**NIH** **Intramural** Research Program
*Our Research Changes Lives*

# SOCcer: Automatic coding of free-text job descriptions to standardized occupation codes

Daniel E. Russ[1], Kwan-Yuet Ho[1], Calvin A. Johnson[1], and Melissa C. Friesen[2]

[1]Division of Computational Bioscience, Center for Information Technology, NIH
[2]Occupational and Environmental Epidemiology Branch, Division of Cancer Epidemiology and Genetics, NCI, NIH

# Overview

- Occupation information
  - Use and collection
  - Standardized classification systems
- SOCcer
  - Framework
  - Performance measures
  - Future advancements and new directions

# Occupation Information



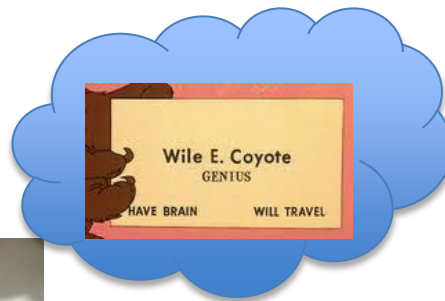Finance, credit, purchase preferences



Medical diagnosis



Surveillance to improve health policy and identify priorities



Primary or confounding factor in **epidemiologic studies**

# What is your job title?

- Many ways of describing an occupation
- Different level of detail
- Changes depending on who is asking

Wile E. Coyote
GENIUS
HAVE BRAIN          WILL TRAVEL

Standardized Occupation Categories
Mechanical Engineer
Misc. Engineer
Physical Scientist
Chief Executive

# Standard Occupational Codes

**US SOC 2010**

| | | |
|---|---|---|
| Major group | **19**-0000 | Life, Physical, and Social Science Occupations |
| Minor Group | **19-1**000 | Life Scientists |
| Broad Group | **19-104**0 | Medical Scientists |
| Detailed Occupation | **19-1041** | Epidemiologists |
| | **19-1042** | Medical Scientists, Except Epidemiologists |

| | | |
|---|---|---|
| Major group | 29-0000 | Healthcare Practitioners and Technical Occupations |
| Minor Group | 29-1000 | Health Diagnosing and Treating Practitioners |
| Broad Group | 29-1060 | Physicians and Surgeons |
| Detailed Occupation | 29-1062 | Family and General Practitioners |
| | 29-1063 | General Internist |
| | 29-1069 | Physicians and Surgeons, All Other |

# Occupation in population-based studies

- Current job, longest job, usual job, all jobs
- Wide variety of occupations, industries
- Open-ended questions:
  - What was your job title?
  - What were your main tasks and activities in this job?
  - Who was your employer?
  - What services were provided or what products were made by your employer?
  - Start year/stop year
- Coded to standardized occupation and industry classification systems: SOC and SIC

# Coding: time-consuming, modestly reliable

- Manual process

- Based on limited information

- No gold standard

- Agreement between 2 coders is poor/moderate
(Koeman et al. 2013)
    - 5-digit level ISCO68 agreement:  36%
    - 3-digit level ISCO68 agreement:  55%

- Preferably independent assignments by 2 coders, resolve discordant assignments

- Costly in large-scale studies

# Multiple recent efforts to automate

- NIOSH: http://wwwn.cdc.gov/niosh-nioccs/
- U. Montreal: www.caps-canada.ca
- Burstyn et al. (2014) Beyond crosswalks: reliability of exposure assessment following automated coding of free-text job descriptions for occupational epidemiology.
- Patel et al. (2012) Performance of automated and manual coding systems for occupational data: A case study of historical records

- Batch vs. job-by-job.
- Most require the user to make the final determination from multiple choices or do not provide an assignment when low confidence/no match.
- Different coding schemes.

# Objective

- Develop a computer algorithm to assign standardized occupation classifications (SOC) based on free-text responses.
  - Reduce but not replace expert coding

- Cross-NIH institute collaboration with Division of Computational Bioscience
  - Expertise in natural language processing and classification

# Our framework

- Adaptable system
  - [Initially] Code to US SOC 2010
    - Electronic knowledge base of job titles and tasks for each SOC (O*NET)
  - [Future] Other classification systems

- Assumption: Better matches by using multiple aspects of job description
  - Job title, tasks, coded industry

# http://soccer.nci.nih.gov

**Intramural Research Program**
*Our Research Changes Lives*

# Knowledge Base Development

| Data Sources | | |
|---|---|---|
| Production database of job titles<br>Crosswalk information | Direct Match Title File | U.S. Census Occupational Index |

~ 62,000 Job Titles

Job Title SOC
Job Title SOC

# O*NET for US SOC 2010

**O*NET OnLine**

*Updated 2010*

**Summary Report for:**
45-4021.00 - Fallers

Use axes or chainsaws...
direction of fall and m...

**Sample of reported job...**
Sawyer, Tree Topper

**Sample of reported job titles:** Timber Faller, Tree Faller, Timber Cutter, Logger, Tree Feller, Cutter Operator, Sawyer, Tree Topper

View report: Summary | Details | Custom

Tasks | Tools & Technolog...
Work Styles | Work Value...

**Tasks**

- Stop saw eng...
- Appraise trees...
  direction of le...
- Saw back-cut...
- Clear brush fr...
  area, chainsa...
- Measure felled...
- Assess logs a...
- Determine po...
- Control the di...
  with chainsaw...
- Trim off the to...
- Select trees t...
  work.

back to top

**Tools & Techno...**

Tools used in this oc...

Air or gas tanks or...
Forestry skidders
Lifting cables — C...
Lumbering equipm...
Ultrasonic examini...

Technology used in this occupation:

## Tasks

- Stop saw engines, pull cutting bars from cuts, and run to safety as tree falls.
- Appraise trees for certain characteristics, such as twist, rot, and heavy limb growth, and gauge amount and direction of lean, to determine how to control the direction of a tree's fall with the least damage.
- Saw back-cuts, leaving sufficient sound wood to control direction of fall.
- Clear brush from work areas and escape routes, and cut saplings and other trees from direction of falls, using axes, chainsaws, or bulldozers.
- Measure felled trees and cut them into specified log lengths, using chain saws and axes.
- Assess logs after cutting to ensure that the quality and length are correct.
- Determine position, direction, and depth of cuts to be made, and placement of wedges or jacks.
- Control the direction of a tree's fall by scoring cutting lines with axes, sawing undercuts along scored lines with chainsaws, knocking slabs from cuts with single-bit axes, and driving wedges.
- Trim off the tops and limbs of trees, using chainsaws, delimbers, or axes.
- Select trees to be cut down, assessing factors such as site, terrain, and weather conditions before beginning work.

## Work Activities

**Performing General Phy...**
arms and legs and movin...
materials.

**Controlling Machines a...**
machines or processes (...

**Handling and Moving O...**
materials, and manipulati...

back to top

## Work Context

**Outdoors, Exposed to W...**

**Wear Common Protectiv...**
Protection, Hard Hats, ...

## Knowledge

**No knowledge met the minimu...**

back to top

## Skills

**Operation and Control** — Con...

back to top

## Abilities

**Reaction Time** — The ability to...
when it appears.

**Multilimb Coordination** — The...
leg and one arm) while sitting, s...
body is in motion.

http://www.onetonline.org

# Overview of SOCcer



**Job Title**

Soft Jaccard Score → Highest

Maximum Entropy

$x_0$    $x_1$    $x_2$    $x_3$

**Tasks Performed**

Fuzzy Fingerprint

$x_4$

**Industry**

SOC Prevalence >0.01

$x_5$

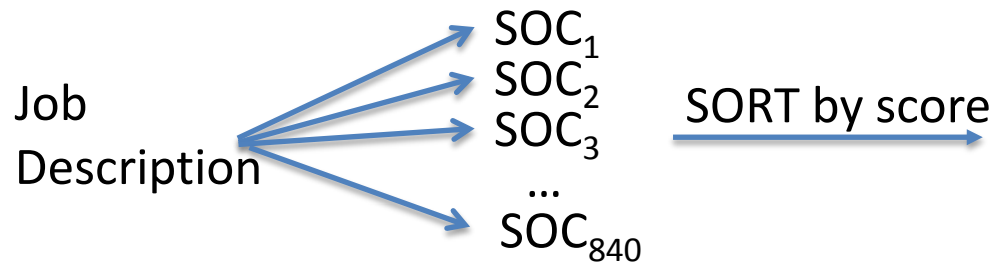Logistic Regression to obtain score

$$\ln\left(\frac{p}{1-p}\right) = f(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5$$

# SOCcer score and SOC code assignment

**Algorithm score:**
probability that an expert
coder would assign that code

Job
Description

SOC$_1$
SOC$_2$
SOC$_3$
...
SOC$_{840}$

SORT by score

SOCcer output:
Top 10 scoring SOCs

**SOC$_1$ – Assign to job description?**
SOC$_2$
SOC$_3$
...
SOC$_{10}$

# SOCcer's performance

Validity assessments at 6-digit level:

"Gold standard": Consensus **expert** SOC assignment

Vs. **Highest scoring SOC-2010 code** from SOCcer

# Overall agreement

| Study | # Jobs | Percent Agreement at SOC Level | | | | Median SOCcer score (IQR) |
|---|---|---|---|---|---|---|
| | | 2-Digit | 3-Digit | 5-Digit | 6-Digit | |
| US Renal | 11,943 | 77 | 64 | 52 | 45 | 0.46 (0.24-0.77) |
| CT Thyroid | 1,942 | 73 | 59 | 50 | 44 | 0.22 (0.08-0.51) |
| Montreal Lung | 829 | 74 | 56 | 46 | 38 | 0.45 (0.22-0.79) |
| **Combined** | **14,519** | **76** | **63** | **51** | **41** | **0.44 (0.23-0.75)** |

| | Coder vs. Coder | Computer vs. Coders | | |
|---|---|---|---|---|
| | Koeman | Burstyn | Patel (NIOSH) | US Renal |
| - 3-digit | 55% | 31-85% | 63% | 63% |
| - 5-digit | 36% | 9-72% | 52% | 51% |

# Did the expert assigned code appear in any of the top 10 codes from SOCcer?

# Agreement by score

Agreement by score distance to 2$^{nd}$ ranked SOC code:

Score distance = score$_1$ − score$_2$

# Do the mis-assignments matter?

- Generic codes may be difficult to correctly code, but may result in same exposure estimate (e.g., welder)

- Linked expert & highest scoring SOC code to CANJEM

- Compared agreement in exposure estimates
  - Generally similar patterns to SOC code
  - Median kappa on exposed/not exposed:        0.56 (IQR 0.52-0.58)
  - Median ICC on continuous probability metric: 0.66 (IQR 0.58-0.73)
  - Median ICC on continuous probability metric: 0.50 (IQR 0.44-0.56)

# Did the top 2 ranked SOC codes assign the same exposure?

# Agreement by confidence in assignment
HIGH: score $\geq 0.3$ _and_ score$_{1-2}$ $\geq 0.15$



**Confidence by metric**

# Agreement at varying hierarchy levels

# Main findings

- Can reduce expert coding task

- Detailed coding not always possible (data quality)

- 6-digit level
  - Manual coding necessary for nontrivial # of jobs
  - Score and JEM-based metrics to prioritize expert assessment

- 3-digit level
  - Excellent overall agreement and at scores ≥0.25

# SOCAssign: Companion Software



NIH⟩ Intramural Research Program
Our Research Changes Lives

# Preview of future advancements

- **SOCcer 1.1**
  - Expand training data to include job descriptions from epidemiologic studies and retrain algorithm
  - Increases **overall** 6-digit agreement to **>50%**
  - March-April 2017

- **SOCcer 2.0**
  - Add and refine classifiers and add training data
  - Increases **overall** 6-digit agreement to **~60%**
  - Late 2017

- **Consider more than one plausible code?**

**NIH** Intramural Research Program
*Our Research Changes Lives*

# Expansion to other systems?

- **Requires training data in target system**
  - Previously coded job descriptions
  - O*NET equivalent data source
  - Crosswalk to US SOC 2010

- **Canadian system?**
  - Current and previous versions
  - English & French?

# Thank you

- *Expert coders*
  - Susan Viet, Pabitra Josse, Sarah Locke
- *CANJEM*
  - Jerome Lavoue, Thomas Remen
- *Epidemiologic studies*
  - Connecticut Thyroid Study
  - Montreal Lung Cancer Study
  - NCI-SEER NHL Study
  - New England Bladder Cancer Study
  - US Renal Cell Cancer Study