# Autocoding the Survey of Occupational Injuries and Illnesses

## Alexander Measure

# Survey of Occupational Injuries and Illnesses

## <u>Example Narrative</u>

**Job title**: sanitation worker

**What was the employee doing just before the incident?**
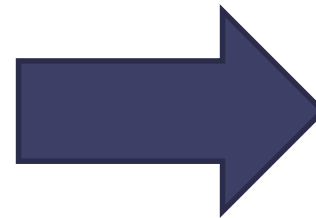mopping floor in gym

**What happened?**
slipped on water on floor and fell

**What part of the body was affected?**
fractured right arm

**What object directly harmed the employee?**
wet floor

## <u>Codes Assigned</u>

**Occup**: 37-2011 (Janitor)

**Nature**: 111 (Fracture)

**Part**: 420 (Arm)

**Event**: 422 (Fall, slipping)

**Source**: 6620 (Floor)
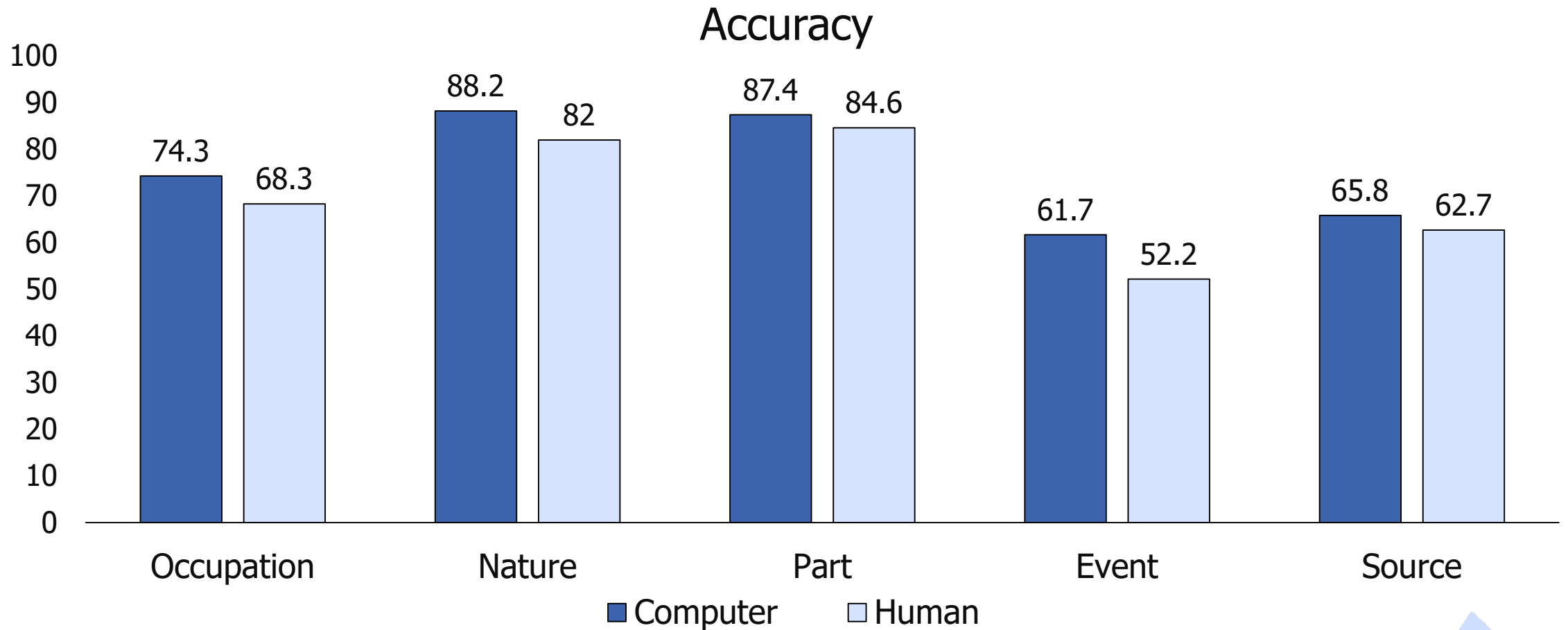
**Secondary**: 9521(Water)

# Supervised Machine Learning

1. Data (training, validation, and test)

2. Determine inputs and numeric representation

3. Choose a learning algorithm

4. Fit to training data, evaluate on validation

5. Modify and repeat
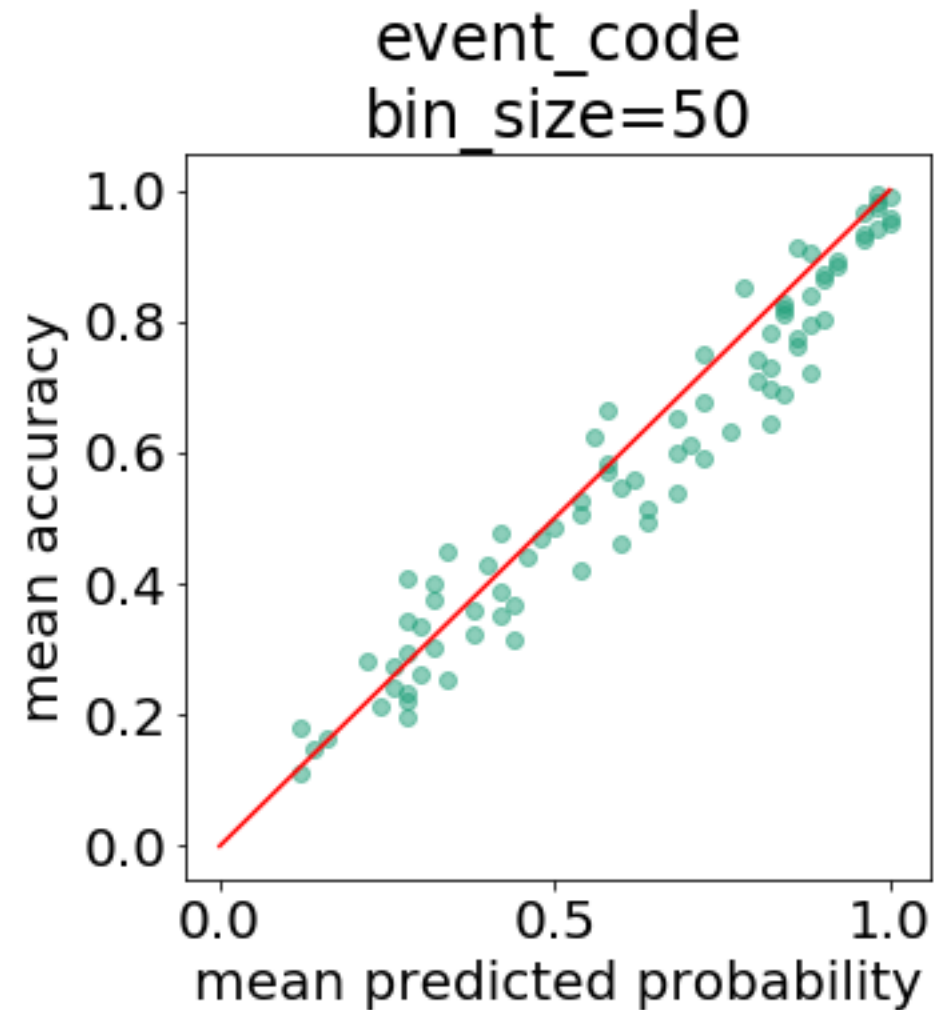
6. At the very end, evaluate on test
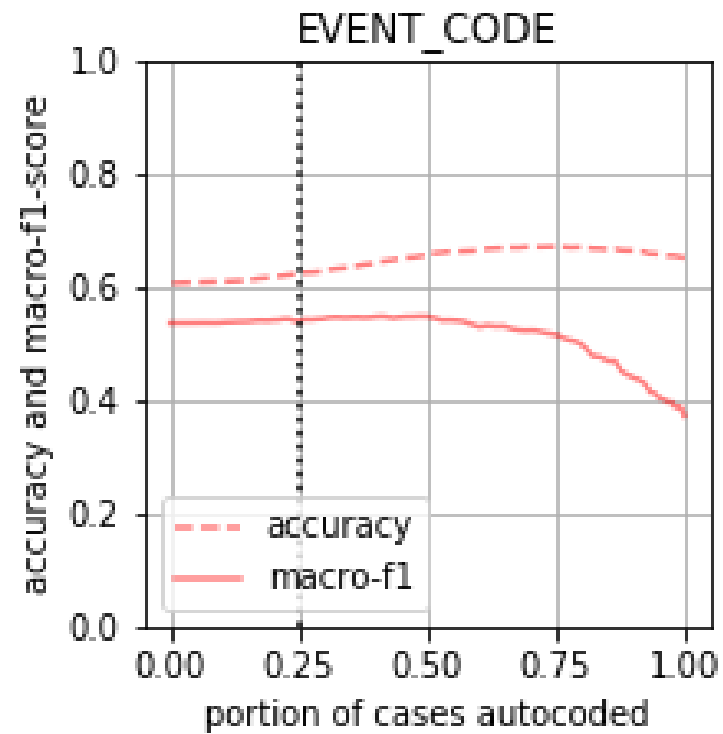
BLS

# Machine Learning vs. Manual Process



Accuracy

# The benefits of probabilistic models

■ Predicted Prob ≈ True Prob

▶ It mostly knows what it doesn't know

▶ Maybe a human knows?



event_code
bin_size=50
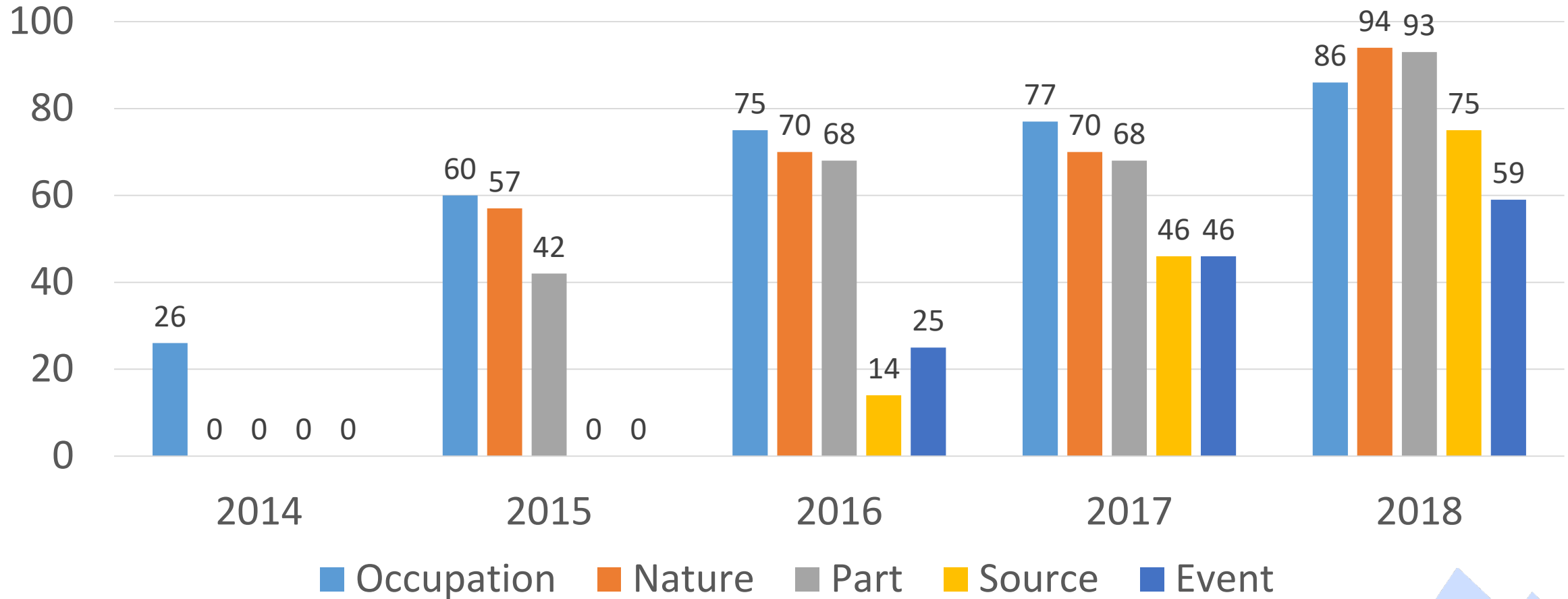
# Maximizing Quality by Simulating Possibilities

■ Gold + Human + Computer codes allows simulation

# % of codes automatically assigned to SOII



Grouped bar chart showing percentage of codes automatically assigned to SOII by year (2014–2018) for categories Occupation, Nature, Part, Source, and Event.

- 2014: Occupation 26, Nature 0, Part 0, Source 0, Event 0
- 2015: Occupation 60, Nature 57, Part 42, Source 0, Event 0
- 2016: Occupation 75, Nature 70, Part 68, Source 14, Event 25
- 2017: Occupation 77, Nature 70, Part 68, Source 46, Event 46
- 2018: Occupation 86, Nature 94, Part 93, Source 75, Event 59

Legend: Occupation, Nature, Part, Source, Event

# The Neural Network Autocoder

# Did it work?

## Accuracy



| | soc | nature | part | event | source | secondary source |
|---|---|---|---|---|---|---|
| human | 68.3 | 80.4 | 84.4 | 52.4 | 62.8 | 79.4 |
| LR autocoder | 75.0 | 88.3 | 88.4 | 61.3 | 66.7 | |
| NN autocoder | 78.6 | 91.1 | 91.9 | 69.8 | 75.8 | 85.0 |

■ human  ■ LR autocoder  ■ NN autocoder

# Things I wish someone had told me

- **Gold standard**
  - Not optional if you care about quality and replacing an existing process
  - It **must be blind** (reviewers 6x more likely to keep codes they see)
- **Not that hard to create**
  - Find an expert (or 2)
  - Ask them to recode your test set (without access to original codes)
  - Bigger is better but even 500 will get you a 95% CI +/- 4.5% accuracy

# Things I wish someone had told me

- You're not done once it's in production
  - ▶ Machine learning models need monitoring and maintenance
  - ▶ Neither is trivial
- Approach that's worked best so far
  - ▶ Hold back a "sample" for humans to code
  - ▶ Then recode with experts, and add to gold standard
  - ▶ Allows updating of human / computer accuracy metrics so you can maintain right mix

BLS

# Things I wish someone had told me

- Don't spend a lot of time trying every preprocessing, feature, and model possible
  - ▶ Most were designed for something else
  - ▶ Most don't matter

- My best model and feature ideas always came from looking carefully at the errors the model was making and working out why that would happen.

BLS

# What's next?

■ State-of-the-art NLP continues to advance rapidly

▶ Transfer learning with pretrained models

▶ BERT, RoBERTa, ALBERT

■ Sharing models with the public

▶ Differential privacy works here too!

■ Training staff

BLS

# Additional Resources

■ Tutorials

▶ Logistic Regression

https://github.com/ameasure/autocoding-class/blob/master/machine_learning.ipynb

▶ Neural Networks

https://colab.research.google.com/drive/1g3MVMCLOYshI_gaqMkDDj9gtG7yQQxib?ts=5c98e613

■ Papers

▶ https://www.bls.gov/osmr/pdf/st140040.pdf

▶ https://www.bls.gov/iif/deep-neural-networks.pdf

– Code: https://github.com/USDepartmentofLabor/soii_neural_autocoder

# Contact Information

**Alexander Measure**

Economist

Office of Safety and Health Statistics

www.bls.gov/iif

202-691-6185

measure.alex@bls.gov

BLS