Enhancing Neural Topic Model with Multi-Level Supervisions from Seed Words

Anonymous ACL submission

Abstract

Efforts have been made to apply topic seed words to improve the topic interpretability of topic models. However, due to the semantic diversity of natural language, supervisions from seed words could be ambiguous, making it hard to be incorporated into the current neural topic models. In this paper, we propose SeededNTM, a neural topic model enhanced with supervisions from seed words on both word and document levels. We introduce a context-dependency assumption to alleviate the 011 ambiguities with context document informa-012 tion, and an auto-adaptation mechanism to au-014 tomatically balance between multi-level infor-015 mation. Moreover, an intra-sample consistency regularizer is proposed to deal with noisy supervisions via encouraging perturbation and semantic consistency. Extensive experiments 019 on multiple datasets show that SeededNTM can derive semantically meaningful topics and outperforms the state-of-the-art seeded topic models in terms of topic quality and classification accuracy.

1 Introduction

024

Unsupervised topic models, despite their efficiency in uncovering the underlying latent topics in text 026 corpora (Blei et al., 2003), may suffer from poor topic interpretability as the semantic interpretability of latent space is poorly explored (Chang et al., 2009; Newman et al., 2011; Eshima et al., 2020) and the generated topics may not match users' desires (Jagarlamudi et al., 2012; Gallagher et al., 2017; Harandizadeh et al., 2022). To address this problem, topic seed words are incorporated as additional prior knowledge to provide richer semantic information and indicate users' preferences. Compared to sample-wise information like document 037 labels, seed words can be easier to access, more widely applicable, and with a milder level of human bias.

Many works in conventional topic models incorporate seed words as guidance. Some works extend Latent Dirichlet Allocation (LDA) into seeded models (Andrzejewski and Zhu, 2009; Jagarlamudi et al., 2012; Li et al., 2016; Eshima et al., 2020), and some draw inspiration from information theory (Gallagher et al., 2017) or word embeddings (Meng et al., 2020). While most of the conventional topic models struggle with the growing number of topics and documents, with the recent development of neural topic models (NTM), keyETM (Harandizadeh et al., 2022) is proposed to incorporate seed words into NTM to combine the advantages of NTM of scalability on large datasets. 041

042

043

044

045

047

049

051

055

057

060

061

062

063

064

065

066

067

069

071

072

073

074

075

076

077

078

079

However, keyETM only focuses on regularizing word-topic relations with seed words and fails to combine document-level topic information, which is essential as the semantics of words may vary under different context documents. As shown in Figure 1(a), under different contexts, the word 'apple' has different semantic meanings and may belong to different topics, even if it co-occurs with the seed word 'company'. This inspires us to incorporate supervisions from seed words into NTM on both word and document level and balance information from both levels for better inference of topics, thus achieving better topic interpretability.

There still remain challenges to effectively combining multi-level supervisions from seed words into the current framework of NTM. Firstly, the **mean-field assumption** made in current NTMs prevents the model from combining topic preferences of words and documents because they are assumed to be conditionally independent. Secondly, as shown in Figure 1(b), document level supervisions from seed words can be noisy due to the semantic ambiguity of natural languages. Previous work (Li et al., 2018) tried to tackle the problem via a neighbor consistency regularization. However, the neighbor-based method can be time-consuming, limiting the scalability on large datasets, and noisy



Figure 1: Examples from UIUC Yahoo Answers dataset. (a) Multiple semantic meanings of the word 'apple' under different contexts. (b) Seed words from three different topics bring noises to each other when estimating document topic preferences.

neighbors may cause cumulative errors.

To address these challenges, we propose a novel neural topic model SeededNTM, which incorporates seed words as supervisions and autoadaptively balances information from both word and document level. During variational inference, we drop the mean-field assumption and make a context-dependency assumption to assist the inference of per-word topic assignment with context document information. Based on this assumption, we implement an auto-adaptation mechanism between multi-level information inspired by the idea of product of experts (Hinton, 2002). Moreover, to deal with the noisy document supervisions, we propose a novel regularizer that encourages intra-sample consistency to avoid time-consuming neighbor finding and cumulative errors. The regularizer encourages consistency between perturbed samples to preserve local structures and consistency between the semantics of outputs from different encoders to improve robustness.

Our contributions are summarized as follows:

- We propose SeededNTM, a novel neural topic model that leverages supervisions from seed words on both word and document level.
- We propose a reasonable context-dependency assumption and develop an auto-adaptation mechanism to automatically balance between word level and document level information.
- We propose an intra-sample consistency regularizer to deal with noises from document level supervisions by encouraging both perturbation and semantic consistency,.
- Extensive experiments on three public datasets show that SeededNTM can derive

semantically meaningful topics and outperforms the state-of-the-art seeded topic models in terms of NPMI and classification accuracy.

2 Related Works

2.1 Neural Topic Model

The recent developments of neural variational inference (Kingma and Welling, 2014; Rezende et al., 2014) enable the application of neural networks on topic models to deal with scalability issues. NVDM (Miao et al., 2016) and ProdLDA (Srivastava and Sutton, 2017) are two representative works. Gaussian and logistic normal distribution are leveraged as approximations of the Dirichlet prior in the original LDA. Subsequently, various works have been proposed (Nan et al., 2019; Dieng et al., 2020; Nguyen and Luu, 2021), aiming for better inference of topics.

Among these works, the most relevant to our work is VRTM (Rezaee and Ferraro, 2020). It explicitly models each word's the topic assignments z_n while other works collapse them for simplicity. However, the mean-field assumption in VRTM prevents the model from combining context document information when inferring words' topic preferences, limiting its performance.

2.2 Topic Model with Prior Knowledge

Introducing prior knowledge into topic models has been a widely adopted way to improve topic interpretability. Sample-wise knowledge, like labels (Blei and Mcauliffe, 2008; Wang and Yang, 2020) and covariates (Eisenstein et al., 2011; Card et al., 2018) are popular choices but can be difficult to acquire and may introduce strong biases. In

2

108

109

110

111

112

113

114

115

116

120 121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

145

146

147

148

149

117

118

contrast, topic seed words, as a kind of topic-wise 150 knowledge, can be easier to access and more appli-151 cable. z-label LDA (Andrzejewski and Zhu, 2009) 152 proposed to use "z-labels" to bias the word-topic 153 distributions in Gibbs sampling. SeededLDA (Ja-154 garlamudi et al., 2012) paired each topic with a 155 seed topic and biased documents to topics if they 156 have corresponding seed words. And keyATM (Es-157 hima et al., 2020) improved upon SeededLDA by 158 allowing topics with no seed word and better empir-159 ical hyperparameters. Anchored CorEx (Gallagher et al., 2017) proposed an information-theoretic 161 framework and incorporates seed words by anchor-162 ing them to topics. CatE (Meng et al., 2020) took 163 category names as seed words and learned a dis-164 criminative embedding space for topics and words.

> Recently, to combine the advantages of NTMs on scalability, keyETM (Harandizadeh et al., 2022) is proposed to incorporate seed words into NTM by regularizing word-topic relations with seed words and pre-trained word embeddings.

168

169

170

171

172

173

174

175

176

177

178

179

180

181

183

187

191

192

193

194

2.3 Dataless Text Classification with Topic Models

Dataless text classification is a branch of classification task which requires building a text classifier with a few relevant words or descriptions for each category and no sample-wise labels. On account of the similar settings with seeded topic modeling, a few topic model-based methods are proposed (Chen et al., 2015; Li et al., 2016, 2018). Despite similar settings, dataless text classification and seeded topic modeling differ in many aspects. While seeded topic modeling aims at discovering latent topics and focuses on the interpretability of learned topics, dataless text classification aims to classify text to pre-defined classes and focuses on the validity of the document-category partitions. Unsupervised topics are allowed in seeded topic modeling, and documents are interpreted as mixtures of multiple topics, while in dataless text classification, every category is assumed to be known in advance, and a document may be assumed to belong to a single category.

3 Background

3.1 Problem Formulation

195 Consider a corpus with D documents, where 196 each document d contains N_d words $w_d =$ 197 $\{w_{d1}, w_{d2}, \ldots, w_{dN_d}\}$, each belonging to a vocab-198 ulary of size V. And suppose that we have K topics, each provided with a set of L_k seed words denoted by $S_k = \{s_{k1}, s_{k2}, \dots, s_{kL_k}\}$. Our goal is to derive topics from the corpus that are semantically coherent with corresponding seed word sets.

199

200

201

202

203

204

205

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

224

225

226

227

228

229

231

232

233

234

235

3.2 Generative Story and Variational Inference

Our model builds on the generative story in (Srivastava and Sutton, 2017), where the Dirichlet prior is approximated via a logistic normal distribution. The generative story is summarized as follows, where α is the parameter for prior distribution and β_k denotes the word distribution for the k-th topic:

For document d, draw topic distribution $\theta \sim$
$\mathcal{LN}(\mu_0(\alpha), \sigma_0^2(\alpha));$
For w_{dn} in this document:
Draw topic $z_{dn} \sim Cat(\theta)$;
Draw word $w_{dn} \sim Cat(\beta_{z_{dn}});$

Based on the generative story, variational inference is used to approximate posterior distribution of latent variables θ_d and $z_d = \{z_{d1}, z_{d2}, \ldots, z_{dN_d}\}$ to maximize the likelihood of observed data. And the evidence lower bound (ELBO) can be derived as

$$(\boldsymbol{w}) = E_{q(\theta, \boldsymbol{z} | \boldsymbol{w})} \log \left(p(\boldsymbol{w} | \theta, \boldsymbol{z}; \beta) \right) - E_{q(\theta, \boldsymbol{z} | \boldsymbol{w})} \log \left(\frac{q(\theta, \boldsymbol{z} | \boldsymbol{w})}{p(\theta, \boldsymbol{z})} \right)$$
(1) 22
$$= - \left(\mathcal{L}_{rec} + \mathcal{L}_{kl} \right),$$

where $q(\theta, \boldsymbol{z} | \boldsymbol{w})$ is the joint variational distribution.

4 Methodology

 \mathcal{L}

In this section, we introduce our proposed *Seed-edNTM*. We start by introducing the model architecture and the designs of multi-level pseudo supervisions. Then we focus on our proposed auto-adaptation mechanism based on contextdependency assumption and our noise-reduction consistency regularizer. Finally, we introduce our training objective and summarize the training procedure with Algorithm 1.

4.1 Model Architecture

4.1.1 Document Encoder

A multi-layer network is used as document encoder to infer the document-topic distributions θ 237 for document d with a word set w. The words are first encoded into word embedding vectors E_d = 239 $\{e_1, e_2, \ldots, e_{N_d}\}$ and then averaged to obtain the document embedding e_d . Then the mean vector μ 241



Figure 2: The overall structure of SeededNTM. The grey boxes indicate the training losses in SeededNTM, and the dashed boxes indicate the variables used in loss computations.

and the diagonal of the covariance matrix σ^2 are further encoded with two sub-networks $\mu = f_{\mu}(e_d)$ and $\sigma^2 = f_{\sigma}(e_d)$, and the document-topic distribution is sampled via the reparameterization trick with $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and $\theta = softmax(\mu + \sigma \cdot \epsilon)$. The above procedure is donoted as $\theta = F_d(d)$.

4.1.2 Word Encoder

243

244

245

247

248

251

254

261

263

265

267

Word encoder encodes words to local word-topic preferences ϕ . For a word w_n , it is first encoded to the embedding vector e_n , followed by a feedforward network activated with a softmax function. The above procedure is donoted as $\phi_n = F_w(w_n)$.

4.1.3 Topic Decoder

The decoder contains topic-word distribution and reconstructs documents with topic mixtures. Inspired by (Eisenstein et al., 2011), we disassemble topics in log-space into three parts, background m, regular topic η^r , and seed topic η^s . The background term is estimated with the overall log frequencies of words from the corpus, and both regular and seed topics act as additional deviations on m. The possibility β_{kv} for word w_v in topic k is

$$\beta_{kv} = \frac{\exp(m_v + \eta_{kv}^r + \eta_{kv}^s)}{\sum_v \exp(m_v + \eta_{kv}^r + \eta_{kv}^s)},$$
 (2)

where η_k^r is a V-dimensional parameter vector whose elements at positions corresponding to S_k are fixed to zero. And η_k^s is defined as

$$\eta_{kv}^{s} = \begin{cases} \kappa, & w_{v} \in S_{k}, \\ 0, & \text{otherwise,} \end{cases} v \in \{1, \cdots, V\}, \quad (3)$$

where κ is a hyperparameter of seeding strength.

4.2 Multi-Level Supervisions

270

271

272

273

274

275

276

277

278

279

281

283

287

289

290

4.2.1 Document Level Supervision

With seed words, we can regularize the inferred document-topic distribution $\hat{\theta}$ with the pseudo distribution $\hat{\theta}$ which is estimated via the *tf-idf* scores of seed words appearing in the document. Formally, for a document *d*, its corresponding $\hat{\theta}$ is

$$\hat{\theta}_k = \frac{\frac{1}{L_k} \sum_{s \in S_k} tfidf(s, d)}{\sum_k \left(\frac{1}{L_k} \sum_{s \in S_k} tfidf(s, d)\right)}, k \in \{1, \dots, K\}.$$
(4)

And we regularize θ by minimizing the KL divergence between θ and $\hat{\theta}$,

$$\mathcal{L}_d(\theta, \hat{\theta}) = KL(\hat{\theta} \| \theta) = \sum_k \hat{\theta}_k \log(\frac{\hat{\theta}_k}{\theta_k}).$$
 (5)

4.2.2 Word Level Supervision

Local word-topic preferences ϕ can also be regularized by seed words. We estimate the pseudo word-topic distribution $\hat{\phi}$ with co-occurrence measured by the conditional possibility p(w|s) =df(w,s)/df(s) of word w and seed word s, where $df(\cdot)$ is the number of documents containing s or both s and w. And the pseudo possibility for word w_n belonging to topic k is

$$\hat{\phi}_{nk} = \frac{\frac{\tau}{L_k} \sum_{s \in S_k} p(w_n | s)}{\sum_k \left(\frac{\tau}{L_k} \sum_{s \in S_k} p(w_n | s)\right)}, \quad (6)$$

where τ is a temperature factor to sharpen the distribution. And we also use KL divergence to mini293

297

301

306

307

309

310

313

314

315

317

318

321

322

323

mize the distance between $\hat{\phi}_n$ and ϕ_n ,

$$\mathcal{L}_w(\phi_n, \hat{\phi}_n) = KL(\hat{\phi}_n \| \phi_n) = \sum_k \hat{\phi}_{nk} \log(\frac{\hat{\phi}_{nk}}{\phi_{nk}}).$$
(7)

4.3 Auto-Adaptation of Multi-Level Information

In previous work (Rezaee and Ferraro, 2020), the inferred posterior distribution $q(\theta, z|w)$ is decomposed with a mean-field assumption as

$$q(\theta, \boldsymbol{z} | \boldsymbol{w}) = q(\theta | \boldsymbol{w}) \prod_{n} q(z_n | w_n), \qquad (8)$$

but as we mentioned before, per-word topic preferences can be ambiguous without context document information. Therefore, instead of mean-field assumption, we introduce a context-dependency assumption by taking document topic distribution θ into consideration,

$$q(\theta, \boldsymbol{z} | \boldsymbol{w}) = q(\theta | \boldsymbol{w}) \prod_{n} q(z_n | w_n, \theta).$$
(9)

As z_n is now conditioned on both w_n and θ , how to properly balance information from word and document remains unsolved. Inspired by the idea of *product of experts* (Hinton, 2002), we propose an auto-adaptation mechanism to automatically combine local word-topic preference ϕ_n and the global document-topic preference θ and implement the combination as products of two distributions,

$$\varphi_{nk} = q(z_n = k | \theta, w_n) = \frac{\phi_{nk} \theta_k}{\sum_k (\phi_{nk} \theta_k)}.$$
 (10)

In this way, we avoid manually weighting the global and local topic preferences and achieve autoadaptation between multi-level information. Potential ambiguities in per-word topic preferences get re-weighted by the global document-topic distributions, and topics with higher probabilities in both distributions are further encouraged.

4.4 Noise-Reduction Consistency Regularizer

Document level supervisions can be biased by seed words' semantic diversity and ambiguity of. To avoid time-consuming nearest neighbor method (Li et al., 2018), inspired by recent works in noisy label learning (Li et al., 2020; Englesson and Azizpour, 2021), we propose a consistency regularizer that encourages intra-sample consistency. In this regularizer, we encourage outputs from the document encoder to be consistent with perturbed samples, $d' \sim \mathcal{A}(d)$, where \mathcal{A} is an data augmentation function. Each perturbed sample can be viewed as a neighbor with the original sample in feature space, and by encouraging perturbation consistency, we can preserve local structures without finding nearest neighbors.

332

333

334

335

337

338

339

341

342

343

345

346

347

349

350

351

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

371

Moreover, we encourage consistency with the outputs from the word encoder. The word encoder takes supervisions from the word-word cooccurrences and contains more fine-grained information than the document level. By encouraging consistency with the predictions of the word encoder on document embeddings, we incorporate semantic information from the word level to help correct the predictions from the document encoder and improve its robustness to noises.

We use the symmetric KL Divergence to measure the distance between two distributions, and our consistency regularizer is summarized as follows.

$$SKL(a, b) = KL(a||b) + KL(b||a),$$

$$\mathcal{L}_{c}(d) = SKL(\theta, F_{d}(d')) + SKL(\theta, F_{w}(d)).$$
(11) 353

4.5 Training Objectives

With the new assumption in Eq.9, \mathcal{L}_{rec} and \mathcal{L}_{kl} in Eq.1 can be further derived as

$$\mathcal{L}_{rec} = -\sum_{n,k} \varphi_{nk} \log \beta_{kw_n},$$

$$\mathcal{L}_{kl} = KL \left(\mathcal{N}(\mu, \sigma^2) \| \mathcal{N}(\mu_0, \sigma_0^2) \right) + \sum_n KL \left(\varphi_n \| \theta \right).$$

(12)

Detailed derivations can be found in Appendix A. Our final training objectives is

$$\mathcal{L}_{tr} = \mathcal{L}_{rec} + \lambda_0 \mathcal{L}_{kl} + \lambda_1 \mathcal{L}_d + \lambda_2 \mathcal{L}_w + \lambda_3 \mathcal{L}_c,$$
(13)

where λ_0 is KL annealing factor and gradually increases to 1 during training and $\lambda_1, \lambda_2, \lambda_3$ are hyperparameters. The overall structure of Seeded-NTM is shown in Figure 2, and the training procedure is described in Algorithm 1.

5 Experiments

5.1 Datasets

We conduct our experiments on three datasets: 20 Newsgroups, UIUC Yahoo Answers, and DB-Pedia. 20 Newsgroups (Lang, 1995) is a dataset that contains around 20,000 newsgroup documents

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

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

Algorithm 1 The SeededNTM training procedure.

Input: corpus \mathcal{D} , topic number K, seed word sets $S = \{S_1, S_1, \ldots, S_K\}$, initial KL annealing factor λ_0 , hyperparameters $\lambda_1, \lambda_2, \lambda_3$, max iteration number T.

for t from 1 to T do randomly sample a batch of B documents; $\mathcal{L}_{batch} \leftarrow 0$; $\lambda_0 \leftarrow \min(\lambda_0 + \frac{1}{T}, 1.0)$; compute β_k for each topic k by Eq.3; for each document d in the batch do compute θ with encoder F_d ; compute ϕ_n for each w_n with encoder F_w ; compute $\varphi_d = \{\varphi_1, ..., \varphi_n\}$ by Eq.10; $\mathcal{L}_{batch} \leftarrow \mathcal{L}_{batch} + \mathcal{L}_{tr}$ by Eq.13 end for update model parameters with $\nabla \mathcal{L}_{batch}$ end for

and is commonly used in the topic modeling field. And to verify our model's scalability, we adopt two other larger datasets, the UIUC Yahoo Answers dataset (Chang et al., 2008) and DBPedia (Zhang et al., 2015), which contain 150,000 and 630,000 samples, respectively. We preprocess each dataset and split them for training and testing. The detailed procedure of preprocessing and the statistical summaries for each dataset can be viewed in Appendix B.

5.2 Seed Words Extraction

372

374

375

376

387

388

390

397

To avoid human biases, we follow (Jagarlamudi et al., 2012; Gallagher et al., 2017) and adopt an automatic approach to extract seed words. For each dataset, we set the topic number K the same as its class number, and use Information Gain (IG) to identify the words having the highest mutual information with the class. Specifically, IG of a word w in class c is

$$IG(w,c) = H(c) - H(c|w),$$
 (14)

where H(c) is the entropy of class c and H(c|w)denotes the conditional entropy of c given w. For each class, we choose the top L words with the highest IG scores as seed words.

5.3 Evaluation of Topic Quality

5.3.1 Evaluation Metrics

We use Topic Coherence, i.e., Normalized Pointwise Mutual Information (NPMI), to evaluate the quality of learned topics. NPMI between words w_i and w_j is defined as:

$$NPMI(w_i, w_j) = \frac{\log \frac{p(w_i, w_j)}{p(w_i)p(w_j)}}{-\log p(w_i, w_j)}.$$
 (15)

As we are dealing with topic models with seed words, we take the top N non-seed words and predefined L seed words for each topic and measure NPMI among the N + L words. For unsupervised methods, we pick the top N + L words. By considering both seed and non-seed words, the NPMI scores can measure how well the learned topics fit the predefined aspects of interests. Also, the score implicitly reflects topic diversity, as topics with a high coherence score with seed words are more likely to be diverse as long as their seed words are distinct. We report NPMI with N = 10, L = 5 on both train and test sets. Results with different seed word numbers can be viewed in Appendix C.

5.3.2 Baselines

We compare SeededNTM with the following baselines. For unsupervised topic models, we compare with LDA (Blei et al., 2003) and prodLDA (Srivastava and Sutton, 2017), which are representative in conventional and neural topic models, and for seed-guided topic models, we compare with z-labels LDA (Andrzejewski and Zhu, 2009), SeededLDA (Jagarlamudi et al., 2012), STM (Li et al., 2016), Anchored Corex (Gallagher et al., 2017), CatE (Meng et al., 2020), keyATM (Eshima et al., 2020) and keyETM (Harandizadeh et al., 2022), which we have introduced in related works.

5.3.3 Performances

The performances on topic qualities are reported in Table 1. As we can see, most seeded topic models achieve better topic coherence than unsupervised ones as the seed words provide additional semantic information. SeededNTM outperforms the baselines in most settings, demonstrating the effectiveness of our approach. Note that the advantages become more significant on the largest datasets, DBPedia, indicating its scalability when facing datasets of huge scale. We can find that keyETM sometimes performs worse performances than conventional methods like STM and keyATM, indicating the necessity to incorporate document level information. Anchor Corex and CatE are strong baselines on some occasions, as Anchor Corex has an information-theory-based objective

	20 Newsgroups			Yahoo Answer				DBPedia				
Methods	NP	PMI	F	1	NF	MI	F	1	NF	PMI	F	'1
	train	test	Macro	Micro	train	test	Macro	Micro	train	test	Macro	Micro
LDA	0.288	0.262	-	-	0.186	0.160	-	-	0.074	-0.027	-	-
ProdLDA	0.289	0.223	-	-	0.225	0.134	-	-	0.116	0.043	-	-
z-labels LDA	0.250	0.223	0.344	0.356	0.149	0.134	0.374	0.394	0.238	0.236	0.791	0.801
Seeded LDA	0.273	0.244	0.346	0.329	0.215	0.208	0.581	0.558	0.266	0.262	0.835	0.837
STM	0.346	0.306	0.485	0.516	0.290	0.280	0.606	0.617	0.309	0.295	0.898	0.899
Anchor Corex	0.360	0.313	0.387	0.357	0.295	0.282	0.502	0.497	0.312	0.295	0.776	0.771
CatE	0.358	0.332	0.238	0.242	0.321	0.239	0.214	0.209	0.178	0.069	0.522	0.521
keyATM	0.294	0.267	0.298	0.293	0.177	0.174	0.610	0.592	0.274	0.269	0.854	0.856
keyETM	0.359	0.329	0.310	0.333	0.242	0.233	0.439	0.425	0.259	0.254	0.754	0.776
SeededNTM	0.368	0.338	0.570	0.576	0.334	0.286	0.629	0.627	0.331	0.311	0.902	0.903

Table 1: The NPMI and F1 scores on three datasets. Results are reported through a single run with a randomly chosen seed word.

	NP	MI	F1		
Methods	train	test	Macro	Micro	
SeededNTM	0.368	0.338	0.570	0.576	
SeededNTM-noise	0.359	0.328	0.559	0.564	
SeededNTM-NN	0.359	0.329	0.566	0.572	
SeededNTM-doc	0.362	0.329	0.567	0.570	
SeededNTM-word	0.358	0.316	0.563	0.568	
SeededNTM-mean	0.279	0.216	0.414	0.525	

Table 2: Results of different variants of SeededNTM on20 Newsgroups.

similar to NPMI, and CatE takes the order words as additional information when learning embeddings.

5.4 Evaluation of Text Classification

5.4.1 Evaluation Metrics

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

Text classification is a prevalent task to test topic models' ability to extract semantic information from documents. Here we adopt the setting of dataless text classification and take the maximum probability in the document topic distribution as the predicted label. We use Macro and Micro F1 scores as the evaluation metrics. As most baselines cannot predict on new data, we report the results on the train set and take the test set for validation.

5.4.2 Baselines

We compare SeededNTM on classification with the aforementioned baselines except for the unsupervised ones. Specifically, we follow CatE's original paper and use a dataless classification method, WeSTClass (Meng et al., 2018), to classify its outputs.

5.4.3 Performances

Table 1 summarizes the F1 scores on three datasets. SeededNTM outperforms other baseline models on most occasions, indicating our model can understand the semantics of the documents and learn more reliable and helpful topic distributions for each document. Among the baselines methods, seededNTM, STM, and keyATM achieve better performances on three datasets, as they incorporate information from seed words on both levels. 470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

5.5 Ablation Studies

We analyze the effects of different modules of SeededNTM by comparing among the following variants: 1) SeededNTM-noise: SeededNTM without the consistency regularizer, 2) SeededNTM-NN: SeededNTM without the consistency regularizer and with a neighbor-based noise-reduction method as in (Li et al., 2018). 3) SeededNTM-doc: Seeded-NTM with supervisions only from document level, 4) SeededNTM-word: SeededNTM with supervisions only from word level, 5) SeededNTM-mean: SeededNTM with the mean-field assumption as in (Rezaee and Ferraro, 2020).

Performances are provided in Table 2, from which we can draw the following conclusions. The effectiveness of the noise-reduction method can be proved by the comparisons between variants with and without noise regularizer. Both SeededNTM-NN and original SeededNTM outperform SeededNTM-noise. And the effectiveness of our intra-sample consistency regularizer can be further demonstrated by the improvements of SeededNTM over SeededNTM-NN. The decreases in SeededNTM-doc and SeededNTM-word indicate the importance of supervisions on both levels. Moreover, the significant decay on SeededNTMmean proves the effectiveness of our proposed assumption and the necessity to balance context doc-

	Topic 1: Game&Recreation	Topic2: Arts	Topic3: Pregnancy&Parenting
Seed words	pokemon, game, diamond, games, trade	book, harry, potter, books, poem	pregnancy, baby, weeks, child, pregnant
z-labels LDA	play, think, best, ps, great	product, black, color, white, read	just, time, day, days, period
Seeded LDA	play, ps, wii, level, code	read, know, names, love, movie	just period time days day
STM	ps, wii, level, code, xbox	read, story, write, series, movie	period, doctor, sex, months, normal
Anchor Corex	play, pearl, playing, fc, ps	read, write, reading, writing, author	months, period, days, week, birth
CatE	gba, ds, nintendo, replay, mew	rowling, hallows, novel, author, deathly	trimester, babies, conception, expecting, womb
KeyATM	play, ps, just, need, wii	read, know, just, good, think	just, know, time, period, day
KeyETM	know, think, good, really, want	question, answer, read, come, called	year, years, old, months,feel
SeededNTM	fc, wii, nintendo, ds, pearl	hallows, deathly, author, rowling, novel	ovulation, period, ttc, ovulating, pill

Table 3: Top five words of part of the topics and corresponding seed words learned by different models on UIUC Yahoo Answers dataset.

Topics	keyATM	KeyETM	SeededNTM
Business&Finance	need, want, work, time, business	phone, card, business, download, video	loan, bank, tax, payment, income
Health	just, know, day, time, good	water, hair, product, cup, add	pregnancy, pregnant, pill, ovulation, period
Education	school, college, know, just, work	god, book, books, world, classes	colleges, classes, degree, gpa, schools
Pets	dog, just, dogs, know, cat	old, wear, house, clean, big	puppy, kitten, puppies, breed, litter
Computer&Internet	just, need, want, download, know	-	wireless, router, vista, phones, cable
New Topic	-	time, long, way, probably, usually	craigslist, ebay, google, shops, sites

Table 4: Top five words learned on UIUC Yahoo Answers dataset while only 3 topics are with seed words.

ument information when modeling per-word topicassignments.

5.6 Qualitative Evaluation

508

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

Besides quantitative evaluations, we hope to demonstrate our model's ability to discover semantically meaningful topics under conditions closer to real-world situations in a more intuitive manner.

5.6.1 Topic Presentation

We first compare part of topics learned by Seeded-NTM on UIUC Yahoo Answer dataset with topics learned by baselines methods using the same seed words in the aforementioned experiments in Table 3. We can find that some baselines, such as z-labels LDA, Anchor Corex, and KeyETM, tend to put high weights on several commonly used words like 'play', 'great', 'good', while SeededNTM tends to pay attention to words that are more specific such as 'nintendo', a Japanese multinational video game company who releases the game 'Pokemon', and 'rowling', the author of Harry Potter, and 'ttc', meaning 'trying to conceive'.

5.6.2 Topic with Incomplete Seed Words

In the above experiments, seed words are assumed to be complete and accurately represent latent topics in the corpus. However, in practical situations, users may only be interested in part of the corpus or have little prior knowledge, leading to incomplete seed words. To simulate such situations, we preserve seed words for only three topics and leave other topics unsupervised. We present the results of SeededNTM along with the two latest baselines, keyATM and keyETM in Table 4.

For three supervised topics, SeededNTM can discover words related to the seed words as it does under complete seed words, while KeyATM and keyETM produce semantically incoherent topics, such as irrelevant words "god" and "world" appearing in the topic 'Education&Reference' from keyETM. SeededNTM can also discover meaningful unsupervised topics similar to the seeded topics in former experiments, such as 'Pets' and 'Computer&Internet', while keyATM and keyETM find incoherent or unrelated topics. Moreover, new topics which are not included in the original seed word sets can also be discovered by SeededNTM, such as 'Craigslist', a famous American classified advertisements website. 536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

6 Conclusions

In this paper, we propose *SeededNTM* to improve topic interpretability together with scalability. We leverage supervisions from seed words on both word and document levels and propose a contextdependency assumption. An auto-adaptation mechanism is designed to balance word and context document information. Moreover, we propose an intra-sample consistency regularizer to deal with noisy document level supervisions. Perturbation consistency and semantic consistency are encouraged to improve the model's robustness to noises. Through quantitative and qualitative experiments on three datasets, we demonstrate that SeededNTM can derive semantically meaningful topics and outperforms state-of-the-art baselines.

References

568

569

570

571

572

573

574

580

581

583

586

587

588

591

593

594

595

611

612

613

614

615

616

617

618

619

- David Andrzejewski and Xiaojin Zhu. 2009. Latent dirichlet allocation with topic-in-set knowledge. In *Proceedings of the NAACL HLT 2009 Workshop on Semi-Supervised Learning for Natural Language Processing*, pages 43–48.
- David Blei and Jon Mcauliffe. 2008. Supervised topic models. In *Advances in Neural Information Processing Systems*, volume 20.
- David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022.
- Dallas Card, Chenhao Tan, and Noah A Smith. 2018. Neural models for documents with metadata. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 2031–2040.
- Jonathan Chang, Sean Gerrish, Chong Wang, Jordan Boyd-Graber, and David Blei. 2009. Reading tea leaves: How humans interpret topic models. *Advances in neural information processing systems*, 22.
- Ming-Wei Chang, Lev-Arie Ratinov, Dan Roth, and Vivek Srikumar. 2008. Importance of semantic representation: Dataless classification. In *AAAI*, volume 2, pages 830–835.
- Xingyuan Chen, Yunqing Xia, Peng Jin, and John Carroll. 2015. Dataless text classification with descriptive lda. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29.
- Adji B Dieng, Francisco JR Ruiz, and David M Blei. 2020. Topic modeling in embedding spaces. *Transactions of the Association for Computational Linguistics*, 8:439–453.
- Jacob Eisenstein, Amr Ahmed, and Eric P Xing. 2011. Sparse additive generative models of text. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 1041–1048. Citeseer.
- Erik Englesson and Hossein Azizpour. 2021. Generalized jensen-shannon divergence loss for learning with noisy labels. *Advances in Neural Information Processing Systems*, 34.
- Shusei Eshima, Kosuke Imai, and Tomoya Sasaki. 2020. Keyword assisted topic models. *arXiv preprint arXiv:2004.05964*.
- Ryan J Gallagher, Kyle Reing, David Kale, and Greg Ver Steeg. 2017. Anchored correlation explanation: Topic modeling with minimal domain knowledge. *Transactions of the Association for Computational Linguistics*, 5:529–542.
- Bahareh Harandizadeh, J. Hunter Priniski, and Fred Morstatter. 2022. Keyword assisted embedded topic model. In Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining, WSDM '22, page 372–380.

Geoffrey E Hinton. 2002. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8):1771–1800. 623

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

673

674

675

676

- Jagadeesh Jagarlamudi, Hal Daumé III, and Raghavendra Udupa. 2012. Incorporating lexical priors into topic models. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 204–213.
- Diederik P Kingma and Max Welling. 2014. Autoencoding variational bayes. In 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings.
- Ken Lang. 1995. Newsweeder: Learning to filter netnews. In *Machine Learning Proceedings 1995*, pages 331–339. Elsevier.
- Chenliang Li, Jian Xing, Aixin Sun, and Zongyang Ma. 2016. Effective document labeling with very few seed words: A topic model approach. In *Proceedings of the 25th ACM international on conference on information and knowledge management*, pages 85–94.
- Junnan Li, Richard Socher, and Steven C.H. Hoi. 2020. Dividemix: Learning with noisy labels as semisupervised learning. In *International Conference on Learning Representations*.
- Ximing Li, Changchun Li, Jinjin Chi, Jihong Ouyang, and Chenliang Li. 2018. Dataless text classification: A topic modeling approach with document manifold. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 973–982.
- Yu Meng, Jiaxin Huang, Guangyuan Wang, Zihan Wang, Chao Zhang, Yu Zhang, and Jiawei Han. 2020. Discriminative topic mining via category-name guided text embedding. In *Proceedings of The Web Conference 2020*, pages 2121–2132.
- Yu Meng, Jiaming Shen, Chao Zhang, and Jiawei Han. 2018. Weakly-supervised neural text classification. In Proceedings of the 27th ACM International Conference on Information and Knowledge Management, pages 983–992.
- Yishu Miao, Lei Yu, and Phil Blunsom. 2016. Neural variational inference for text processing. In *International conference on machine learning*, pages 1727–1736. PMLR.
- Feng Nan, Ran Ding, Ramesh Nallapati, and Bing Xiang. 2019. Topic modeling with wasserstein autoencoders. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6345–6381.
- David Newman, Edwin V Bonilla, and Wray Buntine. 2011. Improving topic coherence with regularized topic models. *Advances in neural information processing systems*, 24.

Thong Nguyen and Anh Tuan Luu. 2021. Contrastive learning for neural topic model. *Advances in Neural Information Processing Systems*, 34:11974–11986.

678

679

681 682

683

684

686

690

694

695

696

697

698 699

700

701

703

704

- Mehdi Rezaee and Francis Ferraro. 2020. A discrete variational recurrent topic model without the reparametrization trick. In *Advances in Neural Information Processing Systems*, volume 33, pages 13831– 13843. Curran Associates, Inc.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. 2014. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR.
- Akash Srivastava and Charles Sutton. 2017. Autoencoding variational inference for topic models. In 5th International Conference on Learning Representations.
 - Xinyi Wang and Yi Yang. 2020. Neural topic model with attention for supervised learning. In *International Conference on Artificial Intelligence and Statistics*, pages 1147–1156. PMLR.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. 2020. Unsupervised data augmentation for consistency training. *Advances in Neural Information Processing Systems*, 33:6256–6268.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28:649–657.

A Derivation of ELBO-based Loss

The Evidence Lower Bound (ELBO) for our model is

$$ELBO(\boldsymbol{w}) = E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log p(\boldsymbol{w}|\theta, \boldsymbol{z}; \beta) - E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log \left(\frac{q(\theta, \boldsymbol{z}|\boldsymbol{w})}{p(\theta, \boldsymbol{z})}\right).$$
(A.1) (A.1)

To maxmize the ELBO, we minimize its opposite number as training loss, which is

$$\mathcal{L}_{elbo} = -E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log p(\boldsymbol{w}|\theta, \boldsymbol{z}; \beta) + E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log \left(\frac{q(\theta, \boldsymbol{z}|\boldsymbol{w})}{p(\theta, \boldsymbol{z})}\right).$$
(A.2)

And we denote

$$\mathcal{L}_{rec} = -E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log p(\boldsymbol{w}|\theta, \boldsymbol{z}; \beta),$$

$$(q(\theta, \boldsymbol{z}|\boldsymbol{w}))$$
(712)

$$\mathcal{L}_{kl} = E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log \left(\frac{q(\theta, \boldsymbol{z}|\boldsymbol{w})}{p(\theta, \boldsymbol{z})} \right), \tag{A.3}$$

$$\mathcal{L}_{elbo} = \mathcal{L}_{rec} + \mathcal{L}_{kl}.$$

For the posterior $q(\theta, \boldsymbol{z} | \boldsymbol{w})$, we have

$$q(\theta, \boldsymbol{z} | \boldsymbol{w}) = q(\theta | \boldsymbol{w}) \prod_{n} q(z_n | \theta, w_n).$$
(A.4) 715

For $p(\boldsymbol{w}|\boldsymbol{\theta}, \boldsymbol{z}; \boldsymbol{\beta})$, we have

$$p(\boldsymbol{w}|\boldsymbol{\theta}, \boldsymbol{z}; \boldsymbol{\beta}) = \prod_{n} p(w_n | \boldsymbol{z}_n; \boldsymbol{\beta}).$$
(A.5) 717

So for \mathcal{L}_{rec} we have

$$\mathcal{L}_{rec} = -E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log p(\boldsymbol{w}|\theta, \boldsymbol{z}; \beta)$$

$$= -E_{q(\theta|\boldsymbol{w})} E_{q(z_1|\theta, w_1)} \dots E_{q(z_N|\theta, w_N)} \log p(\boldsymbol{w}|\theta, \boldsymbol{z}; \beta)$$

$$= -E_{q(\theta|\boldsymbol{w})} \sum_{n} E_{q(z_n|\theta, w_n)} \log p(w_n|z_n; \beta).$$

(A.6) 719

The expectation $E_{q(\theta|w)}$ can be estimated using a sample-based method by sampling $\theta \sim q(\theta|w)$, and given θ , $\varphi_{nk} = q(z_n = k|\theta, w_n)$ can be computed with Eq.10. So we have 721

$$\mathcal{L}_{rec} \approx -\sum_{n,k} \varphi_{nk} \log \beta_{kw_n}. \tag{A.7}$$

For \mathcal{L}_{kl} we have

$$\mathcal{L}_{kl} = E_{q(\theta, \boldsymbol{z}|\boldsymbol{w})} \log\left(\frac{q(\theta, \boldsymbol{z}|\boldsymbol{w})}{p(\theta, \boldsymbol{z})}\right)$$

= $E_{q(\theta|\boldsymbol{w})} \log\left(\frac{q(\theta|\boldsymbol{w})}{p(\theta)}\right) + E_{q(\theta|\boldsymbol{w})} \sum_{n} E_{q(z_n|\theta, w_n)} \log\left(\frac{q(z_n|\theta, w_n)}{p(z_n|\theta)}\right)$ (A.8)
= $KL\left(q(\theta|\boldsymbol{w}) \| p(\theta)\right) + E_{q(\theta|\boldsymbol{w})} \sum_{n} KL\left(q(z_n|\theta, w_n) \| p(z_n|\theta)\right).$

The former term can be approximated using Laplace approximation to the Dirichlet prior, and can be calculated in closed form as $KL\left(\mathcal{N}(\mu, \sigma^2) \| \mathcal{N}(\mu_0, \sigma_0^2)\right)$ (Srivastava and Sutton, 2017). And the latter term can be estimated by Monte Carlo sampling with $\theta \sim q(\theta | \boldsymbol{w})$: 727

$$E_{q(\theta|\boldsymbol{w})} \sum_{n} KL\left(q(z_n|\theta, w_n) \| p(z_n|\theta)\right) \approx \sum_{n} KL(\varphi_n\|\theta).$$
(A.9) 728

707 708

710

714

716

718

B More Details of Datasets

730 B.1 Dataset Descriptions

729

731

733

734

735

739

741

742

743

744

745

746

747

748

751

Three datasets are used in out experiments: **20** Newsgroups, UIUC Yahoo Answers, and DBPedia. 20 Newsgroups (Lang, 1995) is a collection of newsgroup documents containing 11,000 train samples and 7,000 test samples in 20 classes. It is a common dataset that is widely used in topic modeling field. The UIUC Yahoo Answers dataset (Chang et al., 2008) contains 150,000 question-answer pairs belonging to 15 categories. It is a classification dataset and is used in topic models in (Card et al., 2018). DBPedia (Zhang et al., 2015) is extracted from Wikipedia and contains 560,000 train samples and 70,000 test samples belonging to 14 ontology classes. DBPedia is a classification dataset, and to the best of our knowledge, it is the first time that DBPedia has been used for topic modeling, but similar datasets (though much smaller) from Wikipedia have been adopted to test topic models (Nguyen and Luu, 2021).

740 B.2 Preprocess Procedures for Datasets

We preprocess documents in each dataset by tokenizing, filtering out stop words, words with document frequency above 70%, and words appearing in less than around 100 documents (depending on the dataset). The final vocabulary sizes for each dataset after preprocessing vary from 2,000 to 20,000. Then we remove the documents shorter than two words.

Specifically, for the UIUC Yahoo Answer dataset, we follow the approach used in (Card et al., 2018), and drop the *Cars and Transportation* and *Social Science* classes and merge *Arts* and *Arts and Humanities* into one class, producing 15 categories, each with 10,000 documents.

As for the augmentation functions A, we use the word level augmentation method proposed in (Xie et al., 2020) by randomly replacing words with lower tf-idf scores. Around of 10% words are replaced in our experiments.

B.3 Statistics of Datasets

We summarize the statistics for the three datasets after preporcessing in Table.B.1

	20 Newsgroups	Yahoo Answer	DBPedia
Class Number	20	15	14
Vocabulary Size	2,004	7,468	19,975
Train Set Size	10,732	119,747	559,710
Test Set Size	7,105	29,937	69,962
Avg Doc Length	44.308	46.089	22.730
Token Number	790,324	6,898,796	13,682,938

Table B.1: Summary of the statistics of three datasets

753 C More Experimental Details

C.1 Implementation Datails

As for the training environment, we implement our method based on **PyTorch** 1.6.0 with Python 3.7.9 and perform our experiments on 4 GeForce RTX 2080Ti. For model structure, the dimension for our word embedding layer is 300, and the dimension for the hidden layer in the document encoder is 256. We use a 0.2 dropout rate in our encoder during training. We present our choices for hyperparameters in Table.C.1. Hyperparameters are determined by grid search on the smallest dataset, 20 Newsgroups, and fine-tuned on other two large datasets. The final hyperparameters are shown in Table C.1.

761 C.2 Baselines

We give detailed descriptions of our baselines here.

	LR	batch size	λ_1	λ_2	λ_3	au	κ
20 Newsgroups	0.001	64	2.0	10.0	5.0	4.0	3.0
Yahoo Answer	0.001	128	2.0	10.0	5.0	4.0	3.0
DBPedia	0.0005	256	2.0	10.0	1.0	4.0	3.0

Table C.1: The choices of hyperparameters for each dataset.

LDA (Blei et al., 2003): LDA is one of the most popular unsupervised conventional topic models that deduce posterior distribution via Gibbs sampling or variational inference.
 prodLDA (Srivastava and Sutton, 2017): prodLDA is one of the most representative neural topic 765

766

767

768

769

771

772

773

774

775

776

778

779

780

781

782

783

784

785

786

- models. It uses black-box neural variational inference and optimizes the model with stochastic gradient descent, increasing the model's scalability. prodLDA is unsupervised and cannot incorporate seed words.
- **z-labels LDA** (Andrzejewski and Zhu, 2009): z-labels LDA utilizes seed word information by biasing the seed words' choices for topics in Gibbs sampling.
- SeededLDA (Jagarlamudi et al., 2012): SeededLDA pairs each regular topic with a topic containing only seed words and biases documents' topic preferences in Gibbs sampling if they contain seed words.
- **STM** (Li et al., 2016): STM is a topic model-based dataless text classification method that incorporates both document and word level supervisions to improve classification accuracies.
- Anchored Corex (Gallagher et al., 2017): Anchored CorEx is based on an information-theoretic framework and tries to derive maximally informative topics based on seed words.
- **CatE** (Meng et al., 2020): CatE aims at deriving topics with a single seed word for each topic. It uses a word embedding method and tries to learn a discriminative embedding space for both topics and words.
- **keyATM** (Eshima et al., 2020): keyATM improves upon SeededLDA by allowing some seed-word-free topics.
- **keyETM** (Harandizadeh et al., 2022): keyETM incorporates seed words into NTM by regularizing word-topic and topic-word distributions on word level with seed words and pre-trained word embeddings.

	20 Newsgroups			Yahoo Answer				DBPedia				
Methods	NF	PMI	F	1	NF	PMI	F	'1	NI	PMI	F	'1
	train	test	Macro	Micro	train	test	Macro	Micro	train	test	Macro	Micro
LDA	0.292	0.266	-	-	0.195	0.186	-	-	0.083	-0.002	-	-
ProdLDA	0.297	0.236	-	-	0.242	0.153	-	-	0.121	0.054	-	-
z-labels LDA	0.228	0.208	0.272	0.288	0.156	0.145	0.365	0.385	0.270	0.266	0.747	0.765
Seeded LDA	0.302	0.285	0.335	0.341	0.203	0.195	0.583	0.561	0.275	0.265	0.821	0.824
STM	0.358	0.334	0.484	0.507	0.294	0.283	0.592	0.604	0.313	0.302	0.888	0.890
Anchor Corex	0.343	0.314	0.396	0.384	0.309	0.300	0.458	0.450	0.315	0.299	0.746	0.739
CatE	0.360	0.341	0.233	0.227	0.365	0.278	0.233	0.224	0.153	0.035	0.581	0.575
keyATM	0.302	0.269	0.307	0.306	0.174	0.169	0.602	0.584	0.278	0.270	0.830	0.833
keyETM	0.363	0.322	0.323	0.328	0.228	0.222	0.370	0.384	0.260	0.240	0.596	0.624
SeededNTM	0.381	0.331	0.562	0.570	0.367	0.320	0.609	0.606	0.352	0.343	0.896	0.896

C.3 More Quantitative Results

Table C.2: The NPMI and F1 scores on three datasets when N=10,L=3

787 C.4 More Qualitative Results

Due to the space limit, we present here some more qualitative results under settings different from the main paper.

C.4.1 Noisy Seed Words

790

791

792

793

795

796

797

802 803

805

816

The seed word set may contain irrelevant words in real-world practice due to users' mistakes or unfamiliarity with the corpus. To simulate such situations, we manually intrude irrelevant words from other topics into the seed words. The results are shown in Table C.3, from which SeededNTM can still find meaningful topics when there are noisy intrusions in the seed words, while keyATM and keyETM provide topics that are less explicit and coherent.

Topics	noisy word	keyATM	KeyETM	SeededNTM
Society&Culture	company	people, just, think, life, believe	life, believe, world, man, word	christian, religious, beliefs, faith, christianity
Sports	phones	think, good, year, game, best	game, pokemon, play, points, level	baseball, league, win, fans, nfl
Beauty&Style	cat	product, look, color, just, want,	product, cute, black, color, clothes	jpg, shoes, hollister, shirt, curly

Table C.3: The top five words of topics learned on UIUC Yahoo Answers dataset with noisy seed words.

C.4.2 Transferred Seed Words

One way to explore an unfamiliar dataset is to start with topics from another known corpus. In this experiment, we transfer the topical seed words from 20 Newsgroups and DBPedia and use them for training SeededNTM on UIUC Yahoo Answers dataset. Topics learned with the transferred seed words are presented in Table C.4, along with the topics learned in the original topics. We can find that though these datasets are collected from entirely different sources, some semantically meaningful topics can still be discovered with transferred seed words, and some lead to slightly different concepts from the originals. Moreover, the results indicates that topic-wise supervisions are flexible and bear less bias than sample-wise supervisions.

Seed Words	20News	Yahoo
god, atheists, religion	belief, religions, existence	belief, religious, christians
graphics, format, image	files, ftp, screen	picture, jpg, albums
space, launch, orbit	moon, solar,flight	paint, walls, room
	DBPedia	Yahoo
football, league, played	player, professional, team	qb, wr, rb
high, school, students	schools, secondary, grades	degree, college, university
species, family, flowering	endemic, native, habitat	plant, soil, flowers

Table C.4: The top words of topics learned with transferred seed words from 20 Newsgroups and DBPedia.

C.4.3 Exploration on the various aspects of single concept

Due to the ambiguity of natural language, a single word or concept may relate to various topics with different meanings, especially for some common words such as 'apple', 'doctor' or 'card'. In this case, we assume that the users aim at using topic models to understand different topics in the corpus related to a single word. We start with a single word, 'card'. We set only one topic with a single seed word 'card' and leave other topics unsupervised. Then we use the topic model to generate one supervised topic about 'card' and several unsupervised topics. Iteratively, we treat the most related word in the topic 'card' as the seed word for a new topic and train another topic model under new settings. The results are shown in Table C.5. Due to space limitations, we only list the topic 'card' in round 4 and round 5. From the results, SeededNTM shows its ability to distinguish different semantic topics related to the same word, which can be used to assist users with understanding complex concepts.

D Limitations and Potential Risks of SeededNTM

Though SeededNTM achieves good performances in our experiments, there are still some limitations. Firstly, supervisions from seed words, though flexible, are also very weak and vulnerable to noises.

Round	seed words	SeededNTM
1	card	phone, phones, cell, cards, sim, mobile
2	card	itunes, ipods, vista, router, dvd, xp
	phone	phones, cell, verizon, mobile, cingular, motorola
3	card	credit, money, pay, loan, bank, cards
	phone	phones, know, cell, cards, mobile, verizon
	itunes	ipod, download, windows, songs, music, files
4	card	camera, cards, digital, memory, laptop, graphics
5	card	wii, grphics, cards, memory, dell, ram

Table C.5: The top five words of topics learned on UIUC Yahoo Answer dataset with iteratively-given seed words.

Though we introduce some ways to improve the model's robustness, it is still possible that the model may
crash under intentional attacks. Secondly, seed words in our model are used as pseudo supervisions. A
more elegant way is to incorporate it into the generative story. As for potential risks, seeded topic models
can be used to trace a specific topic, so it is possible that it's used to track someone's information from
texts collected from the internet, violating personal privacy.819820
821
822