

# Analysis of selected articles

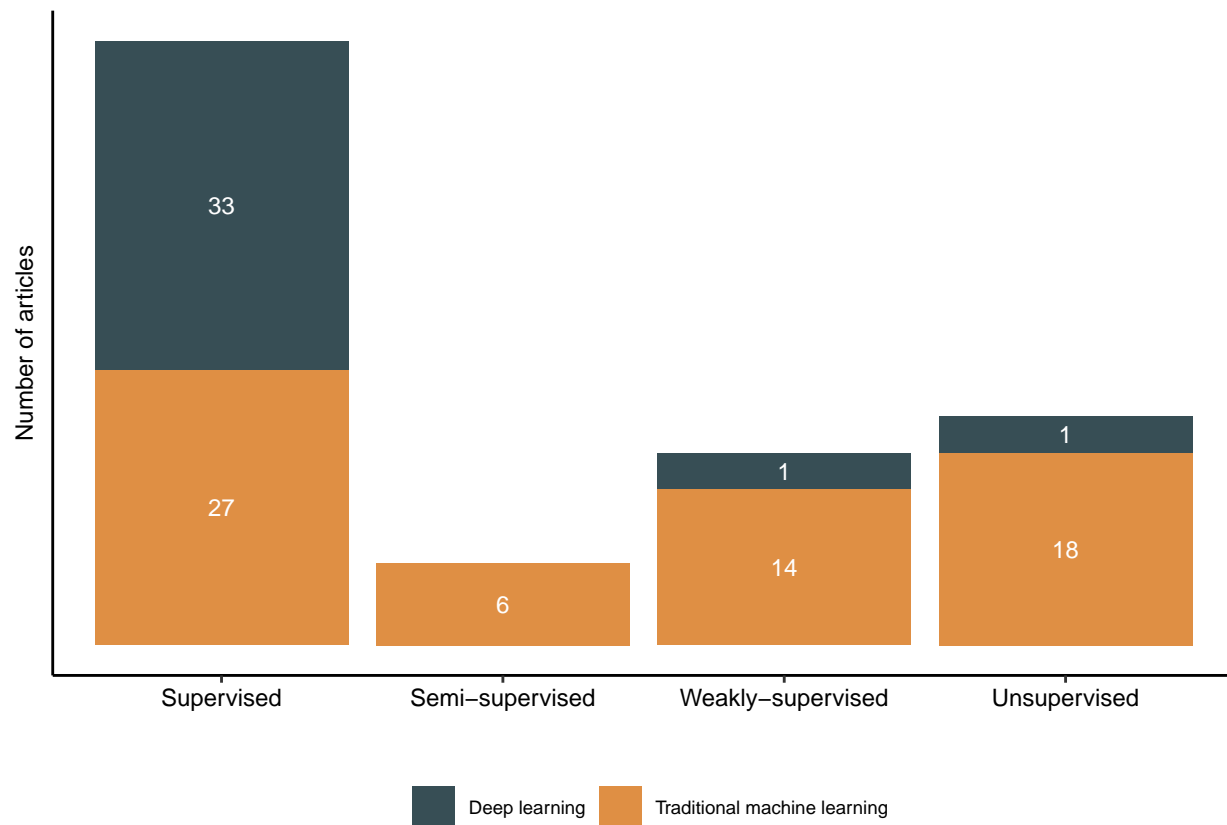
Last Updated: 11/17/2022

## Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Machine learning (ML) methods</b>                         | <b>2</b>  |
| 1.1      | ML paradigms . . . . .                                       | 2         |
| 1.2      | Traditional ML methods . . . . .                             | 2         |
| 1.3      | Deep learning (DL) methods . . . . .                         | 3         |
| <b>2</b> | <b>Phenotypes</b>  | <b>5</b>  |
| 2.1      | Phenotypes considered across ML paradigms . . . . .          | 5         |
| 2.2      | Unstratified summary of phenotypes considered . . . . .      | 6         |
| <b>3</b> | <b>Data sources</b>  | <b>7</b>  |
| 3.1      | Use of structured and unstructured data . . . . .            | 7         |
| 3.2      | Structured and unstructured data types . . . . .             | 8         |
| 3.3      | Terminologies . . . . .                                      | 9         |
| 3.4      | Natural language processing (NLP) software . . . . .         | 10        |
| 3.5      | Embeddings . . . . .   | 10        |
| 3.6      | Openly-available data . . . . .                              | 11        |
| 3.7      | Private data sources and demographics reporting . . . . .    | 13        |
| 3.8      | Institutions . . . . .                                       | 13        |
| 3.9      | Data sources summary across different ML paradigms . . . . . | 13        |
| <b>4</b> | <b>Reporting and evaluation</b>                              | <b>14</b> |
| 4.1      | Traditonal supervised ML vs. rule-based . . . . .            | 14        |
| 4.2      | Weakly-supervised ML vs. rule-based . . . . .                | 15        |
| 4.3      | Weakly-supervised ML vs. traditional . . . . .               | 16        |
| 4.4      | Deep ML vs. traditional . . . . .                            | 17        |

# 1 Machine learning (ML) methods

## 1.1 ML paradigms



## 1.2 Traditional ML methods

Table 1: Common traditional machine learning methods (Count > 1)

| ML                | Traditional ML method         | Count |
|-------------------|-------------------------------|-------|
| Supervised        | Random forest                 | 14    |
| Supervised        | Logistic regression           | 11    |
| Supervised        | SVM                           | 11    |
| Supervised        | L1 logistic regression        | 8     |
| Supervised        | Decision trees                | 4     |
| Supervised        | XGBoost                       | 4     |
| Supervised        | Naive Bayes                   | 3     |
| Weakly-supervised | PheNorm                       | 3     |
| Weakly-supervised | MAP                           | 2     |
| Weakly-supervised | Random forest                 | 2     |
| Unsupervised      | LDA                           | 5     |
| Unsupervised      | K-means                       | 4     |
| Unsupervised      | UPGMA Hierarchical clustering | 2     |

## [1] "There are 18 papers using multiple traditional machine learning methods"

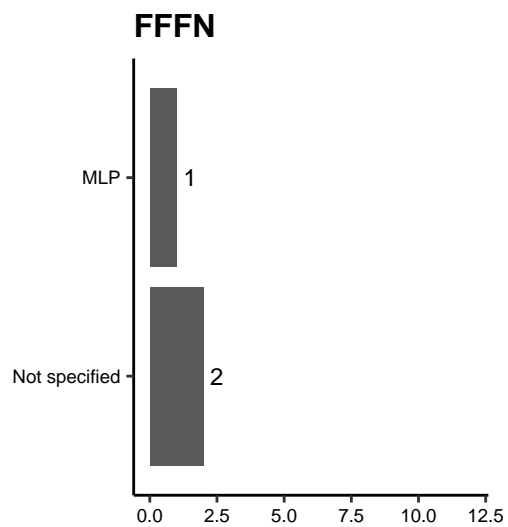
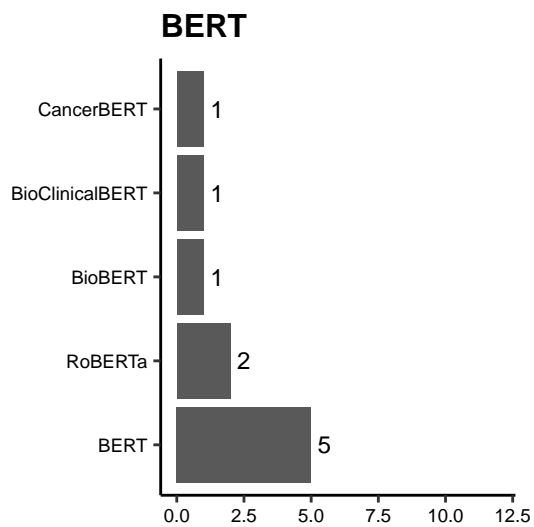
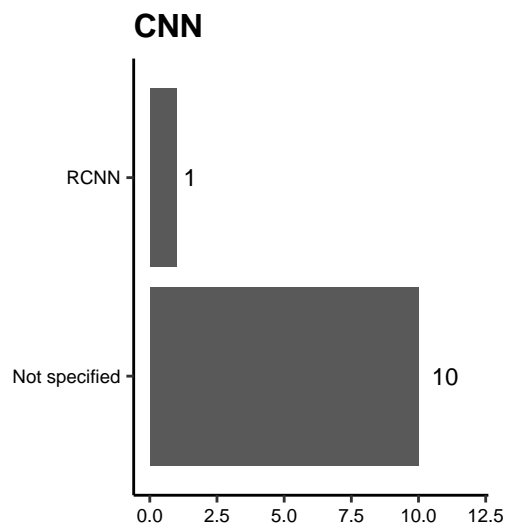
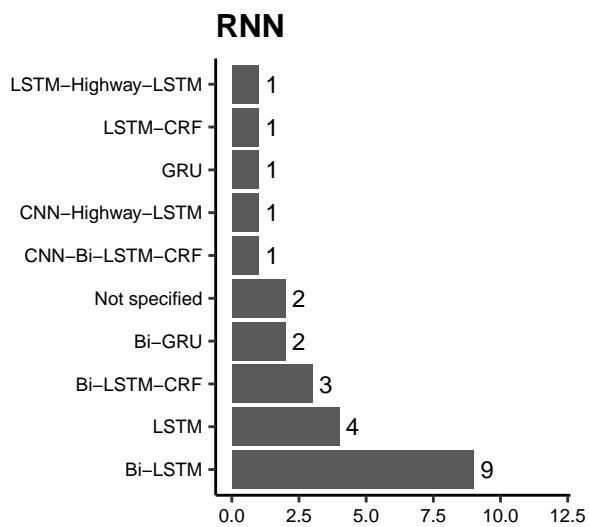
### 1.3 Deep learning (DL) methods

Table 2: Common deep learning methods (Count > 1)

| DL method | ML         | Count |
|-----------|------------|-------|
| BERT      | Supervised | 7     |
| CNN       | Supervised | 11    |
| FFNN      | Supervised | 3     |
| RNN       | Supervised | 19    |

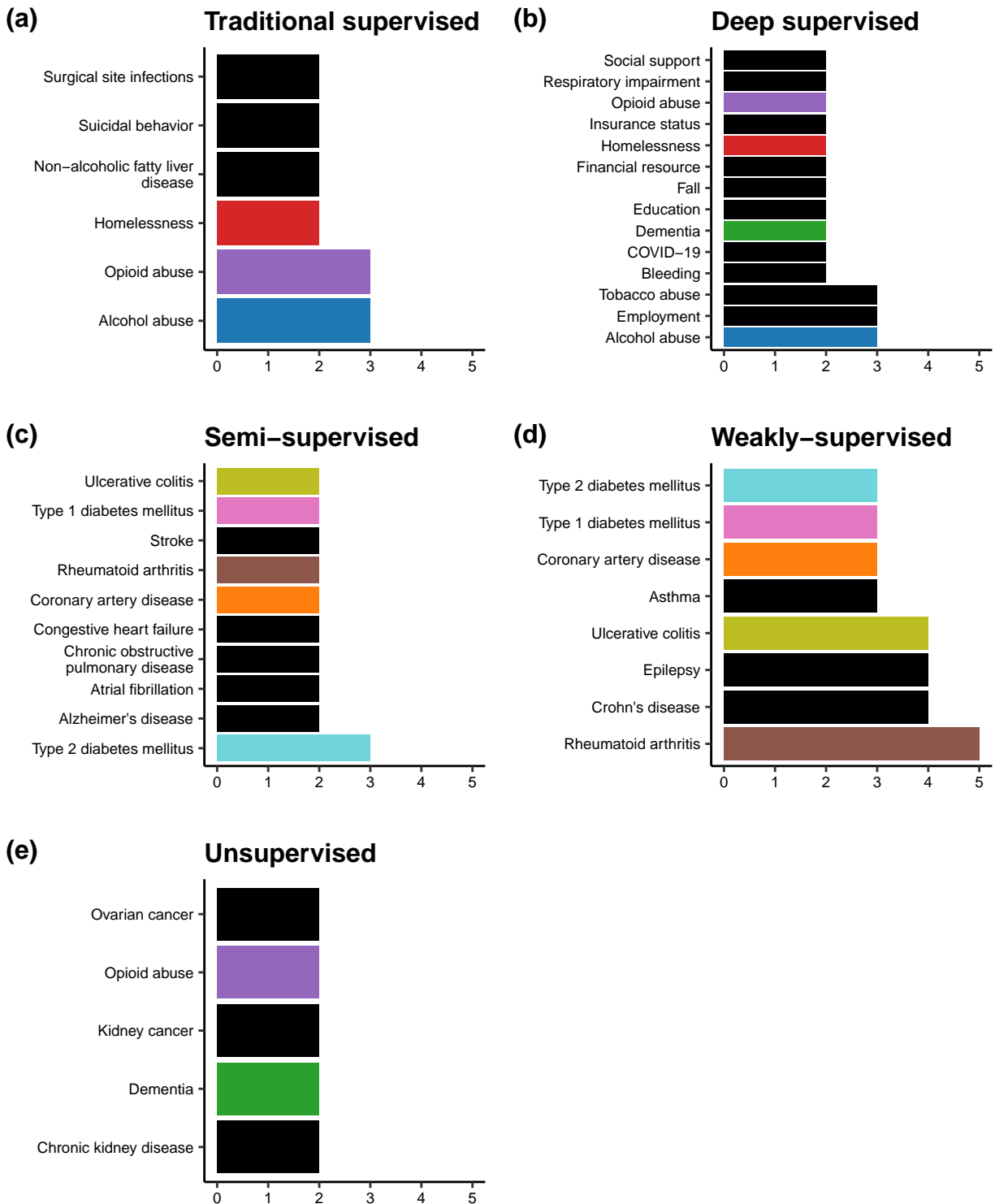
## [1] "There are 5 papers using multiple deep learning methods"

### 1.3.1 Neural network variants

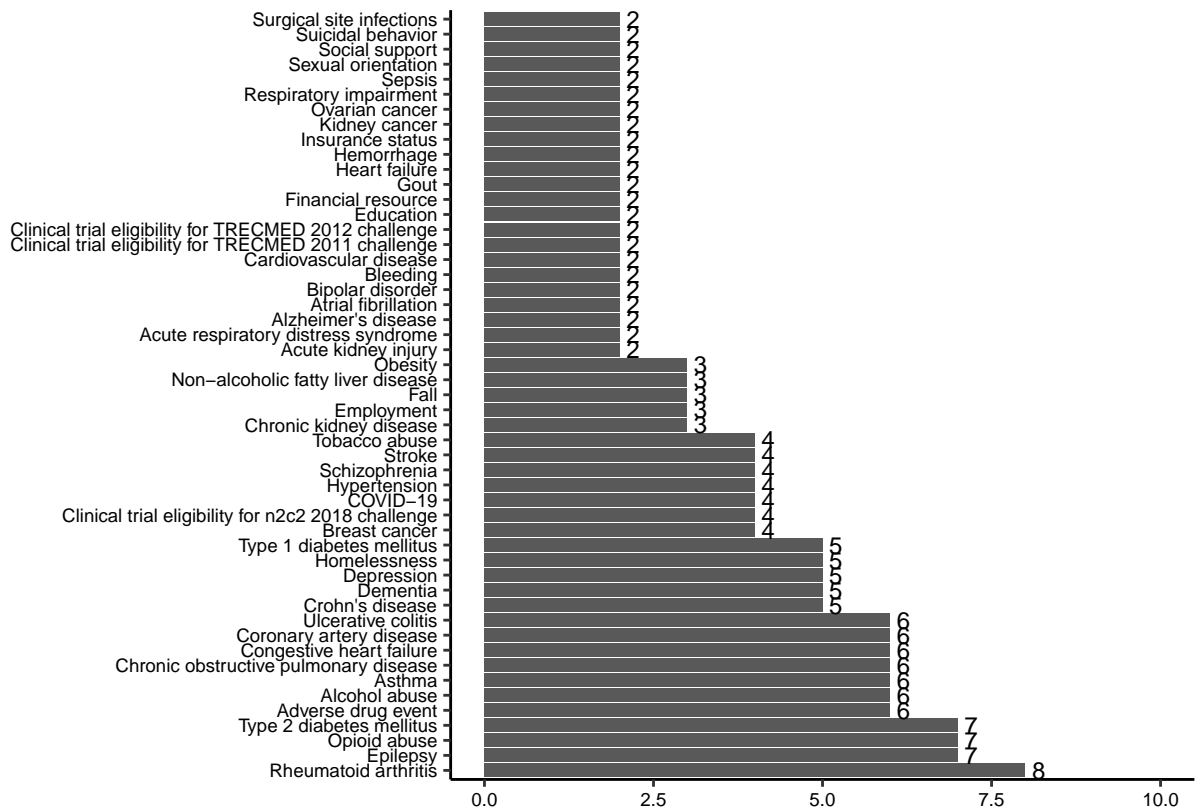


## 2 Phenotypes

### 2.1 Phenotypes considered across ML paradigms



## 2.2 Unstratified summary of phenotypes considered



### 3 Data sources

#### 3.1 Use of structured and unstructured data

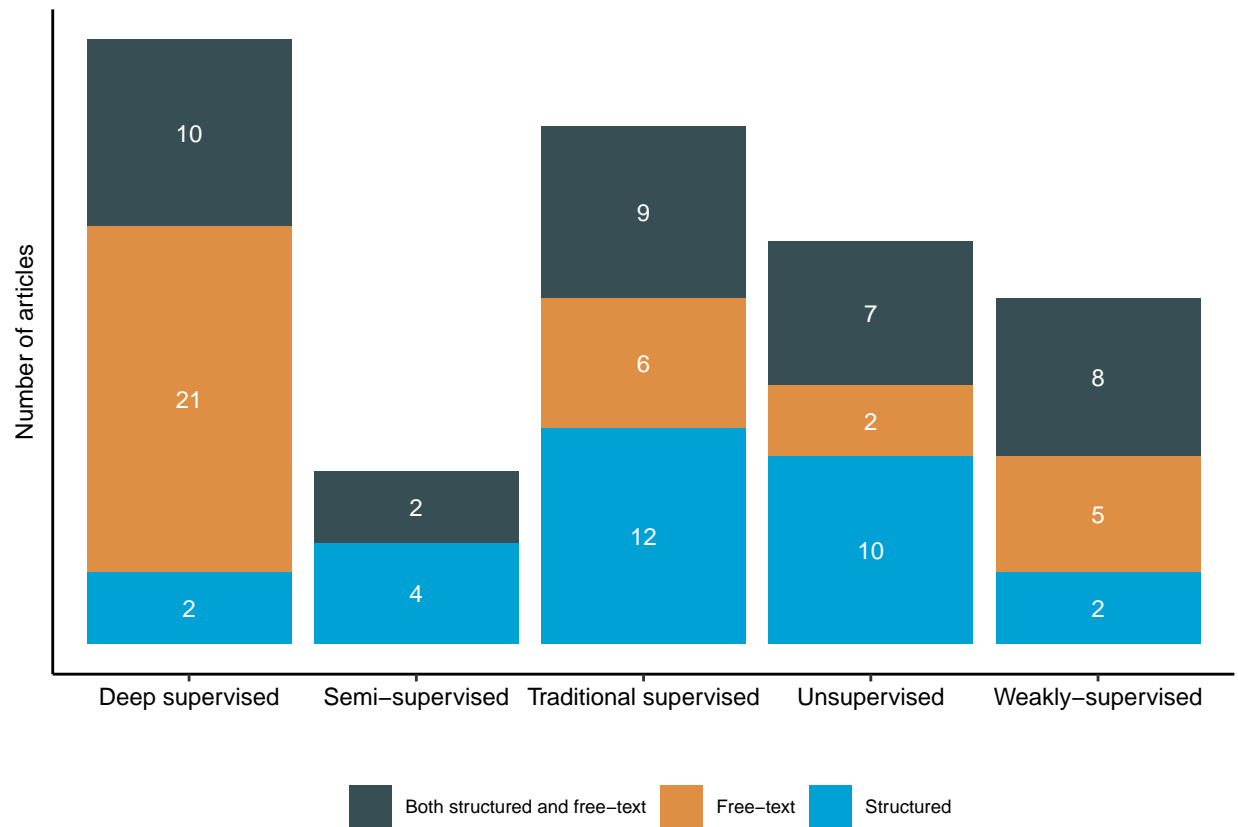


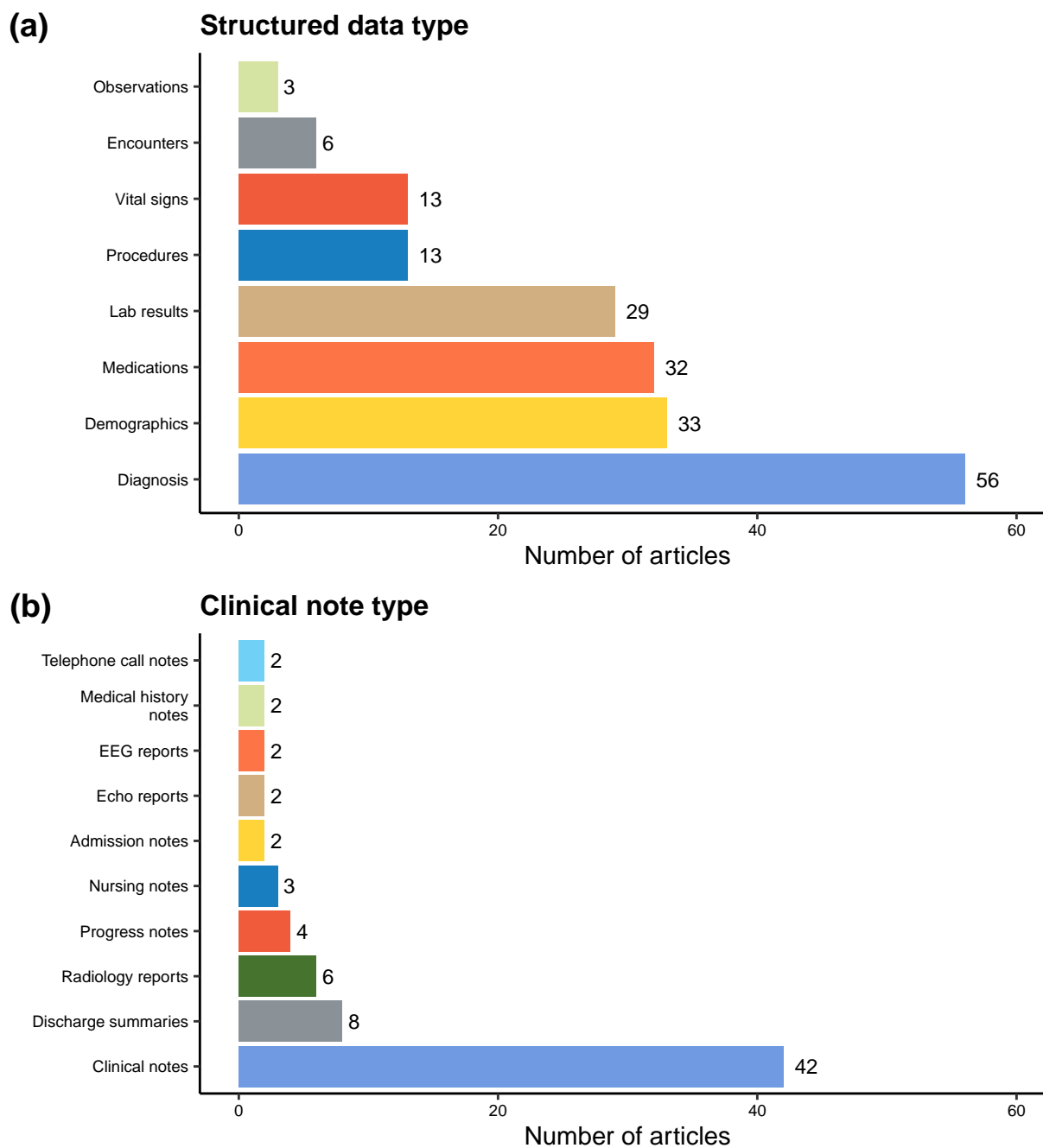
Table 3: Use of structured and unstructured data

| Data                          | Count |
|-------------------------------|-------|
| Both structured and free-text | 36    |
| Free-text                     | 34    |
| Structured                    | 30    |

### 3.2 Structured and unstructured data types

```
## [1] "There are 50 papers using multiple structured data types"
```

```
## [1] "There are 13 papers using multiple unstructured data types"
```





| Terminology<br>unnested                 | Supervised<br>Traditional<br>machine<br>learning | Unsupervised<br>Traditional<br>machine<br>learning | Supervised<br>Deep<br>learning | Weakly-<br>supervised<br>Traditional<br>machine<br>learning | Semi-<br>supervised<br>Traditional<br>machine<br>learning | Count |
|---|--|--|--------------------------------|---|---|-------|
| UMLS                                    | 11   | 3  | 8                              | 8   | 1   | 31    |
| ICD-9                                   | 6  | 5  | 4                              | 4   | 2   | 21    |
| ICD-9/10                                | 11   | 1  | 3                              | 0   | 2   | 17    |
| SNOMED-<br>CT                           | 2  | 3  | 4                              | 3   | 0   | 12    |
| RxNorm                                  | 3  | 1  | 2                              | 2   | 1   | 9     |
| CPT                                     | 2  | 0  | 3                              | 2   | 0   | 7     |
| Phecode                                 | 0  | 2  | 0                              | 3   | 2   | 7     |
| ICD                                     | 0  | 1  | 0                              | 4   | 0   | 5     |
| ICD-9-CM                                | 1  | 2  | 0                              | 1   | 0   | 4     |
| LOINC                                   | 3  | 0  | 0                              | 1   | 0   | 4     |
| ICD-10                                  | 0  | 0  | 1                              | 1   | 1   | 3     |
| ATC                                     | 2  | 0  | 0                              | 0   | 0   | 2     |
| (Anatomical<br>therapeutic<br>chemical) |  |  |                                |   |   |       |
| NDC                                     | 2  | 0  | 0                              | 0   | 0   | 2     |
| (National<br>drug<br>codes)             |  |  |                                |   |   |       |

### 3.3 Terminologies

## [1] "There are 37 papers using multiple terminologies"

| NLP software        | Supervised<br>Deep<br>learning | Weakly-<br>supervised<br>Traditional<br>machine<br>learning | Supervised<br>Traditional<br>machine<br>learning | Semi-<br>supervised<br>Traditional<br>machine<br>learning | Unsupervised<br>Traditional<br>machine<br>learning | Count |
|---------------------|--------------------------------|---|--|---|--|-------|
| cTAKES              | 8                              | 0   | 8  | 1   | 2  | 19    |
| NegEx               | 0                              | 2   | 3  | 0   | 1  | 6     |
| NILE                | 0                              | 5   | 1  | 0   | 0  | 6     |
| NLTK                | 4                              | 0   | 0  | 0   | 1  | 5     |
| MetaMap             | 1                              | 0   | 3  | 0   | 0  | 4     |
| Stanford<br>CoreNLP | 2                              | 0   | 0  | 0   | 0  | 2     |

### 3.4 Natural language processing (NLP) software

## [1] "There are 7 papers using multiple NLP software"

### 3.5 Embeddings

Embeddings were only used in deep supervised articles.

| Embedding training data       | Count |
|-------------------------------|-------|
| Unstructured EHR              | 11    |
| Biomedical literature         | 10    |
| MIMIC-III database (internal) | 7     |
| MIMIC-III database (external) | 6     |
| Wikipedia                     | 6     |
| Structured EHR                | 2     |

## [1] "There are 7 papers using multiple embedding training data"

| Embedding       | Count |
|-----------------|-------|
| Word2vec        | 19    |
| GloVe           | 6     |
| BERT            | 5     |
| RoBERTa         | 3     |
| BioBERT         | 2     |
| BioClinicalBERT | 2     |
| FastText        | 2     |
| Not specified   | 2     |

## [1] "There are 11 papers using multiple embedding training methods"

## 3.6 Openly-available data

### 3.6.1 Competition data

## [1] "There are 2 papers using multiple competition data"

| Competition<br>data name       | Supervised<br>Traditional<br>machine<br>learning | Supervised<br>Deep<br>learning | Count |
|--------------------------------|--|--------------------------------|-------|
| 2018 n2c2<br>track 2           | 0  | 6                              | 6     |
| 2018 n2c2<br>track 1           | 1  | 3                              | 4     |
| TRECMED<br>2011                | 1  | 1                              | 2     |
| TRECMED<br>2012                | 1  | 1                              | 2     |
| 2008 i2b2                      | 1  | 0                              | 1     |
| 2012<br>physionet<br>Challenge | 0  | 1                              | 1     |

| Data<br>source        | Supervised<br>Deep<br>learning | Supervised<br>Traditional<br>machine<br>learning | Weakly-<br>supervised<br>Deep<br>learning | Weakly-<br>supervised<br>Traditional<br>machine<br>learning | Unsupervised<br>Traditional<br>machine<br>learning | Count |
|-----------------------|--------------------------------|--|---|---|--|-------|
| MIMIC-III<br>database | 9                              | 1  | 1   | 1   | 3  | 15    |
| MTSamples<br>database | 1                              | 0  | 0   | 0   | 0  | 1     |

### 3.6.2 Other publicly available data sources

### 3.7 Private data sources and demographics reporting

## [1] "71 articles did not use openly available data"

## [1] "Among these 71 articles, 38 articles considered temporal phenotypes"

### 3.8 Institutions

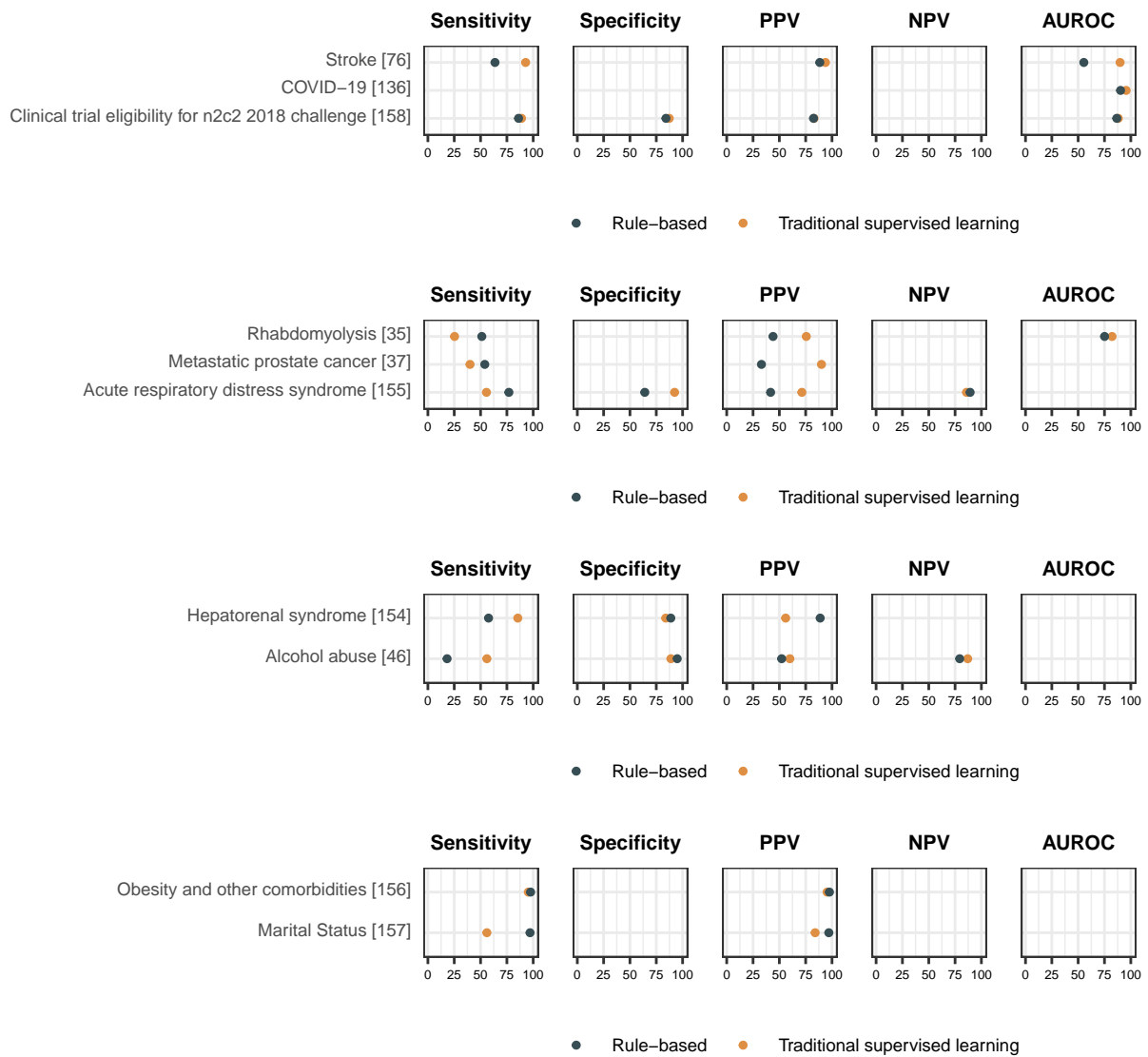
| Country    | Count |
|------------|-------|
| US         | 94    |
| France     | 2     |
| Canada     | 1     |
| China      | 1     |
| Germany    | 1     |
| Israel     | 1     |
| Italy      | 1     |
| Korean     | 1     |
| Netherland | 1     |
| Singapore  | 1     |
| Spain      | 1     |

### 3.9 Data sources summary across different ML paradigms

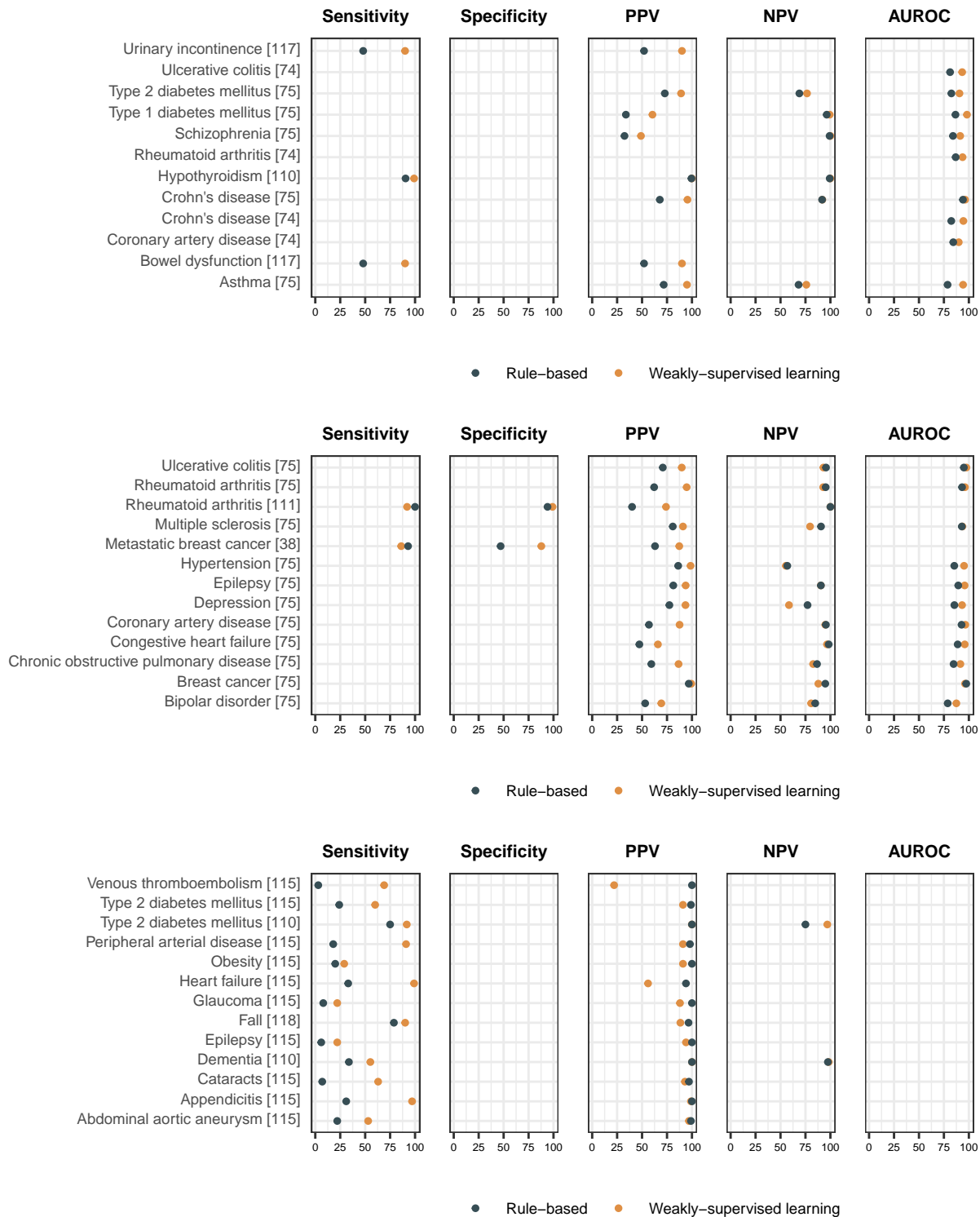
|       | Total<br>number<br>of papers | Used<br>free-text | Used<br>NLP<br>software | Used<br>competi-<br>tion data | Used<br>multisite<br>data | Used<br>open<br>data | Used<br>private<br>single-<br>site data | Compared<br>to rule-<br>based<br>algo-<br>rithms | Comapred<br>to tradi-<br>tional<br>ML | Reported<br>patient<br>demo-<br>graphic | Released<br>open<br>code |
|-------|------------------------------|-------------------|-------------------------|-------------------------------|---------------------------|----------------------|---|--|---------------------------------------|---|--------------------------|
| TSL   | 27                           | 15                | 14                      | 3                             | 1                         | 1                    | 22                                      | 10   | 0                                     | 13                                      | 4                        |
| DSL   | 33                           | 31                | 18                      | 11                            | 1                         | 9                    | 12                                      | 2  | 20                                    | 5                                       | 9                        |
| SSL   | 6                            | 2                 | 1                       | 0                             | 0                         | 0                    | 6                                       | 1  | 0                                     | 3                                       | 0                        |
| WSL   | 15                           | 13                | 10                      | 0                             | 3                         | 2                    | 10                                      | 8  | 1                                     | 4                                       | 3                        |
| USL   | 19                           | 9                 | 4                       | 0                             | 3                         | 3                    | 13                                      | 0  | 0                                     | 13                                      | 4                        |
| Total | 100                          | 70                | 47                      | 14                            | 8                         | 15                   | 63                                      | 21   | 21                                    | 38                                      | 20                       |

## 4 Reporting and evaluation

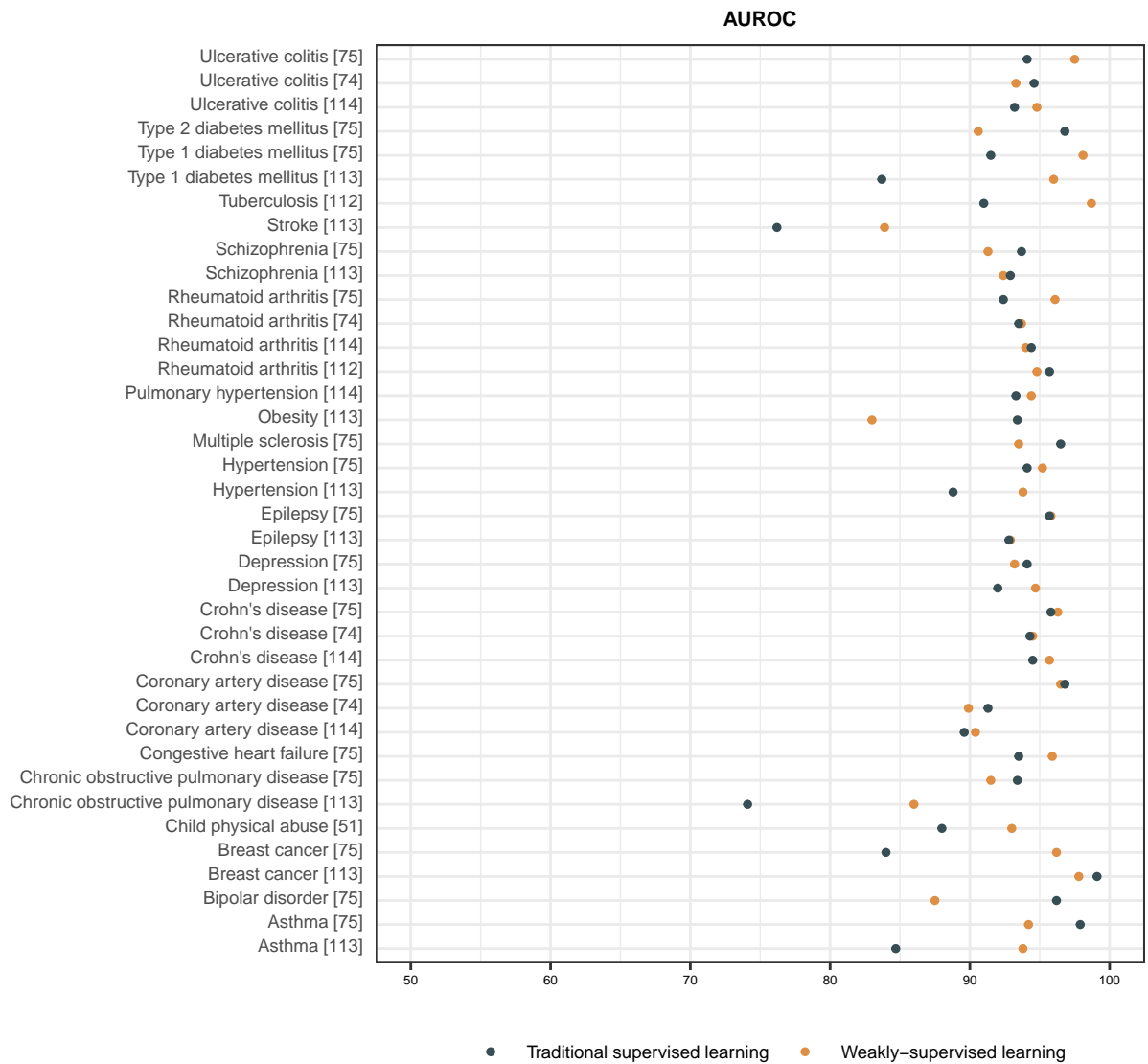
### 4.1 Traditonal supervised ML vs. rule-based



## 4.2 Weakly-supervised ML vs. rule-based

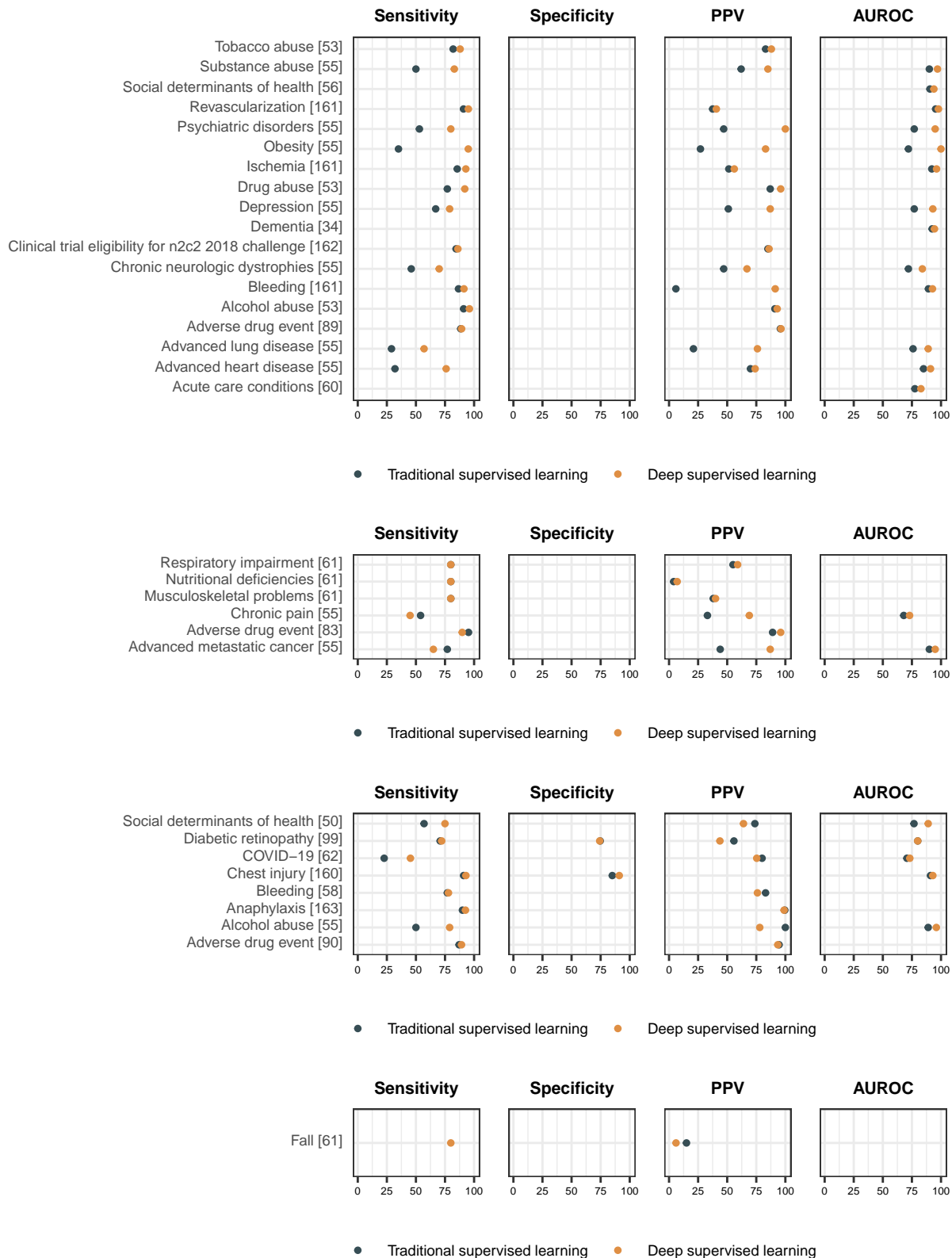


### 4.3 Weakly-supervised ML vs. traditional





## 4.4 Deep ML vs. traditional



| Model<br>performance<br>metrics                       | Supervised<br>Deep<br>learning | Supervised<br>Tradi-<br>tional<br>machine<br>learning | Weakly-<br>supervised<br>Deep<br>learning | Weakly-<br>supervised<br>Tradi-<br>tional<br>machine<br>learning | Semi-<br>supervised<br>Tradi-<br>tional<br>machine<br>learning | Count |
|---|--------------------------------|---|---|--|--|-------|
| Precision   | 26                             | 23  | 0   | 8  | 4  | 61    |
| Recall  | 25                             | 23  | 1   | 7  | 2  | 58    |
| AUROC   | 11                             | 15  | 1   | 10   | 5  | 42    |
| F-score   | 26                             | 9   | 0   | 7  | 0  | 42    |
| Specificity   | 6                              | 11  | 1   | 1  | 0  | 19    |
| Accuracy  | 4                              | 8   | 1   | 4  | 0  | 17    |
| NPV   | 1                              | 7   | 0   | 5  | 2  | 15    |
| AUPRC   | 4                              | 2   | 0   | 2  | 0  | 8     |
| Calibration<br>plots                                  | 2                              | 3   | 0   | 0  | 0  | 5     |
| Log loss  | 1                              | 1   | 0   | 0  | 1  | 3     |
| Brier<br>score  | 1                              | 1   | 0   | 0  | 0  | 2     |
| Hamming<br>loss                                       | 2                              | 0   | 0   | 0  | 0  | 2     |
| Matthews<br>Correla-<br>tion<br>Coeffi-<br>cient      | 1                              | 1   | 0   | 0  | 0  | 2     |
| Normalized<br>dis-<br>counted<br>cumula-<br>tive gain | 1                              | 1   | 0   | 0  | 0  | 2     |