Detecting Sockpuppetry on Wikipedia Using Meta-Learning

Anonymous ACL submission

Abstract

Malicious sockpuppet detection on Wikipedia is critical to preserving access to reliable information on the internet and preventing the spread of disinformation. Prior machine learning approaches rely on stylistic and meta-data features, but do not prioritise adaptability to author-specific behaviours. As a result, they struggle to effectively model the behaviour of specific sockpuppet-groups, especially when text data is limited. To address this, we propose the application of meta-learning, a ma-013 chine learning technique designed to improve performance in data-scarce settings by training models across multiple tasks. Meta-learning optimises a model for rapid adaptation to the writing style of a new sockpuppet-group. Our results show that meta-learning significantly en-019 hances the precision of predictions compared to pre-trained models, marking an advancement in combating sockpuppetry on open editing platforms. We release an updated dataset of sockpuppet investigations to foster future research in both sockpuppetry and meta-learning fields.

1 Introduction

011

017

021

024

025

027

034

039

042

Over recent years, social media sites have seen a steady increase in the presence of fake accounts (Khaled et al., 2018). These accounts are often used to spread fake news and seed distrust for political gain (Shu et al., 2017). Wikipedia is not immune to such attacks: Saez-Trumper (2019) investigate political and religious groups imposing their narratives on articles. Attacks on Wikipedia are particularly threatening as articles often serve as the ground-truth for automated fact checking systems; used to combat disinformation on other platforms (Thorne et al., 2018).

Changes to articles on Wikipedia are a collaborative process, where decisions are made via the consensus of editors. Malicious users undermine this process through Sockpuppetry: the use of multiple accounts to stack votes, fake majority support

of a view or make a counter-perspective look absurd (Saez-Trumper, 2019). They have been used to vandalise articles (Solorio et al., 2013a), support political views (Kumar et al., 2017) or improve personal standing (Stone and Richtel, 2007).

043

045

047

051

053

054

057

059

060

061

062

063

064

065

066

067

069

070

071

072

073

074

076

077

078

079

081

Many machine learning approaches have been proposed, including linking sockpuppet accounts through their writing style (Solorio et al., 2013a; Sakib and Spezzano, 2022), called authorshipattribution. However, when the available text is scarce, it is difficult to profile an author accurately (Eder, 2013). In the case of sockpuppet detection, where available text samples are short (Solorio et al., 2013b), this makes achieving good performance difficult. Previous approaches (Solorio et al., 2013a; Sakib and Spezzano, 2022) manage this challenge by merging the corpus of sockpuppet investigations into a single dataset of sockpuppet behaviour. A model trained on this dataset then learns the writing style of sockpuppets as a whole. Whilst the model can be later fine-tuned, it will not be sensitive to author-specific features.

Meta-learning instead leverages prior experience to perform well on limited data. Rather than merging the corpus of investigations together, it considers each a separate learning task. The meta-model learns a general understanding of sockpuppets that can adapt to the behaviour of an unseen puppetmaster, the user behind a group of sockpuppets.

In this study, we are the first to apply metalearning to the problem of sockpuppet detection. Our work makes three main contributions:

Evaluate the application of meta-learning to the task of sockpuppet detection on Wikipedia We find that learning over a distribution of tasks significantly improves prediction precision over pretrained approaches. This outcome is valuable for sockpuppet detection, where confidence in positive identifications is paramount. Our approach is applicable to other online communities.

Construct and publicly release a dataset of sockpuppet investigations on Wikipedia Our dataset¹ improves upon existing datasets which are either outdated (Solorio et al., 2013a), unreleased publicly (Kumar et al., 2017) or do not preserve investigation structure (Sakib and Spezzano, 2022).

Formulate a more realistic task definition Previous approaches (Sakib and Spezzano, 2022) train on data from any number of accounts within a sockpuppet-group, preemptively revealing any deceptive efforts made by a puppetmaster to the model. Our model is fine-tuned on just one accused user, as it would be when deployed.

2 Related Work

084

094

096

098

101

102

103

105

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

2.1 Sockpuppetry

Sockpuppetry is typically described as the use of multiple accounts by a single user for deceptive or malicious purposes (Zheng et al., 2011; Solorio et al., 2013a; Bu et al., 2013; Liu et al., 2016; Sakib and Spezzano, 2022), however, not all sockpuppets are malicious. Kumar et al. (2017) provide a more general definition: A sockpuppet is any account controlled by a user with at least one other account. The set of these accounts is referred to as a sockpuppet-group, and their controlling user their *puppetmaster*. This is the definition used here, with one small amendment: that these accounts be, at some point, operated concurrently. This adjustment distinguishes the task of sockpuppet detection from that of ban evasion, where secondary accounts are created strictly after the primary accounts are banned (Niverthi et al., 2022).

2.1.1 Motivation

Malicious users use sockpuppets to vandalise Wikipedia pages (Solorio et al., 2013a), propagandise political views (Kumar et al., 2017; Afroz et al., 2012), or improve their own public image (Owens, 2013). Sockpuppets undermine collaboration on Wikipedia through false majority opinions, vote stacking (Solorio et al., 2013a) and *Straw man socks*, which argue easily refuted opposing arguments to discredit opposition (Kumar et al., 2017).

2.1.2 Detection

The current approach to detecting sockpuppetry on Wikipedia is manual²: Users argue their case

²https://en.wikipedia.org/wiki/Wikipedia: Sockpuppet_investigations before a presiding administrator, who may supplement evidence with technical logs. Once a verdict is reached, guilty accounts are suspended and the investigation is archived. 128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

Many automated approaches have been proposed to support this process. Authorship Attribution (AA) determine whether the intent and writing style of accounts are similar enough to be the same user. Linear classifiers with manually selected authorship features have achieved consistent results (Solorio et al., 2013a; Bu et al., 2013; Liu et al., 2016; Sakib and Spezzano, 2022). They use lexical, structural and syntactic features (Bu et al., 2013). AA classifiers struggle when only given short pieces of text (Shrestha et al., 2017), which is typically all that is available in the case of sockpuppet detection (Solorio et al., 2013a). Features may also require domain specific selection (Kotsiantis et al., 2007), and typically act under the assumption that the authors are not attempting to evade detection (Solorio et al., 2013a). Faced with adversarial authors, commonly used authorship features can be easily evaded (Brennan et al., 2012).

Meta-data approaches focus on the behaviour of users. Tsikerdekis and Zeadally (2014) identified that the number and time between edits deviates from that of normal users over time. Metadata approaches typically do not require pairwise comparison between accounts, reducing computational complexity. Kumar et al. (2017) highlight six points of divergence from usual user activity. These features can be combined with authorship features for improved performance (Solorio et al., 2013a; Sakib and Spezzano, 2022). In all prior literature, approaches consider a single model that classifies users or edits as belonging to a sockpuppet or not, rather than each investigation being a separate task (Solorio et al., 2013a; Bu et al., 2013; Liu et al., 2016; Sakib and Spezzano, 2022).

2.2 Meta-learning

Meta-learning research focuses on the problem of "Learning to Learn". In this setting, a machine learning model gains experience over a collection of tasks, rather than just one, and in doing so improves its performance on future tasks (Hospedales et al., 2020). Hospedales et al. (2020) define *Baselearning* as the inner learning algorithm solving a task, such as authorship attribution. *Meta-learning* is an outer learning algorithm which updates the inner algorithm according to its own *meta-objective*, typically quick adaption to new tasks (Finn et al.,

¹Anonymous link

2017; Snell et al., 2017; So, 2021).

179

180

181

188

189

192

193

194

196

197

199

205

210

211

212

213

214

215

216

217

218

219

221

224

227

Meta-learning has been successful in many domains, such as image classification (Antoniou et al., 2019), sentiment analysis (Liang et al., 2023), and text classification (Bansal et al., 2021). Tian et al. (2023) investigated the approach to detect statesponsored trolls. Beyond reducing data dependence, other limitations of deep neural networks, such as unsupervised learning performance, may also be improved (Hospedales et al., 2020).

Meta-Learning approaches are made up of several families. Gradient-based approaches use gradient descent to update a model's parameters to minimise the loss according to a meta-objective. These approaches are model-agnostic, making them advantageous over Metric and Model-based approaches that make restrictions on model architecture. The foremost approach is MAML (Finn et al., 2017) and its successors (Antoniou et al., 2019; Triantafillou et al., 2020; Finn et al., 2018; Rajeswaran et al., 2019). A related approach is Reptile (Nichol et al., 2018) that requires significantly less compute whilst achieving similar performance (Vinyals et al., 2016). An advantage of Reptile over MAML is that it does not require a train-test split for each training task (Nichol et al., 2018), allowing more data to be used in training.

3 Method

We provide two task definitions. The first is a description of the *base-learning* problem, which considers training and evaluating a classifier on a single task. This is also referred to as the *inner-loop* in the context of meta-learning. The second description is of the meta-learning problem, and describes how a meta-model is learnt across a distribution of baselearning tasks. This is also called the *outer-loop*.

3.1 Base-learning

The base-learning task is a binary classification problem. As input, the model will receive two data sources: the article *page* and *message* describing the contribution. The model outputs a classification, identifying the contribution as either a positive (made by a sockpuppet), or negative sample.

In a deployed setting, there is no given list of confirmed sockpuppets, only a set of accused accounts. Therefore, a model may only be trained on the contributions of a single accused user, which may then be tested on the contributions of the remaining accused accounts. We make similar restrictions on



Figure 1: Dataset Topology

our training data: For each investigation, we define the *puppetmaster* as the sockpuppet with the most contributions. A model is then trained on their contributions. This model is assessed by its ability to distinguish the contributions of the remaining sockpuppets from the samples of non-sockpuppets. 228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

253

254

255

257

258

259

260

261

262

263

264

265

We call the set of puppetmaster samples the train set, and the set of sockpuppet samples the test set. Negative samples are split between the two sets proportionally. Due to how these sets are created, test sets much larger than the train set are common.

A validation set is split from the train set, containing 20% of the available samples and maintaining the same proportion of positive and negative samples. This set is used during base-learning to provide feedback as the model is being trained, and to prevent over-fitting via early stopping.

3.2 Meta-learning

The meta-learning problem considers a distribution of tasks (Finn et al., 2017). This distribution is split into two sets, meta-train and meta-test. The meta-train set is used for the meta-learning process, whilst the meta-test set will be used to evaluate how well the meta-learned model performs. The performance of the models on these tasks is what will be reported in Section 7.

Figure 1 depicts the dataset topology, where each row represents the samples of a task. The tasks are split into the meta-train and meta-test sets, and each task is split into train, validation and test sets.

We make several restrictions on the shape of this distribution. To ensure that each task has an over representation of negative samples, we limit the distribution to only include tasks with a negative to positive ratio of at least one. We also ensure that each task has at least ten puppetmaster samples and five sockpuppet samples. These restrictions are derived from the model architecture, which uses a 266triplet contrastive loss function that requires at least267two positive samples in each task. By ensuring each268task contains at least ten puppetmaster samples, we269guarantee a valid validation set. These restrictions270reduce the total number of tasks from 23, 610 to27113, 549.272the meta-train set, and the rest (1, 355) become the273meta-test set. To compare the model with non-meta-274learning approaches, the meta-train set will either275be used for the pre-trained approach, or, where no276pre-training is necessary, will not be used at all.

Whilst training on the meta-train set of tasks, approaches are not required to maintain the distinction between task specific train, test and validation sets. These sets are only preserved to the extent that they are required for the meta-learning or pretrained approach. The training process on each meta-test task is kept constant throughout each approach. Each model is given a maximum of ten epochs to train on the new task before predictions must be made. The metrics of each task in the meta-test set are then averaged to create the overall metrics for the approach. Each approach is run three times, and their mean and standard deviation reported in Section 7.

3.3 Meta-learning Strategy

279

281

286

290

291

294

295

304

306

307

311

312

We use Reptile as our meta-learning strategy. This is because it is less computationally complex and similarly performant to MAML (Nichol et al., 2018; Rajeswaran et al., 2019). Compared to metric and model-based approaches, Reptile has the additional benefit of being model-agnostic, accommodating versatile model architecture.

We use the serial implementation that updates the parameters directly through linear interpolation. It works by adapting a clone of the meta-model to a task, and then moving the parameters of the original meta-model toward the adapted parameters through linear interpolation. The Reptile algorithm is provided in Algorithm 1 in Appendix A.

4 Model Architecture

A diagram of our approach is given in Figure 2. The sample represents a positive or negative contribution. The two textual inputs of the contribution, *page* and *message*, are concatenated using the appropriate separator token. We use RoBERTa³ (Liu et al., 2019), a pre-trained transformer with frozen



Figure 2: Model Topology

parameters to generate a matrix of contextualised word embeddings. This approach is typical of recent authorship attribution classifiers (De Langhe et al., 2024; Huertas-Tato et al., 2022; Ai et al., 2022; Rivera-Soto et al., 2021). 313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

330

331

332

333

334

335

This matrix is then fed to a transformer encoder, optimised using Adam (Kingma and Ba, 2014). This encoder is where the majority of task learning takes place, and is the model that will be trained using meta-learning. During meta-learning, the β_1 parameter of Adam is set to 0, as recommended (Nichol et al., 2018). The interpolation rate for the Reptile algorithm is set at 0.2, and a total of five gradient steps are taken on each task.

Once the authorship embeddings are produced, a neural network interprets the embedding to produce a logit, which is later interpreted to produce discrete classifications. The classifier has two fully connected layers of dimension 768 with dropout. It also uses the Adam optimiser, and is trained using a cross-entropy loss function. The classifier is trained on the embeddings produced after the encoder has been trained.

³https://huggingface.co/sentence-transformers/ all-distilroberta-v1

4.1 Loss Functions

336

339

341

343

345

352

354

364

367

371

374

376

378

379

The encoder model is trained using triplet margin loss (Schroff et al., 2015). A contrastive loss function is typical in natural language tasks comparing document similarity (Pennington et al., 2014; Devlin et al., 2019), and also has prior success in contrastive authorship models (Huertas-Tato et al., 2022). The triplets are created by iterating through all the positive samples as the anchor, and randomly selecting another negative and positive. Typically, the anchor may be drawn from both classes, however we limit triplets to the positive sockpuppet class. This is because the authors in the negative samples are all different, and therefore should occupy different regions in the embedding space only positive samples should be clustered together. For the classifier models that interpret the embeddings, we use binary cross-entropy loss.

4.2 Hyper-parameters

We tuned model hyper-parameters using the Optuna⁴ framework. The encoder and classifier models were tuned together, optimising for AUROC. 100 trials were run, where the model was trained over ten randomly selected tasks from the metatrain distribution of tasks. The performance of each on their respective tasks was averaged and provided to the optimiser as feedback. We will release hyperparameters along with the implementation.

We performed three optimisations, one for the encoder and classifier models, a second for the classifier component of the RoBERTa baseline (Section 5.1), and a third for the Reptile hyper-parameters. The tuned hyper-parameters are provided in Appendix A.

4.3 Training Parameters

When training both the classifier and encoder on the meta-test set of tasks, each was given a maximum of ten epochs, where each epoch is one complete pass through the train set. This is a limitation of the time and computing resources available.

We used a variable batch size strategy that scaled with the length of the task. This catered for smaller tasks whilst still allowing larger tasks to benefit from the stability and speed of larger batch sizes. We used early stopping based on the validation loss with a patience of 3 epochs. During the metalearning stage, the model was trained over five epochs of the meta-train set of tasks. On each task, Reptile performed five gradient steps before the parameters were updated using an interpolation rate of 0.2. At the end of each epoch, the model was saved along with the sum of the training loss of each task in that epoch. The best performing model relative to the validation loss was selected.

385

386

389

390

391

392

393

394

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

5 Metrics

Reported metrics undergo an aggregation process. For each experiment, the results of the approach on the test set of each task in the meta-test set (see Figure 1) are computed. The results of each task are then averaged to find the overall result of the approach for that experiment. Three experiments of each approach are run. The metrics across each experiment are averaged, and the standard deviation provided as a confidence interval.

The main metrics for comparison between approaches should be the area-based metrics, AU-ROC and AUPRC. This is because they do not require a specific threshold to be decided, which may distort the appearance of classifier performance. We also provide the F1-Score and F0.5-Score. The first is justified through use in previous literature (Sakib and Spezzano, 2022; Solorio et al., 2013a,b), whilst the second presents a balance between recall and precision more relevant to the deployment environment, where false positives are strongly discouraged⁵. We also provide the accuracy, precision, and recall of each model as supplementary metrics.

5.1 Baselines

We consider several baselines, including two trivial baselines (random and majority classifiers) previously used by (Solorio et al., 2013a). In the case of the majority baseline, the class predicted is based on the training dataset.

RoBERTa Baseline To assess whether the encoder model itself provides a significant improvement, we train a simple binary classifier on the sentence-level RoBERTa embeddings. This changes the model architecture by reducing the output from the frozen RoBERTa model from a two-dimensional matrix to a one-dimensional vector. The vector is then fed directly into a fully connected neural network classifier. Huertas-Tato et al. (2022) employed a similar baseline to assess their authorship representation learner.

⁴https://optuna.org/

⁵https://en.wikipedia.org/w/index.php? title=Wikipedia:Sockpuppet_investigations/SPI/ Administrators_instructions&oldid=1173289303

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

474

475

476

430 Non-meta-learning Approach To isolate the ef431 fects of meta-learning, we also test our approach
432 without it. This model will follow the same training
433 approach as the Meta-learned model on test tasks.

434 Pre-trained Approach The pre-trained approach
435 trains our model on a merged dataset of the meta436 train set of tasks. This is the approach used in
437 prior literature (Solorio et al., 2013a; Sakib and
438 Spezzano, 2022). At test time, this approach will
439 be fine-tuned on tasks in the meta-test set.

Upper Limit As a significant portion of all con-440 tributions do not have any text in their message 441 feature (24.45% of all contributions, 63.64% of 442 which are positive samples), these contributions 443 are indistinguishable from one another using the 444 provided features. Therefore, the upper limit of 445 446 performance is more accurately defined as the perfect classifier on all contributions where a message 447 value is present, and a random classifier otherwise. 448

6 Dataset

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

To test our meta-learning approach, we create a dataset of sockpuppet detection tasks. Each task is a discrete problem, where writing samples from a single sockpuppet-group must be separated from writing samples of non-sockpuppets.

Tasks consist of writing samples from both sockpuppet and non-sockpuppet users. Edits to Wikipedia are called *contributions*, and contain a *message* component where the user can describe their edit. As in previous approaches (Solorio et al., 2013a; Sakib and Spezzano, 2022), we use these contribution messages as the writing samples. Negative samples are contributions from nonsockpuppet accounts drawn from the same time and article distribution as the sockpuppet-group.

The final dataset consists of 23, 610 tasks. For each contribution, we provide the timestamp, revision ID, ID of the preceding contribution, user name, article title, contribution message, and a binary label. The dataset is made publicly available⁶.

6.1 Data collection

To collect the positive samples, the confirmed Wikipedia sockpuppets page⁷ was crawled using a combination of Pywikibot⁸ and MediaWiki⁹ API

⁷https://en.wikipedia.org/wiki/Category: Wikipedia_sockpuppets calls, which extracted each investigation page and the contributions of each confirmed sockpuppet.

Negative samples for each task were collected from the same articles and within the active time period (first and last contribution) of the task's sockpuppet-group. For each positive sample, two random timestamps were drawn, and the next ten contributions made after each was collected. From each set of ten, the first valid (non-duplicate, nonsockpuppet) sample was selected. Multiple samples from the set of ten were not collected, so as to maintain a uniformly random temporal distribution. If a set of ten contained no valid negative samples, the attempt was abandoned. This occurred in cases where the active time period was very short, or the articles were inactive or new.

In reality, there is a class imbalance between the number of sockpuppet and non-sockpuppet accounts. Negative contributions were thus over sampled. Arbitrarily, an ideal ratio of two negative samples to each positive was set. It is unclear how an informed estimate might be reached: The true ratio of genuine users to sockpuppets would not only be too extreme to replicate or learn with, but investigations only occur on accused users, not all users. Ideally one might ascertain the ratio of accused sockpuppets to confirmed sockpuppets, however as the investigations of falsely accused accounts are not archived, this is impractical to obtain.

For investigations that did not reach the ideal ratio of two negatives for each positive, a second identical pass was performed. This strategy oversampled articles with more non-sockpuppet editors to make up for the shortfall. This yielded some tasks with up to four negatives for each positive.

For a variety of reasons, 664 investigations failed to collect any negative samples. These investigations were retained in the dataset, but were excluded from any experiments. 984 investigations contained no positive samples. These were removed, and a list of empty investigations provided alongside the rest of the dataset.

7 Results

Our results are presented in Table 1. Metrics are measured as the mean across all three test runs with the standard deviation as error bounds. The classification threshold used to compute predictions from logits for applicable metrics were computed at a task level, using the optimal threshold relative to the F0.5-Score on the task validation set. A T-test

⁶Anonymouslink

⁸https://github.com/wikimedia/pywikibot

⁹https://www.mediawiki.org/wiki/MediaWiki

Approach	AUROC	AUPRC	F1-Score	F0.5-Score	Accuracy	Precision	Recall
Random	50.10 ± 0.14	50.85 ± 0.09	40.34 ± 0.11	36.52 ± 0.12	50.10 ± 0.16	34.46 ± 0.12	50.05 ± 0.15
Majority	-	-	-	-	65.60 ± 0.00	-	-
RoBERTa	65.70 ± 0.00	50.45 ± 0.03	57.97 ± 0.01	57.52 ± 0.06	66.98 ± 0.06	59.54 ± 0.13	67.63 ± 0.17
Standard Enc.	68.33 ± 0.09	50.67 ± 0.33	60.05 ± 0.18	58.73 ± 0.16	68.90 ± 0.07	59.72 ± 0.32	69.88 ± 0.34
Pretrained Enc.	62.74 ± 0.02	44.80 ± 0.19	57.49 ± 0.13	52.90 ± 0.12	62.79 ± 0.05	51.45 ± 0.25	74.76 ± 0.28
Reptile Enc.	$78.98 \pm 0.12 \ast$	$62.21 \pm 0.08 \ast$	$67.46 \pm \mathbf{0.53^*}$	$67.89 \pm 0.17^*$	$77.51 \pm \mathbf{0.19*}$	$69.43 \pm 0.26 \ast$	70.81 ± 0.82
Upper Limit	96.73 ± 0.00	93.56 ± 0.00	86.48 ± 0.00	91.11 ± 0.00	92.01 ± 0.00	95.38 ± 0.00	81.66 ± 0.00

Table 1: Results for the sockpuppet prediction task. Asterisks indicate statistical significance compared to the standard encoder.

was used to evaluate the statistical significance of the meta-encoder compared to the standard encoder, indicated with asterisks. The averaged ROC and PR curves are provided in Appendix A.

525

526

530

531

533

536

537

539

541

543

545

546

547

549

550

551

553

555

557

558

559

560

561

562

565

The meta-learning approach significantly outperforms other approaches (P << 0.05) in AU-ROC, AUPRC, F1-score, F0.5-score and accuracy, with substantial improvements of approximately 10%. Recall did not improve significantly, tying the overall improvement to an increase in precision. This suggests that meta-learning helps the classifier make fewer false positive predictions whilst preserving its ability to identify true positives. This is desirable for sockpuppet detection, where a high confidence in positive predictions is paramount.

The metrics of the pre-trained encoder are unexpected. The additional pre-training should have provided the model with a general understanding of the task prior to fine-tuning, however the approach falls behind the non-pre-trained encoder model on most metrics. In prior work (Sakib and Spezzano, 2022; Solorio et al., 2013a), the pretrained approach performed well. This reduction in performance may be due to the harder task setting, however there may be other causes. Sakib and Spezzano (2022) combined their authorship attribution features with behavioural features, which may be more consistent across tasks, and therefore better for the pre-trained approach. This suggests approaches that focus on authorship attribution may require a model to have greater adaptability.

The lowest performing metric is the AUPRC result. This is likely due to the class imbalance, which overall consisted of 65.60% negative and 34.40% positive samples. The AUPRC is more sensitive to 'hard' negative samples, as precision decreases substantially for each false positive.

This result suggests that despite an increase in the precision of predictions being the principal benefit of meta-learning, it still remains the model's main flaw. This conclusion is further corroborated in Section 7.1. Intriguingly, both non-metalearning approaches achieved marginally worse scores than the random baseline. The similarly low precision scores corroborate the earlier statement that the principal improvement of meta-learning in this domain is the reduction in false positives.

Surprisingly, the performance difference between the basic RoBERTa classifier and the encoder model is small. The encoder model was expected to perform better as it is trained on the word level embeddings produced by RoBERTa, and therefore should have had a richer understanding of user writing style than the semantic sentence level embeddings used in the RoBERTa classifier. In all metrics the encoder performs slightly better, suggesting there is some truth to the hypothesis, however, the small training set sizes may have prevented a significant divergence.

Whilst prior Wikipedia sockpuppet detection approaches report higher F1-scores of 73 (Solorio et al., 2013a) and 82 (Sakib and Spezzano, 2022), the difference in datasets and task construction (neither study distinguishes between sockpuppets and puppetmasters) make a fair comparison difficult.

7.1 Error Analysis

Figure 3 provides insight into the effect of metalearning on the embeddings. The embeddings of two test investigations¹⁰ have been projected to two dimensions using Principal Component Analysis (PCA), with positive samples being coloured in blue. The left-hand column contains the embeddings of the test samples produced by the standard encoder model after training on the task. The righthand column are the embeddings produced with the meta-encoder.

The embeddings learnt by the meta-encoder appear tighter, with less overlap between the clusters. This is supported by the results of these particular investigations: Investigation A received an AU-ROC of 62% with the standard encoder, which

602

603

604

566

¹⁰Investigations of *Film_Fan* and *Al_aman_kollam*, respectively.







Figure 4: PCA of Low Performing Embeddings.

improved to 83% using the meta-encoder. Investigation B had a similar improvement, from 69% to 91%. Both standard and meta-encoders were able to cluster positive samples together, however the meta-encoder exhibits better separation from negative samples. This aligns with the overall results, where the meta-encoder saw small improvements in recall, but large improvements in precision.

607

611

613

615

616

617

619

621

623

To contrast the successful examples, Figure 4 presents two investigations (*Amirharbo* and *Cameronfree*) that performed poorly. Investigation A achieved AUROCs of 52% (standard encoder) and 54% (meta-encoder). Investigation B was similar, with a small improvement from 58% to 61%. Whilst less defined, the right-angle structure is still evident, and positive samples are still clustered within a single arm, suggesting the encoder has no issues identifying positive samples. The difference then is the proportion of negative samples that ap-

pear in the 'positive' arm. This again aligns with the overall results, where recall is largely consistent between approaches, and most of the improvement is in the precision of positive classifications. In these two cases, the poor AUROC performance can be attributed to the failure of the meta-encoder to improve upon the precision. 624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

In Investigation B, contribution messages are characteristically short, typically the name of the article section edited. The messages of many negative samples are similar. This convention is easily detected, explaining why both classifiers were able to cluster the positive samples, but found distinguishing them from negative samples using the convention difficult. This may explain why the recall is acceptable, but precision is low. The sockpuppetgroup in investigation A at most use a message of just a couple of words, but typically use no message at all. As most Wikipedians add a contribution message (only 8.89% of negative samples collected had empty message fields, compared to 15.56% of positive samples), a consistently empty one would identify the sockpuppet to the encoder, but would be indistinguishable from legitimate message-less contributions. This leads to the following conclusion: where a sockpuppet's behaviour is characterised by empty or conventional messages, positive recall is strong, but precision suffers.

8 Conclusion

We study the problem of detecting malicious sockpuppetry on Wikipedia. We are the first to propose meta-learning to address the data-scarcity challenge in detecting sockpuppet accounts through writing style. Our results demonstrate significant performance improvements when compared to pretrained approaches, especially in prediction precision. We attribute this to our approach's ability to quickly adapt to distinct authorship styles with limited samples. In doing so, we defined a more realistic task definition that provides a more accurate measure of performance, and released an updated, verifiable and adaptable dataset of sockpuppet investigations appropriate for future meta-learning research. Our findings extend to any online social platform where users engage in sockpuppetry.

Limitations

There are several limitations of our model that could benefit from further research.

As discussed in Section 7.1, our model is limited

in cases where sockpuppet contributions contain lit-673 tle or no message data. In these cases, the encoder 674 requires additional information. One approach to 675 do so would be to include the edit data of contributions, that is the changes made to the article itself. This would allow a model to understand the intent and implications of a contribution even 679 when a description is absent. A contribution must make changes to the article, and edits themselves are likely to be far more diverse in nature than the messages, providing a model with a strong distinguishing signal. The additional signal would also further improve high performing investigations.

Another limitation is in safety. There are several legitimate reasons why a user might have several accounts. One of these reasons may be for the safety of editors editing politically contentious articles. Whilst our approach was trained solely on malicious examples of sockpuppetry, no efforts were made to ensure this approach could not reveal benign sockpuppet-groups by mistake. Additional work may focus on providing a safeguard measure that ensures the sockpuppet behaviour being observed is malicious.

695

702

706

707

Our approach should also be evaluated against Generative language models, which are becoming increasingly effective at creating text that looks human (Liu et al., 2023). It is likely that future sockpuppet-groups might utilise generative models to edit Wikipedia, rewriting edits to hide the author's writing style, or automating edits entirely. Many approaches are already focusing on the detection of text generated by prolific models (Dhaini et al., 2023). Future work may evaluate how robust the meta-learning approach is to authorship obfuscation using generative models.

Considering the performance of the approach, it is unable to replace human-led sockpuppet investi-710 gations. When found to be guilty of sockpuppetry, accounts are blocked. Some users assign great 712 value to their accounts, and incorrect sockpuppet 713 classifications would be damaging to the commu-714 nity. The approach could serve as an additional 715 716 source of evidence in open investigations, or as a detection method that triggers human-led investi-717 gation on suspicious accounts. To occupy a larger 718 role in investigations, the precision of the approach must improve further. 720

Ethical Considerations

There is a valid concern for privacy when releasing this dataset. Usernames are important to Wikipedia editors, and may be used to represent a person's real identity, contain some personally identifiable information, or obscure their identity completely.

Arguments against anonymisation are numerous. Whilst this study does not use username data, previous approaches have (Sakib and Spezzano, 2022), and future approaches may too. Wikipedia moderators are currently debating what the best practice should be¹¹, however, as the data is public, any anonymisation attempts would ultimately be circumventable. For these reasons, the data was not anonymised.

References

- Sadia Afroz, Michael Brennan, and Rachel Greenstadt. 2012. Detecting hoaxes, frauds, and deception in writing style online. In 2012 IEEE Symposium on Security and Privacy, pages 461–475.
- Bo Ai, Yuchen Wang, Yugin Tan, and Samson Tan. 2022. Whodunit? learning to contrast for authorship attribution. In *Proceedings of the 2nd Conference* of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing.
- Antreas Antoniou, Harrison Edwards, and Amos Storkey. 2019. How to train your maml. In *Seventh International Conference on Learning Representations*.
- Trapit Bansal, Karthick Gunasekaran, Tong Wang, Tsendsuren Munkhdalai, and Andrew McCallum. 2021. Diverse distributions of self-supervised tasks for meta-learning in nlp. In *EMNLP*, pages 5812– 5824.
- Michael Brennan, Sadia Afroz, and Rachel Greenstadt. 2012. Adversarial stylometry: Circumventing authorship recognition to preserve privacy and anonymity. *ACM Trans. Inf. Syst. Secur.*, pages 12:1–12:22.
- Zhan Bu, Zhengyou Xia, and Jiandong Wang. 2013. A sock puppet detection algorithm on virtual spaces. *Knowledge-Based Systems*, pages 366–377.
- Loic De Langhe, Orphee De Clercq, and Veronique Hoste. 2024. Unsupervised authorship attribution for medieval Latin using transformer-based embeddings. In Proceedings of the Third Workshop on Language Technologies for Historical and Ancient Languages (LT4HALA) @ LREC-COLING-2024, pages 57–64.

730 731 732

733

734

735

721

722

723

724

725

726

727

728

729

736

737 738

739

740

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

760

763

764

765

766

767

768

¹¹https://meta.wikimedia.org/wiki/Research: Wikimedia_Research_Best_Practices_Around_ Privacy_Whitepaper/Draft#2.4_Wikipedia_(user) names

- 769 776 777 782 783 786 789 791 797 798 799 800 805 812 813 814 815 816 817 818 819

822 823

- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4171–4186.
- Mahdi Dhaini, Wessel Poelman, and Ege Erdogan. 2023. Detecting ChatGPT: A survey of the state of detecting ChatGPT-generated text. In Proceedings of the 8th Student Research Workshop associated with the International Conference Recent Advances in Natural Language Processing, pages 1–12.
- Maciej Eder. 2013. Does size matter? authorship attribution, small samples, big problem. Digital Scholarship in the Humanities, pages 167–182.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In Proceedings of the 34th International Conference on Machine Learning - Volume 70, page 1126-1135.
- Chelsea Finn, Kelvin Xu, and Sergey Levine. 2018. Probabilistic model-agnostic meta-learning. In NeurIPS.
- Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. 2020. Meta-learning in neural networks: A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence, pages 5149-5169.
- Javier Huertas-Tato, Alvaro Huertas-Garcia, Alejandro Martin, and David Camacho. 2022. Part: Pre-trained authorship representation transformer.
- Sarah Khaled, Neamat El-Tazi, and Hoda M. O. Mokhtar. 2018. Detecting fake accounts on social media. In 2018 IEEE International Conference on Big Data.
- Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. International Conference on Learning Representations.
- Sotiris B Kotsiantis, Ioannis Zaharakis, P Pintelas, et al. 2007. Supervised machine learning: A review of classification techniques. Emerging artificial intelligence applications in computer engineering, pages 3-24.
- Srijan Kumar, Justin Cheng, Jure Leskovec, and V.S. Subrahmanian. 2017. An army of me: Sockpuppets in online discussion communities. In Proceedings of the 26th International Conference on World Wide Web, page 857-866.
- Bin Liang, Xiang Li, Lin Gui, Yonghao Fu, Yulan He, Min Yang, and Ruifeng Xu. 2023. Few-shot aspect category sentiment analysis via meta-learning. ACM Trans. Inf. Syst., pages 22:1–22:31.
- Dong Liu, Quanyuan Wu, Weihong Han, and Bin Zhou. 2016. Sockpuppet gang detection on social media sites. Frontiers of Computer Science, pages 124-135.

Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, Zihao Wu, Lin Zhao, Dajiang Zhu, Xiang Li, Ning Qiang, Dingang Shen, Tianming Liu, and Bao Ge. 2023. Summary of chatgpt-related research and perspective towards the future of large language models. Meta-Radiology.

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach.
- Alex Nichol, Joshua Achiam, and John Schulman. 2018. On first-order meta-learning algorithms.
- Manoj Niverthi, Gaurav Verma, and Srijan Kumar. 2022. Characterizing, detecting, and predicting online ban evasion. In Proceedings of the ACM Web Conference 2022, page 2614–2623.
- Simon Owens. 2013. The battle to destroy wikipedia's biggest sockpuppet army.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In EMNLP, pages 1532–1543.
- Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. 2019. Meta-learning with implicit gradients. In NeurIPS.
- Rafael A. Rivera-Soto, Olivia Elizabeth Miano, Juanita Ordonez, Barry Y. Chen, Aleem Khan, Marcus Bishop, and Nicholas Andrews. 2021. Learning universal authorship representations. In *EMNLP*.
- Diego Saez-Trumper. 2019. Online disinformation and the role of wikipedia.
- Mostofa Najmus Sakib and Francesca Spezzano. 2022. Automated detection of sockpuppet accounts in wikipedia. In 2022 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), pages 155–158.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 815-823.
- Prasha Shrestha, Sebastian Sierra, Fabio González, Manuel Montes, Paolo Rosso, and Thamar Solorio. 2017. Convolutional neural networks for authorship attribution of short texts. In EACL, pages 669-674.
- Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. SIGKDD Explor. Newsl.
- Jake Snell, Kevin Swersky, and Richard S. Zemel. 2017. Prototypical networks for few-shot learning. In NeurIPS, page 4080-4090.

Chaehan So. 2021. Exploring meta learning: Parameterizing the learning-to-learn process for image classification. In 2021 International Conference on Artificial Intelligence in Information and Communication (ICAIIC), pages 199–202.

877

878

881

882

894

896

897

898

899

900

901

902

904

909

910

911

912

913 914

915

916

917

918

919

921

- Thamar Solorio, Ragib Hasan, and Mainul Mizan. 2013a. A case study of sockpuppet detection in wikipedia. In *Proceedings of the Workshop on Language Analysis in Social Media*.
- Thamar Solorio, Ragib Hasan, and Mainul Mizan. 2013b. Sockpuppet detection in wikipedia: A corpus of real-world deceptive writing for linking identities. In *International Conference on Language Resources and Evaluation*.
- Brad Stone and Matt Richtel. 2007. The hand that controls the sock puppet could get slapped. *The New York Times*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In *NAACL-HLT*.
- Lin Tian, Xiuzhen Zhang, and Jey Han Lau. 2023. Metatroll: Few-shot detection of state-sponsored trolls with transformer adapters. In *Proceedings of the ACM Web Conference 2023*.
- Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle. 2020. Meta-dataset: A dataset of datasets for learning to learn from few examples.
- Michail Tsikerdekis and Sherali Zeadally. 2014. Multiple account identity deception detection in social media using nonverbal behavior. *IEEE Transactions on Information Forensics and Security*, pages 1311– 1321.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. 2016. Matching networks for one shot learning. In *NeurIPS*.
- Xueling Zheng, Yiu Ming Lai, K.P. Chow, Lucas C.K. Hui, and S.M. Yiu. 2011. Sockpuppet detection in online discussion forums. In 2011 Seventh International Conference on Intelligent Information Hiding and Multimedia Signal Processing.

A Appendix

Algorithm 1 Reptile Algorithm (Serial)

Input: Learning rate α , number of inner steps k, task distribution $p(\mathcal{T})$ Initialise θ , the vector of initial model parameters

for iteration = 1, 2, ... do Sample a task $\mathcal{T}_i \in p(\mathcal{T})$ Compute $\theta' = U_{\mathcal{T}_i}^k(\theta)$, denoting k steps of SGD or Adam Update $\theta \leftarrow \theta' + \epsilon(\theta' - \theta)$ end for



Figure 5: Average ROC curves of models on test tasks.



Meta Precision Recall Curves

Figure 6: Average PR curves of models on test tasks.

Model	Hyper-parameter	Value
	Number of Attention Heads	2
	Number of Layers	6
Encoder Model	Learning Rate	0.0001
	Loss Margin	0.2
	Optimiser	Adam
	Dropout Chance	0.35
	Learning Rate	0.001
Classification Model	Layer 0 Nodes	768
	Layer 1 Nodes	768
	Optimiser	Adam
	Dropout Chance	0.7615
	Learning Rate	0.0008
	Layer 0 Nodes	768
	Layer 1 Nodes	512
RoBERTa Classifier	Layer 2 Nodes	512
	Layer 3 Nodes	256
	Layer 4 Nodes	256
	Layer 5 Nodes	128
	Optimiser	Adam
Dontilo	Interpolation Rate	0.2
керше	Number of Steps	5

Table 2: Tuned model hyper-parameters.