# Automatic Methods for Coding Historical Occupation Descriptions to Standard Classifications

Graham Kirby, Jamie Carson, Fraser Dunlop, Chris Dibben, Alan Dearle, Lee Williamson, Eilidh Garrett, Alice Reid

digitisingscotland@lscs.ac.uk

digitisingscotland.cs.st-andrews.ac.uk

University of St Andrews

THE UNIVERSITY of EDINBURGH

UNIVERSITY OF CAMBRIDGE

digitising SCOTLAND
understanding Scotland's people

E·S·R·C ECONOMIC & SOCIAL RESEARCH COUNCIL

wellcome trust

# Motivation

- Increasing number of digitised registration records for the 19th and 20th centuries.

- Varying forms of data

- Scale of data prevents manual analysis

# Challenges

- Significant methodological issues:
  - How can we consistently code occupational data so that researchers can explore changing patterns and trends?
  - How can we automate this process so that the majority of records do not need to be manually coded?

# Digitising Scotland

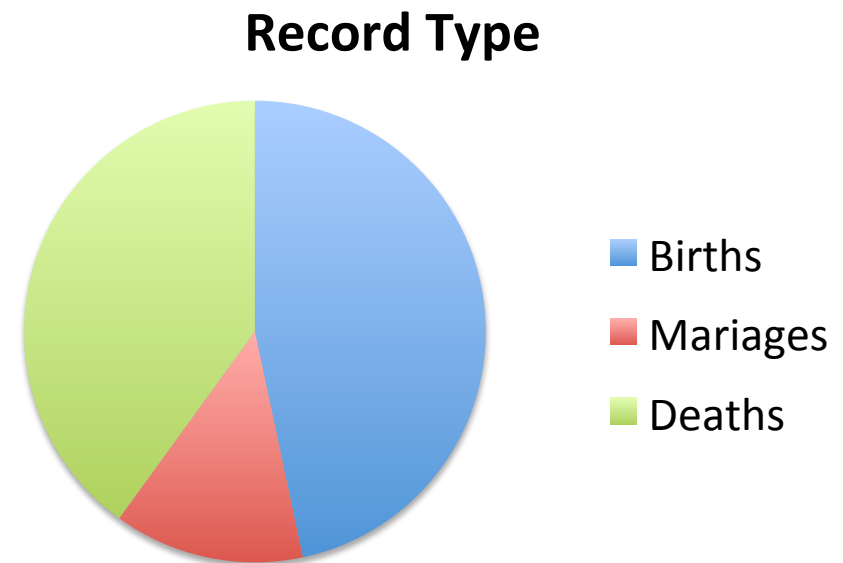- Records of births, marriages and deaths recorded in Scotland from 1855 to present day.

# Digitising Scotland

- Approximately 29 million records
- Approximately 50 million occupation strings, 8 million causes of death
- Classify occupations to Historical International Standard Classification of Occupations (HISCO)
- Cause of death to ICD10

**Record Type**



- Births
- Mariages
- Deaths

**MARRIAGE** | District No. _107_ | Year _1972_ | Entry No. _1_

IN THE DISTRICT OF _Nesting_

**1. When and where married**

19.72  December ~~eighth~~ Eighth

The Manse, Nesting

|  | BRIDEGROOM | BRIDE |
|---|---|---|
| Surname | Williamson | Pottinger |
| Name(s) | George Angus | Agnes Anne |
| (Signed) | George A Williamson | Agnes Anne Pottinger |
| 3. Occupation | Civil Servant | Farm Worker |

| 4. Marital status | | 5. Date of birth Year | Month | Day | 4. | | 5. Year | Month | Day |
|---|---|---|---|---|---|---|---|---|---|
| Bachelor | ~~1951~~ | 1951 | 6 | 17 | Spinster | | 1952 | 7 | 26 |

| 6. Birthplace | ~~St. Peter Bridge of Walls~~ Lerwick | Benston, |
|---|---|---|
| 7. Usual residence | Belmoe, Bridge of Walls | Benston, Nesting |
| 8. Father's name(s) surname and occupation | — | John Charles Pottinger Farmer |

6

# MARRIAGE

| | | | |
|---|---|---|---|
| **MARRIAGE** | District No. 107 | Year 1972 | Entry No. 1 |

IN THE DISTRICT OF _Nesting_

1. When and where married

19.72 December ~~eight~~ Eighth
The Manse, Nesting

| | BRIDEGROOM | BRIDE |
|---|---|---|
| Surname | Williamson | Pottinger |
| Name(s) | George Angus | Agnes Anne |
| (Signed) | George A Williamson | Agnes Anne Pottinger |
| 3. Occupation | Civil Servant | Farm Worker |

| | | 5. Date of birth | Year | Month | Day | | | 5. Year | Month | Day |
|---|---|---|---|---|---|---|---|---|---|---|
| 4. Marital status | Bachelor | | 1951 | 6 | 17 | Spinster | | 1952 | 7 | 26 |

| | | |
|---|---|---|
| 6. Birthplace | ~~St. Andr~~ ~~Bridge of Walls~~ Lerwick | Benston, |
| 7. Usual residence | Selivoe, Bridge of Walls | Benston, Nesting |
| 8. Father's name(s) surname and occupation | — | John Charles Pottinger Farmer |

7

# MARRIAGE

| District No. 107 | Year 1972 | Entry No. 1 |

IN THE DISTRICT OF **Nesting**

**1. When and where married**

19.72 December ~~Eighth~~ Eighth

The Manse, Nesting

| | BRIDEGROOM | BRIDE |
|---|---|---|
| Surname | Williamson | Pottinger |
| Name(s) | George Angus | Agnes Anne |
| (Signed) | George A. Williamson | Agnes Anne Pottinger |
| 3. Occupation | Civil Servant | Farm Worker |

| | | 5. Date of birth Year | Month | Day | | | 5. Year | Month | Day |
|---|---|---|---|---|---|---|---|---|---|
| 4. Marital status | Bachelor | ~~1951~~ 1951 | 6 | 17 | | Spinster | 1952 | 7 | 26 |

| 6. Birthplace | ~~St Ringo~~ Bridge of Walls Lerwick | Benston, |
| 7. Usual residence | Belivoe, Bridge of Walls | Benston, Nesting |
| 8. Father's name(s) surname and occupation | — — | John Charles Pottinger Farmer |

8

## MARRIAGE

| | | | |
|---|---|---|---|
| **MARRIAGE** | District No. 107 | Year 1972 | Entry No. 1 |

IN THE DISTRICT OF Nesting

1. When and where married

19.72 December ~~Eighth~~ Eighth

The Manse, Nesting

| | BRIDEGROOM | BRIDE |
|---|---|---|
| Surname | Williamson | Pottinger |
| Name(s) | George Angus | Agnes Anne |
| (Signed) | George A Williamson | Agnes Anne Pottinger |
| 3. Occupation | Civil Servant | Farm Worker |

| 4. Marital status | 5. Date of birth | Year | Month | Day | 4. | 5. Year | Month | Day |
|---|---|---|---|---|---|---|---|---|
| Bachelor | ~~1951~~ | 1951 | 6 | 17 | Spinster | 1952 | 7 | 26 |

| | | |
|---|---|---|
| 6. Birthplace | ~~St Boots Bridge of Walls~~ Lerwick | Benston |
| 7. Usual residence | Belivoe | Benston |
| | Bridge of Walls | Nesting |
| 8. Father's name(s) surname and occupation | — | John Charles Pottinger |
| | — | Farmer |

9

# MARRIAGE

| District No. | Year | Entry No. |
|---|---|---|
| 107 | 1972 | 1 |

**IN THE DISTRICT OF** Nesting

**1. When and where married**

19.72 December eighth Eighth

The Manse, Nesting

|  | BRIDEGROOM | BRIDE |
|---|---|---|
| Surname | Williamson | Pottinger |
| Name(s) | George Angus | Agnes Anne |
| (Signed) | George A. Williamson | Agnes Anne Pottinger |
| 3. Occupation | Civil Servant | Farm Worker |

| 4. Marital status | 5. Date of birth | Year | Month | Day | 4. | 5. Year | Month | Day |
|---|---|---|---|---|---|---|---|---|
| Bachelor | | 1951 | 6 | 17 | Spinster | 1952 | 7 | 26 |

| 6. Birthplace | Lerwick | Benston, |
|---|---|---|
| 7. Usual residence | Belivoe | Benston, |
| | Bridge of Walls | Nesting |
| 8. Father's name(s) surname and occupation | — | John Charles Pottinger |
| | | Farmer |

10

# Experimental Dataset

- Vital event records currently being transcribed

- Use a dataset with similar content for experiments

- 60,000 records from the Cambridge Family History Study (records from 1800-1990)

- Occupation descriptions and associated HISCO codes

- HISCO coding done by historians

- Dataset contains 330 different HISCO codes

# HISCO Hierarchy Example



**Major Groups**

6 Agricultural, animal husbandry and forestry workers, fishermen and hunters

**Minor Groups**

62 Agricultural And Animal Husbandry Workers

63 Forestry Workers

**Unit Groups**

624 Livestock Workers

631 Loggers

**Micro Groups**

62460 Horse Worker

# Classification Example

| String from record | Gold Standard Classification | Automatic Classification Output |
|---|---|---|
| Farm horseman | 62460 | 62460 |
| Shoe maker | 80110 | 80110 |
| Fireman (railway) | 98330 | 98330 |
| Fireman | 58100 | 58100 |
| Stationer | 41000 | 91000 |

# Classification Example

| String from record | Gold Standard Classification | Automatic Classification Output |
|---|---|---|
| Farm horseman | 62460 Horse Worker | 62460 Horse Worker |
| Shoe maker | 80110 Shoemaker, General | 80110 Shoemaker, General |
| Fireman (railway) | 98330 Railway Steam-Engine Fireman | 98330 Railway Steam-Engine Fireman |
| Fireman | 58100 Fire-Fighter | 58100 Fire-Fighter |
| Stationer | 41000 Working Proprietors (Wholesale and Retail Trade) | 91000 Paper and Paperboard product makers |

# Classification Example

| String from record | Gold Standard Classification | Automatic Classification Output |
|---|---|---|
| Farm horseman | 62460 Horse Worker | 62460 Horse Worker |
| Shoe maker | 80110 Shoemaker, General | 80110 Shoemaker, General |
| Fireman (railway) | 98330 Railway Steam-Engine Fireman | 98330 Railway Steam-Engine Fireman |
| Fireman | 58100 Fire-Fighter | 58100 Fire-Fighter |
| Stationer | 41000 Working Proprietors (Wholesale and Retail Trade) | 91000 Paper and Paperboard product makers |

# Approach

- Text analysis
- Supervised machine learning
  - Apache Mahout framework.
- Combination of these techniques.

# Supervised Machine Learning

**Training Data** → **Machine Learning** → **Prediction Model**

**Unseen Data** → **Prediction Model** → **Predicted Classification**

# Supervised Machine Learning

**Training Data** ➜ **Machine Learning** ➜ **Prediction Model**

| | |
|---|---|
| Farm horseman | 62460 |
| Shoe maker | 80110 |
| Fireman | 58100 |
| Stationer | 41000 |

**Unseen Data** **Prediction Model** ➜ **Predicted Classification**

# Supervised Machine Learning

**Training Data** ➡ **Machine Learning** ➡ **Prediction Model**

| | |
|---|---|
| Farm horseman | 62460 |
| Shoe maker | 80110 |
| Fireman | 58100 |
| Stationer | 41000 |

**Unseen Data** ➡ **Prediction Model** ➡ **Predicted Classification**

Farm horseman
Boot maker
Fireman
Painter

# Supervised Machine Learning

**Training Data** ➡️ **Machine Learning** ➡️ **Prediction Model**

| | |
|---|---|
| Farm horseman | 62460 |
| Shoe maker | 80110 |
| Fireman | 58100 |
| Stationer | 41000 |

**Unseen Data** ➡️ **Prediction Model** ➡️ **Predicted Classification**

| |
|---|
| Farm horseman |
| Boot maker |
| Fireman |
| Painter |

?

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot and shoe maker | | |
| B | Boot and shoe dealer | | |
| C | Fireman | | |
| D | Cattle (& sheep) farmer | | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|---|---|---|---|
| A | Boot **and** shoe maker | boot shoe maker | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot **and** shoe maker | boot shoe maker | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | | | | | | | | | |
| B | | | | | | | | | |
| C | | | | | | | | | |
| D | | | | | | | | | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot **and** shoe maker | **boot** shoe maker | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | **1** | | | | | | | | |
| B | | | | | | | | | |
| C | | | | | | | | | |
| D | | | | | | | | | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot **and** shoe maker | boot shoe maker | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | 1 | 0 | | | | | | | |
| B | | | | | | | | | |
| C | | | | | | | | | |
| D | | | | | | | | | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot **and** shoe maker | boot shoe maker | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | 1 | 0 | 0 | 0 | 0 | 0 | | | |
| B | | | | | | | | | |
| C | | | | | | | | | |
| D | | | | | | | | | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot **and** shoe maker | boot shoe **maker** | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | 1 | 0 | 0 | 0 | 0 | 0 | **1** | | |
| B | | | | | | | | | |
| C | | | | | | | | | |
| D | | | | | | | | | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot **and** shoe maker | boot **shoe** maker | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | **1** |
| B | | | | | | | | | |
| C | | | | | | | | | |
| D | | | | | | | | | |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|----------------|---------------|--------|
| A | Boot **and** shoe maker | boot shoe maker | |
| B | Boot **and** shoe dealer | boot shoe dealer | |
| C | Fireman | fireman | |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| B | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| C | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| D | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |

# Machine Learning

- Inputs are split into features and converted to high dimension vectors

| Record | Original Input | Cleaned input | Vector |
|--------|---------------|---------------|--------|
| A | Boot **and** shoe maker | boot shoe maker | 100000101 |
| B | Boot **and** shoe dealer | boot shoe dealer | 101000001 |
| C | Fireman | fireman | 000010000 |
| D | Cattle **(&** sheep**)** farmer | cattle sheep farmer | 010100010 |

| | boot | cattle | dealer | farmer | fireman | horse | maker | sheep | shoe |
|---|------|--------|--------|--------|---------|-------|-------|-------|------|
| A | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| B | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| C | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| D | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |

# Approach to Classification

Training Data

# Approach to Classification
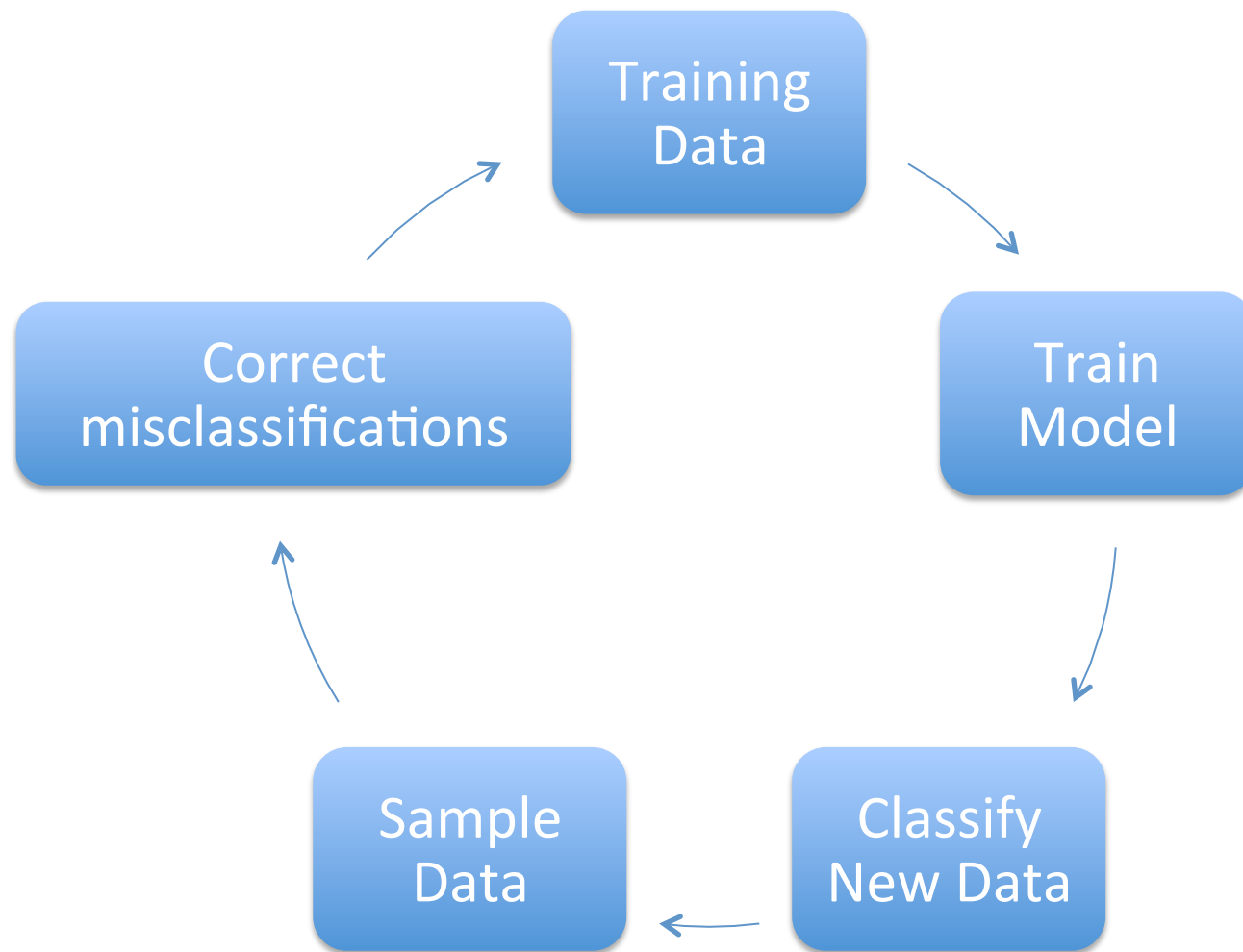


Training Data → Train Model

# Approach to Classification

# Approach to Classification

# Approach to Classification

# Feature Selection

- Changes to the data input to the classification system
- Feature Selection
  - Selecting the most appropriate features to use in the training data
- Analyse the input to identify features are likely to harm the quality of the classification
- Common words that appear in lots of different output codes.
- Example: "farmers daughter", "Butchers daughter"…
  - Remove the word daughter.

# Gold Standard Misclassification

- Variations in coding of unique strings make it harder to calculate a good model

- Different coders, extra data, mistakes

- Try removing strings coded to multiple codes

- Try changing less common codes to most common

# Edit Distance Classifier

- Relatively simple string similarity classifier
- HISCO uses numerical codes, so compare with code description
- Assume similar inputs have similar descriptions
- Similarity measured using edit distance
  - Number of single-character insertions, deletions or replacements needed to transform
- Look for highest number of exact matches between words, fall back to similarity if equal number of matches.

# Edit Distance Example

| Occupation | Gold Standard Output | Edit Distance Output |
|---|---|---|
| Hotel proprietor | Working Proprietor (Hotel and Restaurant) | Working Proprietor (Hotel and Restaurant) |
| Taxi driver | Taxi driver | Taxi driver |
| Tax clerk | Tax collector | Tax collector |
| Painter & decorator | Painters, Construction | Sign Painter |
| File Cutter | Machinery Fitters, Machine Assemblers and Precision Instrument Makers (except Electrical) NEC | Stock Clerks |

# Edit Distance Example

| Occupation | Gold Standard Output | Edit Distance Output |
|---|---|---|
| Hotel proprietor | Working Proprietor (Hotel and Restaurant) | Working Proprietor (Hotel and Restaurant) |
| Taxi driver | Taxi driver | Taxi driver |
| Tax clerk | Tax collector | Tax collector |
| Painter & decorator | Painters, Construction | Sign Painter |
| File Cutter | Machinery Fitters, Machine Assemblers and Precision Instrument Makers (except Electrical) NEC | Stock Clerks |

# Edit Distance Example

| Occupation | Gold Standard Output | Edit Distance Output |
| --- | --- | --- |
| Hotel proprietor | Working Proprietor (Hotel and Restaurant) | Working Proprietor (Hotel and Restaurant) |
| Taxi driver | Taxi driver | Taxi driver |
| Tax clerk | Tax collector | Tax collector |
| Painter & decorator | Painters, Construction | Sign Painter |
| File Cutter | Machinery Fitters, Machine Assemblers and Precision Instrument Makers (except Electrical) NEC | Stock Clerks |

# Individual Machine Learning Classifiers

- ## Naïve Bayes
    - Probabilistic classifier
    - Co-occurrence of features

- ## Stochastic Gradient Descent
    - Optimisation of logistic regression

# Ensemble Approaches

- Majority voting
  - Pick the most frequent classification

- Confidence threshold technique
  - Pick the SGD classification unless its likelihood value is below a given threshold

- Pseudo confidence threshold
  - Produce a pseudo measure of likelihood for the Naive Naïve classifier. Pick the best classification from Naïve Bayes and SGD.

# Experiments

- Which single classification technique produces the highest accuracy when classifying occupations to HISCO?

- Which ensemble technique produces the highest accuracy?

- What difference, if any, does using feature selection make?

- What difference does fixing or removing multiple codings make?

- What effect does classifying to different HISCO levels make?
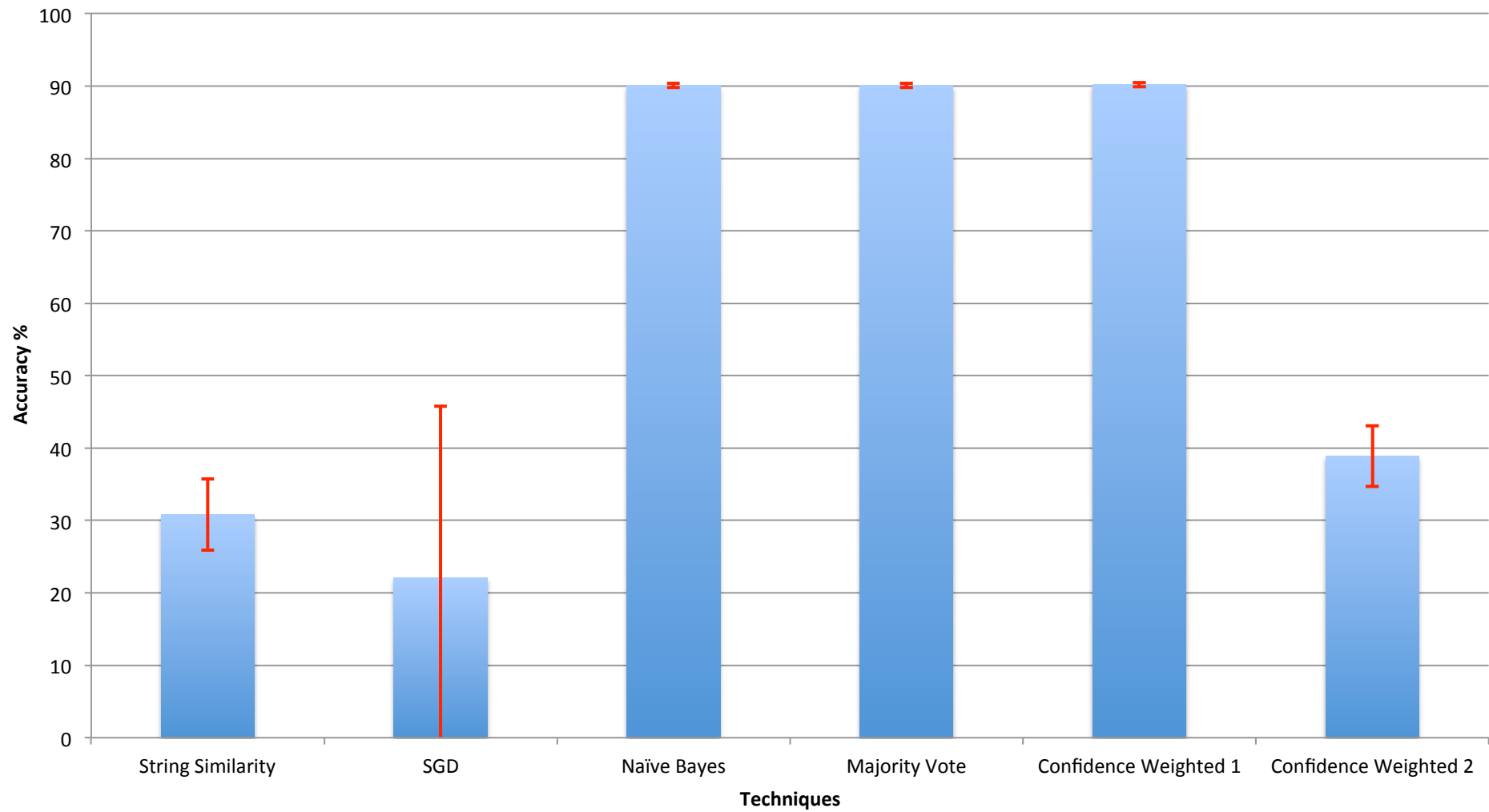
# Evaluation

- Need to assess the quality of the automatic coding
- Hold out method
- Split data into two sets, a training set and a validation set
- 80% chosen for training, 20% for validation
- Pick a new training/test set each repetition
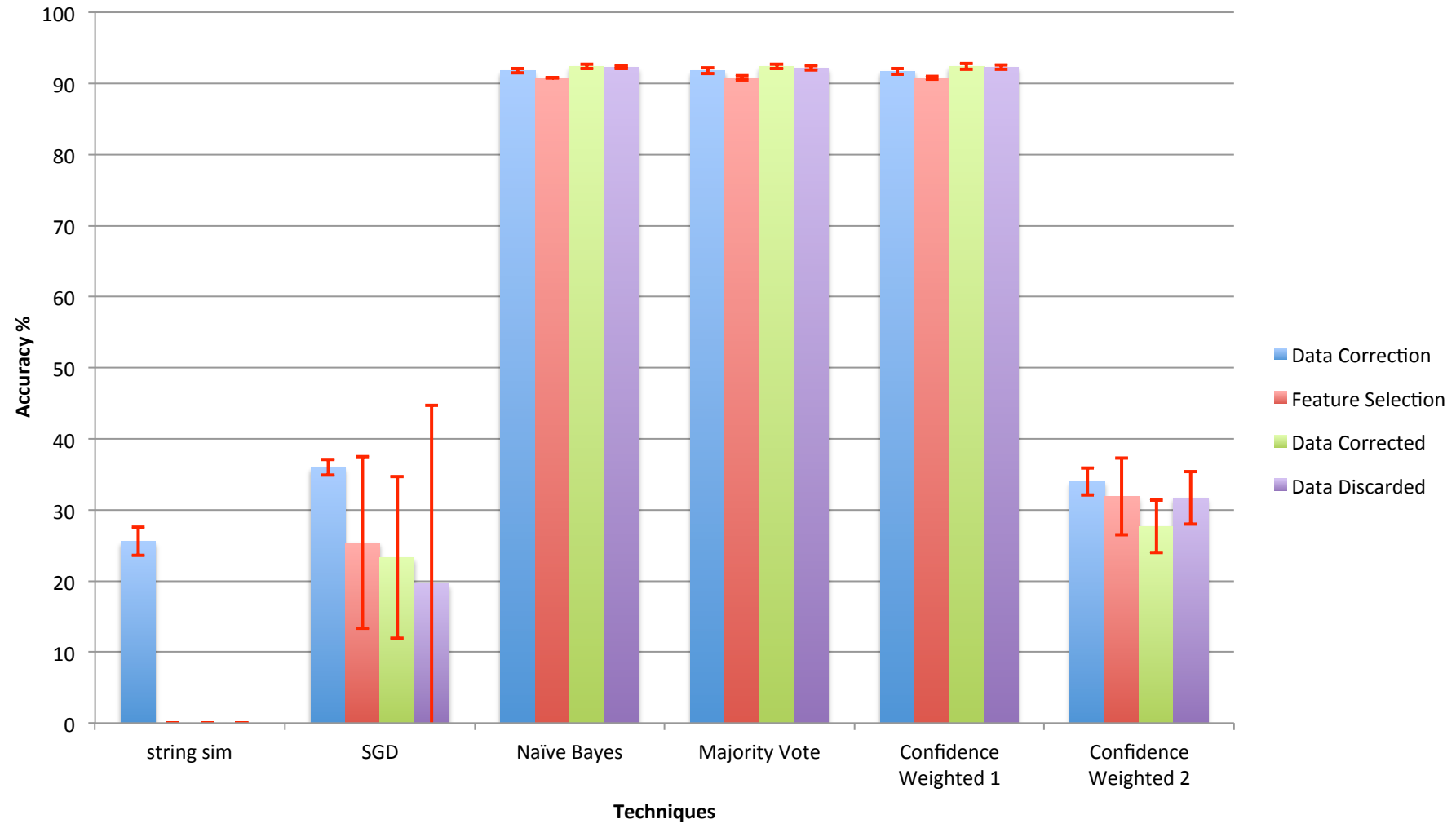- Correct classification is gold standard code matches output code

# Accuracy Measures

- HISCO employs a hierarchical structure
- If we are only interested in coarse classifications we can relax the closeness of the match required
- Match unit group
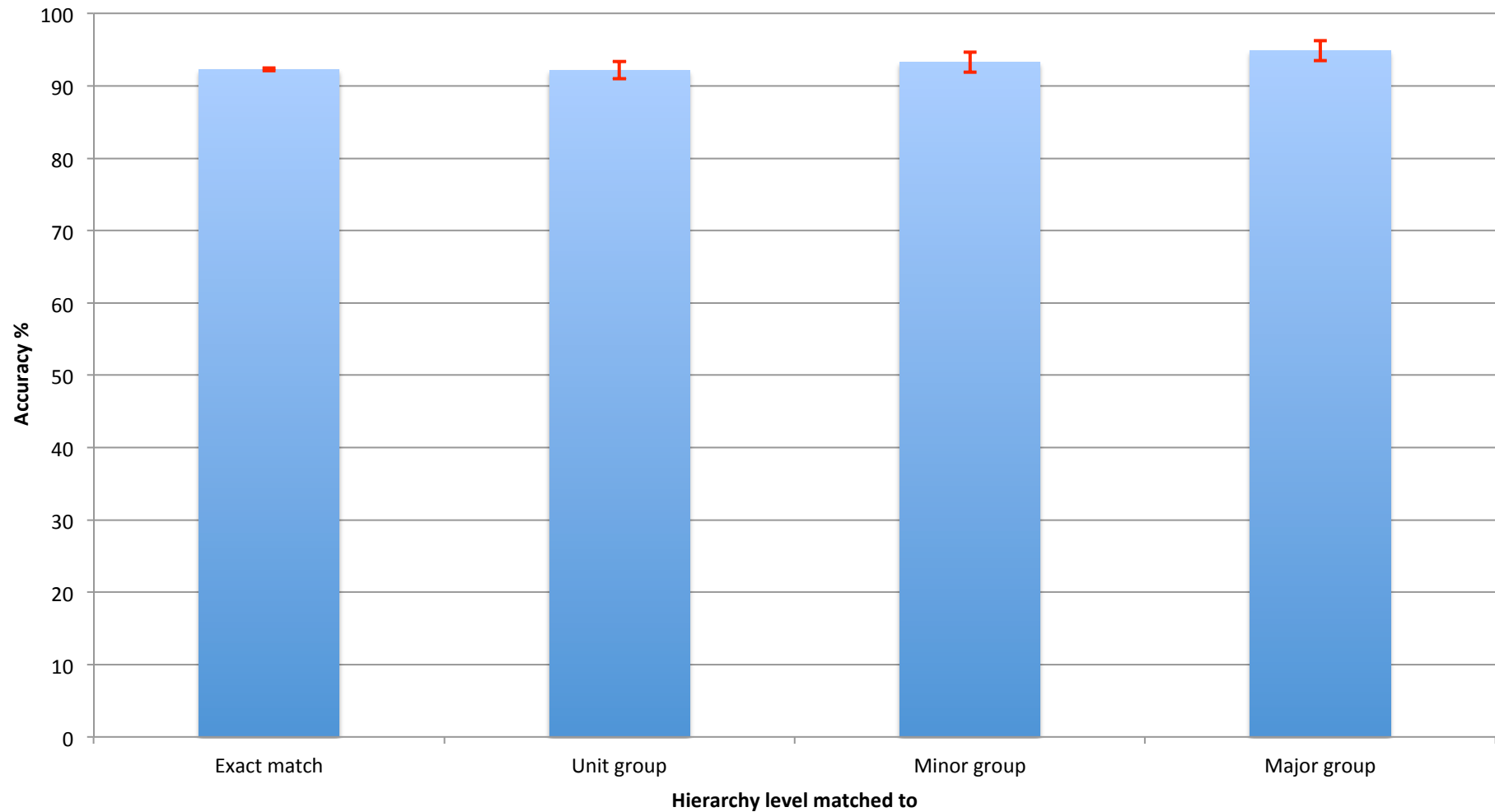- Match minor group
- Match major group

Classification Accuracy

**Comparison of different data manipulations techniques**

Varying levels of HISCO hierarchy with Naïve Bayes Classifier

# Summary

- Highest accuracy: Naïve Bayes classifier with feature selection and correction of multiply coded descriptions.

- Exact match accuracy using this technique was 92.3 ± 0.2%

- Considering only major group matching 94.9 ± 1.4% was achieved.

- Although the ensemble did not improve performance, addition of another high performance algorithm should yield gains.

# Discussion

- Previous results classifying cause of death using ensemble methods showed improvement of 2-3%

- Run times:
  - String Similarity: a few minutes
  - Naïve Bayes: a few minutes
  - SGD: 3-4 hours depending on learning parameters

- SGD has been reworked, preliminary results: 88-94%

# Future Work

- Continue machine learning and string matching development to classify cause of death and occupations

- Continue to examine behaviour of SGD algorithm to try and achieve better performance.

- Add further machine learning models, such as support vector machines into the ensemble