | 0.59 | 0.59 | 0.58 | 0.58 | 0.58 | 0.58 |
| | F1-score | 0.75 | 0.74 | 0.74 | 0.74 | 0.74 | 0.73 | 0.73 | 0.73 | 0.73 | 0.73 | 0.73 | 0.73 |
| | Precision | 0.96 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.6 | Recall | 0.55 | 0.54 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 | 0.53 |
| | F1-score | 0.70 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 | 0.69 |
| | Precision | 0.97 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 |
| 0.7 | Recall | 0.47 | 0.46 | 0.45 | 0.45 | 0.45 | 0.45 | 0.45 | 0.45 | 0.45 | 0.45 | 0.45 | 0.45 |
| | F1-score | 0.63 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 | 0.62 |
| | Precision | 0.97 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.8 | Recall | 0.34 | 0.33 | 0.33 | 0.33 | 0.33 | 0.33 | 0.33 | 0.33 | 0.33 | 0.33 | 0.33 | 0.33 |
| | F1-score | 0.51 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 | 0.50 |
| | Precision | 0.98 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.9 | Recall | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 | 0.16 |
| | F1-score | 0.28 | 0.27 | 0.27 | 0.27 | 0.27 | 0.27 | 0.27 | 0.27 | 0.27 | 0.27 | 0.27 | 0.27 |

Table 4.12. Model performance according to match selection parameters for the sub test set provided by the NHDC (Extended feature based model). The highest F1-scores across the table and in unique matches are shown in color.

| Absolute cutoff | | Relative cutoff | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 1.1 | 1.2 | 1.3 | 1.4 | 1.5 | 1.6 | 1.7 | 1.8 | 1.9 | 2 | Unique |
| | Precision | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.98 | 0.99 |
| 0.1 | Recall | 0.90 | 0.89 | 0.88 | 0.87 | 0.86 | 0.85 | 0.84 | 0.83 | 0.83 | 0.82 | 0.81 | 0.62 |
| | F1-score | 0.94 | 0.93 | 0.93 | 0.92 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.76 |
| | Precision | 0.98 | 0.98 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.2 | Recall | 0.90 | 0.88 | 0.87 | 0.86 | 0.85 | 0.84 | 0.83 | 0.83 | 0.82 | 0.81 | 0.80 | 0.69 |
| | F1-score | 0.94 | 0.93 | 0.92 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.89 | 0.82 |
| | Precision | 0.98 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.3 | Recall | 0.89 | 0.87 | 0.86 | 0.85 | 0.84 | 0.83 | 0.83 | 0.82 | 0.81 | 0.80 | 0.80 | 0.73 |
| | F1-score | 0.93 | 0.93 | 0.92 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.88 | 0.84 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.4 | Recall | 0.88 | 0.87 | 0.85 | 0.84 | 0.84 | 0.83 | 0.82 | 0.81 | 0.80 | 0.80 | 0.79 | 0.76 |
| | F1-score | 0.93 | 0.92 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.88 | 0.88 | 0.86 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.5 | Recall | 0.87 | 0.86 | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.80 | 0.80 | 0.79 | 0.78 | 0.78 |
| | F1-score | 0.93 | 0.92 | 0.91 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.88 | 0.88 | 0.88 | 0.88 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.6 | Recall | 0.86 | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 | 0.80 |
| | F1-score | 0.92 | 0.92 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 |
| | Precision | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 | 0.99 |
| 0.7 | Recall | 0.85 | 0.84 | 0.83 | 0.82 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 |
| | F1-score | 0.92 | 0.91 | 0.90 | 0.90 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 |
| | Precision | 0.99 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| 0.8 | Recall | 0.84 | 0.82 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 | 0.81 |
| | F1-score | 0.91 | 0.90 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 | 0.89 |
| | Precision | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| 0.9 | Recall | 0.80 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 | 0.79 |
| | F1-score | 0.89 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 | 0.88 |