

Appendix for Putting Back the Stops: Integrating Syntax with Neural Topic Models

Author Name
Affiliation
email@example.com

- Appendix A: Proofs and Training
- Appendix B: Context Network
- Appendix C: Decision Module
- Appendix D: Datasets and Preprocessing
- Appendix E: Modeling Details
- Appendix F: Ablation studies for SyConNTM
- Appendix G: Syntactic and Semantic Topics
- Appendix H: Motivation from Cognitive Science and Linguistics
- Appendix I: Limitations
- Appendix J: Things we tried that did not work

A Proofs and Training

In this section, we present the proof for Theorem 3.1, ultimately leading to the objective of the C-VAE model.

For clarity, the plate notation for C-VAE is given in Figure 1. In the figure, x_i represents part of the data from feature class i (e.g. x_1 might be the semantic words, and x_2 the syntactic words). Each x_i is a BoW vector and $\sum_i x_i = x$ is one document. z_i are the corresponding latent vectors. C represents the context network which decides how each document is divided into n parts x_1, \dots, x_n . Assuming the data features are divided into n parts, where each part is considered independent, the log-likelihood of the data can be given as:

$$\log p(x) = \log \prod_{i=1}^n p(x_i)$$

Using the Bayes' rule to rewrite the log-likelihood in terms of the latent variables:

$$\log p(x) = \log \int \prod_{i=1}^n p(x_i|z_i)p(z_i) dz_1 \dots dz_n$$

Introducing the approximate posterior distributions for z_1, z_2, \dots, z_n :

$$\log p(x) = \log \int \prod_{i=1}^n q(z_i|x)p(x_i|z_i)p(z_i) \cdot \prod_{i=1}^n \frac{1}{q(z_i|x)} dz_1 \dots dz_n$$

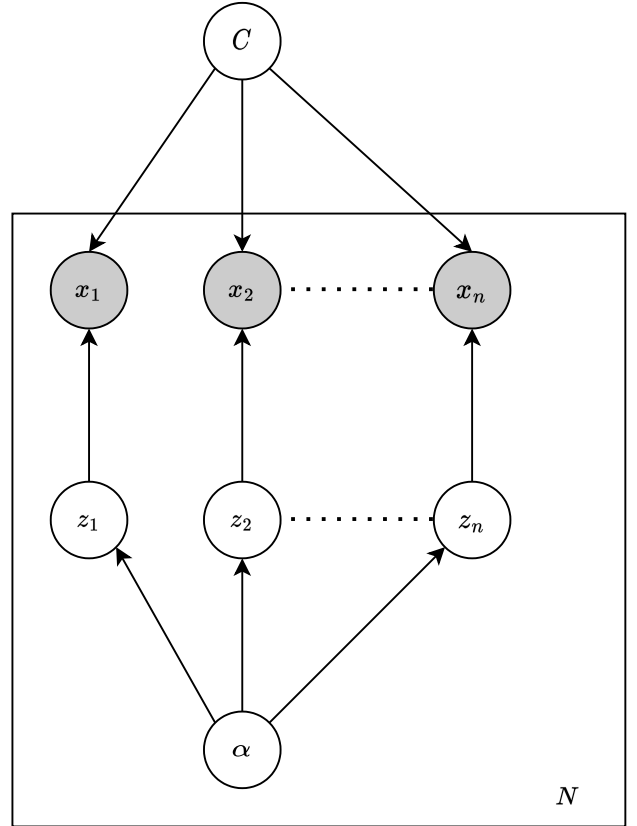


Figure 1: C-VAE in plate notation

Using the property of the log function to split the integral:

$$\log p(x) = \sum_{i=1}^n \mathbb{E}_{q(z_i|x)} [\log p(x_i|z_i)] + \sum_{i=1}^n \log \int q(z_i|x)p(z_i) dz_1 \dots dz_n$$

Substituting the KL-divergence for the last term:

$$ELBO = \sum_{i=1}^n \mathbb{E}_{q(z_i|x)} [\log p(x_i|z_i)] - \sum_{i=1}^n D_{KL}(q(z_i|x)||p(z_i))$$

28

29

Algorithm 1 Training SyConNTM

```

1: Input: documents ( $D$ ), context type ( $C_{\text{type}}$ ), and size ( $c$ ),
   decision type  $\text{dt}$ , and decision hparam  $t/M$ 
2: Initialize hyper-parameters and model parameters.
3: for batch  $\mathcal{B}$  in  $D$  do
4:   Get  $x_{\text{BoW}}$  and  $x_{\text{seq}}$  from  $\mathcal{B}$ 
5:    $s = \text{context\_net}(x_{\text{seq}}, C_{\text{type}}, c, E)$ 
6:   Update parameters (Equation 2)
7:    $x_{\text{BoW}}^{\text{syn}}, x_{\text{BoW}}^{\text{sem}} = \text{dec\_module}(s, \text{dt}, t/M, x_{\text{seq}}, x_{\text{BoW}})$ 
8:    $x_{\text{BoW}}^{\text{syn}}, x_{\text{BoW}}^{\text{sem}}, z_{\text{syn}}, z_{\text{sem}} = \text{SC-NTM}(x_{\text{BoW}})$ 
9:   Update SC-NTM parameters (Equation 4)
10: end for

```

30 Re-writing *ELBO* in terms of Equation 1:

$$\mathcal{L}_{\text{c-vae}}(\theta, \phi; x) = \sum_{i=1}^n \mathcal{L}(\theta, \phi_i; x)$$

A.1 Training:

The algorithm for training SyConNTM is shown in Algorithm 1. The input to the algorithm includes a set of documents (D), a context type (C_{type}) and context size c for the context network, and a decision type dt for the decision module. For each batch \mathcal{B} in D , a BoW representation, x_{BoW} , and a sequential representation x_{seq} , is constructed. These representations are then passed through the *context_net* (Algorithm 3, Appendix B) which gives probabilities s for each word in x_{seq} . The parameters of the context network are updated based on Equation 2. The prediction probabilities s are then passed through the *dec_module* (Algorithm 4, Appendix C) which generates two bag-of-word vectors, $x_{\text{BoW}}^{\text{syn}}$ and $x_{\text{BoW}}^{\text{sem}}$ with syntax and semantic words respectively. Finally, x_{BoW} is passed through the SC-NTM model, which updates its parameters based on Equation 4. We also investigated pretraining the context network which resulted in a faster convergence of SyConNTM at the cost of the pretraining overhead, but did not improve results.

A.2 Generative process of the network:

We describe the generative process of a document in Algorithm 2 using a combination of semantic and syntactic topics. Firstly, for a given document, we sample the semantic topic distribution, denoted as $\theta_d^{\text{semantic}}$, from a Dirichlet distribution with parameter α . This distribution captures the distribution of semantic topics within the document. Similarly, we sample the syntactic topic distribution, $\theta_d^{\text{syntactic}}$, using the same procedure. Next, for each word in the document, we draw a class, c_i , from the context_net, represented by π_{context} . If the class is determined to be of the semantic type ($c_i = 1$), we draw a topic, z_i , from the semantic topic distribution. Conversely, if the class is syntactic, we draw a topic, z_i , from the syntactic topic distribution. Finally, we draw the word, w_i , from the word distribution ϕ_{z_i} , which is parameterized by the VAE decoder. This generative process allows us to generate documents and capture both semantic and syntactic characteristics.

Algorithm 2 Generative process of the network

```

1: Sample semantic topic distribution for document  $d$ ,
    $\theta_d^{\text{semantic}} \sim \text{Dirichlet}(\alpha)$ 
2: Sample syntactic topic distribution for document  $d$ ,
    $\theta_d^{\text{syntactic}} \sim \text{Dirichlet}(\alpha)$ 
3: for each word  $w_i$  in  $D$  do
4:   Draw class  $c_i \sim \pi_{\text{context}}$  (here  $\pi$  comes from the context_net)
5:   if  $c_i = 1$  (semantic class) then
6:     Draw topic  $z_i \sim \theta_d^{\text{semantic}}$ 
7:   else
8:     Draw topic  $z_i \sim \theta_d^{\text{syntactic}}$ 
9:   end if
10:  Draw  $w_i \sim \phi_{z_i}$  (word distributions in  $\phi$  are parameterized by the VAE decoder)
11: end for

```

B Context Network

The context network is similar to the feed forward n-gram model by [Chelba *et al.*, 2017] and takes as input a concatenation of embeddings of c context words, i.e. $X = \text{concat}(E(w_{i-(c-1)}), \dots, E(w_{i-1}))$.

Our context network is defined as follows:

$$X_1 = \text{dropout}(X, p_{\text{keep}}) \quad (1)$$

$$X_2 = \tanh(L^1 \cdot X_1 + L_{\text{bias}}^1) \quad (2)$$

$$X_3 = \text{BatchNorm}(\text{dropout}(X_2, p_{\text{keep}})) \quad (3)$$

$$\text{out} = \text{softmax}(L^2 \cdot X_3 + L_{\text{bias}}^2) \quad (4)$$

The above equations represent the context network in a series of operations that are performed on the input data, denoted by the variable X . The first operation in Equation 1 applies the dropout regularization technique [Srivastava *et al.*, 2014] to the input data, where p_{keep} is the probability of retaining a given input element during the dropout operation. The second operation in Equation 2 applies the hyperbolic tangent activation function to the dot product of the first layer weights, L^1 , and the dropout-regularized input data, X_1 , plus the first layer bias term, L_{bias}^1 . The third operation in Equation 3 applies batch normalization [Ioffe and Szegedy, 2015] and dropout regularization to the output of the previous operation, X_2 . Finally, the fourth operation in Equation 4 applies the softmax activation function to the dot product of the second layer weights, L^2 , and the batch-normalized, dropout-regularized input, X_3 , plus the second layer bias term, L_{bias}^2 . The output of the final operation, *out*, represents the model's prediction. The hyperparameters for the model are: p_{keep} and the dimensionality of the weight vectors L^1 and L^2 .

The context network predicts a target word w_i given a context size c . We experiment with two types of context: symmetric ($w_{i-(c-1)}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+(c+1)}$) and asymmetric ($w_{i-(c-1)}, \dots, w_{i-1}$) for a target word w_i .

We define the Algorithm 3 for generating target probability r and predicted probability s from an input sequence x_{seq} of length n , a context type C_{type} , a context size C_{size} , and an embedding look-up table E . s is initialized with zero to store the prediction probabilities. The algorithm iterates over each

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

Algorithm 3 context_net

```

1: Input:  $x_{\text{seq}}$  of length  $n$ , context type  $C_{\text{type}}$ , context size  $c$ ,
   embedding look up table  $E$ 
2: Return: prediction probability  $s$ 
3: Initialize  $s \in \{0\}^{n \times V}$  as an empty array.
4: for target words  $w_i, \dots, w_n$  do
5:   if  $C_{\text{type}} == \text{'symmetric'}$  then
6:      $cv \leftarrow \text{concat}(E(w_{i-(c-1)}), \dots, E(w_{i-1}),$ 
        $E(w_{i+1}), \dots, E(w_{i+(c+1)}))$ 
7:   else if  $C_{\text{type}} == \text{'asymmetric'}$  then
8:      $cv \leftarrow \text{concat}(E(w_{i-(c-1)}), \dots, E(w_{i-1}))$ 
9:   end if
10:   $s[i] \leftarrow \text{ContextNetwork}(cv)$  (Equations 1 to 4)
11: end for

```

target word, w_i , in the input sequence, where i ranges from one to n . For each target word, the algorithm retrieves context words, w_c , based on the specified context type, C_{type} , and context size, C_{size} . The algorithm then concatenates the embeddings of the context words, $E(w_c)$, based on whether the context type is symmetric or asymmetric. The concatenated embeddings, cv , are then passed through the context network model to obtain a prediction probabilities $s[i]$ for target word w_i . The process is repeated until all target words in the input sequence have been processed. The algorithm returns the prediction probabilities s at the end.

C Decision Module

Decision module is used to distinguish each word in syntactic or semantic class based on context network output.

We define Algorithm 4 which takes four inputs: the predicted probabilities s , the decision type dt , the decision hyperparameter t/M and the sequence of words x_{seq} . It returns two outputs: $x_{\text{BoW}}^{\text{syn}}$ and $x_{\text{BoW}}^{\text{sem}}$. The algorithm starts by initializing a zero decision vector, $X \in \{0\}^V$. It then iterates through each target word, w_i , for i ranging from one to n , the number of words in the document. Based on the decision type input, the algorithm checks whether the predicted probability for each target word falls within a certain range, either the top M values, or above a probability threshold. If the predicted probability falls within the specified range, the target word is classified as syntax, otherwise it is classified as content. The decision vector is updated accordingly, with syntax words having a positive value, and content words having a negative value. Syntactic (X_{syn}) and semantic (X_{sem}) binary decision vectors are then constructed based on positive and negative values of words in decision vector. The final output is calculated by taking the element-wise product of X_{syn} and X_{sem} with the BoW representation.

D Datasets and Preprocessing

For dataset preparation, we use SpaCy [Honnibal and Montani, 2017] to tokenize the data. We do not perform any additional pre-processing for SyConNTM. For other models, we remove common stop words, punctuations, and high and low frequency words. We determine high and low frequency words by remov-

Algorithm 4 dec_module

```

1: Input: predicted probability  $s$ , decision type  $dt$ , decision
   hyperparameter  $t/M$ , input word sequence  $x_{\text{seq}}$ , BoW in-
   put  $x_{\text{BoW}}$ 
2: Return:  $x_{\text{BoW}}^{\text{syn}}$  and  $x_{\text{BoW}}^{\text{sem}}$ 
3: Initialize a zero decision vector  $X \in \{0\}^V$ .
4: for target words  $w_i, \dots, w_n$  do
5:   if  $dt = \text{top}_n$  then
6:     if  $w_i \in \text{top}_n(N, s[i])$  then
7:        $w_i^{\text{cls}} \leftarrow \text{syntax}$ 
8:     else
9:        $w_i^{\text{cls}} \leftarrow \text{content}$ 
10:    end if
11:  else if  $dt = \text{probability\_threshold}$  then
12:    if  $w_i \in \text{prob\_threshold}(t, s[i])$  then
13:       $w_i^{\text{cls}} \leftarrow \text{syntax}$ 
14:    else
15:       $w_i^{\text{cls}} \leftarrow \text{content}$ 
16:    end if
17:  end if
18:   $\text{index} \leftarrow \text{index}(w_i, \text{vocabulary})$ 
19:  if  $w_i^{\text{cls}} == \text{syntax}$  then
20:     $X[\text{index}] \leftarrow X[\text{index}] + 1$ 
21:  else
22:     $X[\text{index}] \leftarrow X[\text{index}] - 1$ 
23:  end if
24: end for
25:  $X_{\text{syn}} \leftarrow (X_{\text{syn}1}, \dots, X_{\text{syn}V})$  where  $X_{\text{syn}i} = 1$  if  $X_i > 0$ 
    $\text{else } X_{\text{syn}i} = 0 \forall i = 1, \dots, V$ 
26:  $X_{\text{sem}} \leftarrow (X_{\text{sem}1}, \dots, X_{\text{sem}V})$  where  $X_{\text{sem}i} = 1$  if  $X_i \leq 0$ 
    $\text{else } X_{\text{sem}i} = 0 \forall i = 1, \dots, V$ 
27:  $x_{\text{BoW}}^{\text{syn}} \leftarrow X_{\text{syn}} \cdot x_{\text{BoW}}$ 
28:  $x_{\text{BoW}}^{\text{sem}} \leftarrow X_{\text{sem}} \cdot x_{\text{BoW}}$ 

```

ing the words which occur in more than 75% of the documents and in less than 50 documents.

E Model Details

Let x_{BoW} be the BoW representation of a document. The common encoder used across different models is:

Encoder:

$$\begin{aligned}
X_1 &= \text{ReLU}(L^1 \cdot x_{\text{BoW}} + L_{\text{bias}}^1) \\
X_2 &= \text{dropout}(X_1, p_{\text{keep}}) \\
X_3 &= \text{BatchNorm}(L^2 \cdot X_2 + L_{\text{bias}}^2) \\
\text{enc_out} &= \text{Softplus}(X_3)
\end{aligned}$$

For D-VAE we use the following decoder with Dirichlet prior.

Decoder:

$$\begin{aligned}
X_4 &= \text{BatchNorm}(L^3 \cdot z + L_{\text{bias}}^3) \\
\text{recon} &= \text{LogSoftmax}(X_4)
\end{aligned}$$

For ETM we use the following decoder with both LogNormal and Dirichlet prior. α and δ represent topic and word embeddings respectively.

Table 1: Statistics for the datasets used in our experiments. Given is the average length of one document, the size of the vocabulary before preprocessing and after preprocessing.

	AVG. LENGTH	VOCAB	VOCAB (PREPROCESSED)
20NG	228	25,296	2,213
AMAZON REVIEWS	86	31,667	1,524
AG NEWS	44	25,736	203
IMDB REVIEWS	270	37,552	3,937
ROTTEN TOMATOES	23	6,124	159
YELP REVIEWS	152	30,065	2,690
GOV. REPORT SUMMARY	526	20,795	1,439

Decoder ETM:

$$X_4 = z \cdot \alpha \cdot \delta$$

$$\text{recon} = \text{LogSoftmax}(\text{BatchNorm}(X_4))$$

For SyConNTM we have two decoders, for syntax and semantics of the document.

Decoder SyConNTM:

Syntax:

$$X_4 = \text{BatchNorm}(L^3 \cdot z_{\text{syn}} + L_{\text{bias}}^3)$$

$$\text{recon}_{\text{syn}} = \text{LogSoftmax}(\text{BatchNorm}(X_4))$$

Semantics:

$$X_5 = \text{BatchNorm}(L^4 \cdot z_{\text{sem}} + L_{\text{bias}}^4)$$

$$\text{recon}_{\text{sem}} = \text{LogSoftmax}(\text{BatchNorm}(X_4))$$

F Ablation studies for SyConNTM

To decide on hyperparameter settings for the context network and decision module, we use the total number of stop words that are incorrectly identified as semantic words by the decision module. For this, we use the commonly-used stop word list made available by [Loper and Bird, 2002, NLTK] which contains 179 stop words. Throughout the experiments, we keep $\beta = 2$. All hyperparameters are given in Table 2.

F.1 Finding the optimal context size (c)

In order to determine the most appropriate context size, a series of experiments were conducted utilizing context sizes ranging from one to five for all the datasets. The context type and decision type were fixed at symmetric and top-5, respectively, throughout the experiment. The findings of the aforementioned experiment are presented in Figure 2. As can be observed, as the context size increases, there is a decrease in the number of stop words incorrectly identified as semantic words. Hence, for the rest of experiments, we select context size of five.

F.2 Finding the optimal context type (C_{type})

In order to evaluate the performance of different context types, a series of experiments were conducted utilizing symmetric and asymmetric context types. The results, shown in Figure 3, indicate that the symmetric context type performs significantly better across all datasets. The decision type was fixed at top-5 and the context size was kept constant at 5 throughout the experiment. Based on these results, the symmetric context type was selected.

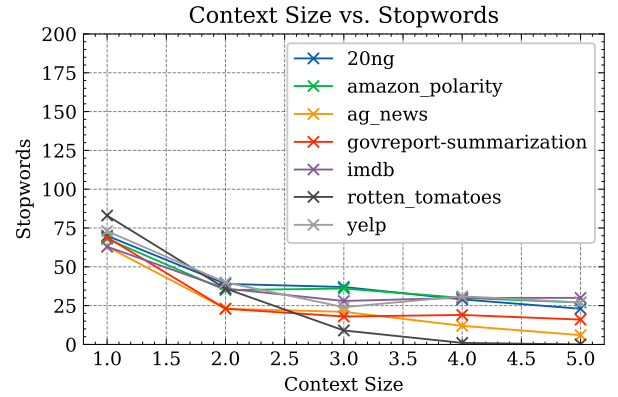


Figure 2: Finding the optimal context size. As the context size increases, there is a decrease in the total number of stop words incorrectly classified as semantic words.

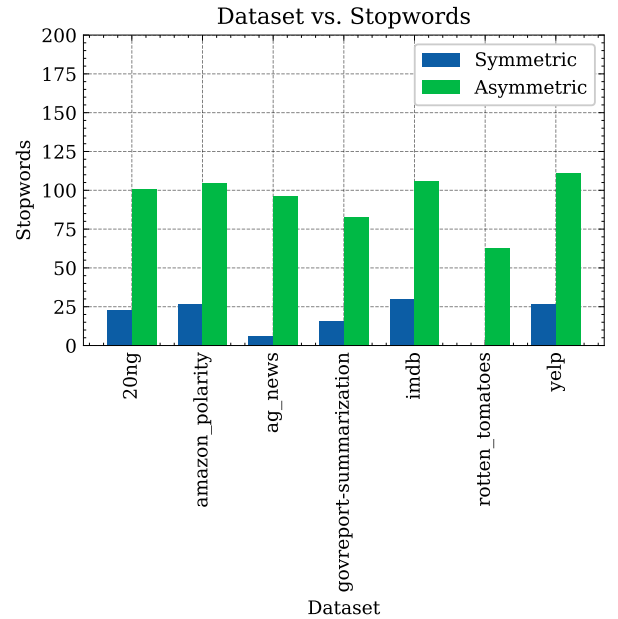


Figure 3: Finding the optimal context type. The results demonstrate that the symmetric context type consistently performs better across all datasets, with a lower number of stop words incorrectly classified as semantic words.

Table 2: Hyperparameter settings for the experiments.

BATCH SIZE	128
TRAIN:VAL:TEST	70:15:15
α	0.02
λ	0.5
β	2
LEARNING RATE	0.001
MAX EPOCHS	100
CONTEXT SIZE	5
CONTEXT TYPE	SYMMETRIC
DECISION TYPE	TOP-M
M	5
L^1	$\mathbb{R}^{\text{VOCAB_SIZE} \times 500}$
L^2	$\mathbb{R}^{500 \times Z}$
L^3	$\mathbb{R}^{Z_{syn} \times \text{VOCAB_SIZE}}$
L^4	$\mathbb{R}^{Z_{sem} \times \text{VOCAB_SIZE}}$
p_{keep}	0.25
WORD EMBEDDINGS	[PENNINGTON <i>et al.</i> , 2014, GLoVe]

F.3 Finding the optimal Top-M/Probability Threshold

In order to evaluate the performance of different decision types, a series of experiments were conducted utilizing top-m ($M = 1, 3, 5$) and probability threshold ($t = 0.1, 0.3, 0.5$) decision types. The results, as shown in Figure 4, indicate that top-m decision type generally performs better, with top-5 performing the best across all datasets. The context type was fixed at symmetric and the context size was kept constant at five throughout the experiment. Based on these results, the top-5 decision type was selected for experiments.

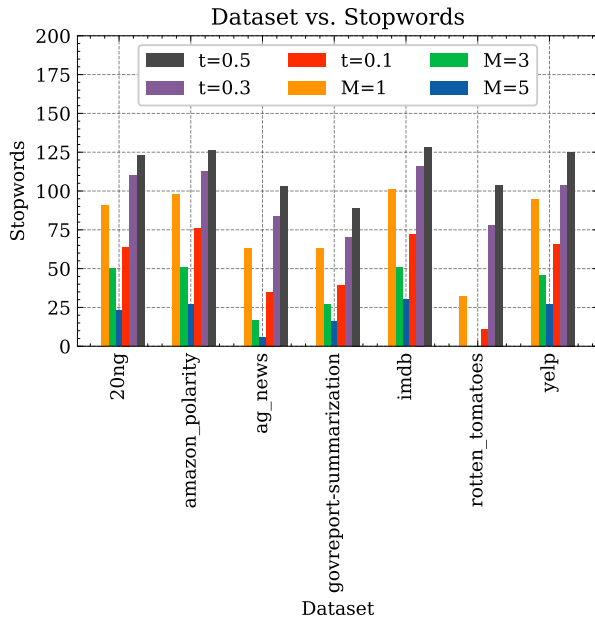


Figure 4: Finding the optimal decision type. The results demonstrate that the top-m decision type generally performs better, with top-5 decision type achieving the highest level of accuracy across all datasets.

G Syntax and Semantic Topics

G.1 Yelp:

Semantic Topics

- casino pool lounge club clubs dj bathrooms crowded security bathroom
- poke vermicelli viet smoothie ra regulars asians gyoza juices ordinary
- lobster creamy ravioli brunch gnocchi coffee crepes pas-try desserts dinner
- stores museum shopping products gear jewelry books target changing classes
- excessive quicker glorified laces borderline helpful lingering patron musky pa
- pour oder manger frugal slip roughly shocking largely thier hubs
- thier kings utterly overlooked quicker flash struggle remained snap anthony
- pizza wings buffalo burgers gyro pita toppings combo joint sandwiches
- place & great maybe sushi try still water chef fish
- business today work customers issues lady needed office delivery situation

Syntactic Topics

- rude gone bad worse charge drive anyone received hair done
- enjoyed great delicious amazing awesome excellent dry thai pork salsa
- un le et les pas est la en du de
- fabled sidney yellowtail revoir instantly bleck spicy swim experiece hiring
- . , () " _ * -
- i me my she in and have in . for

- 230 • always has strip area can vegas selection places may feel
- 231 • were we was with like but good of it and
- 232 • und die ich der ist es das nicht war kann
- 233 • said she back when we after asked service time us

234 LDA Topics

- 235 • . the , a and to you is of it
- 236 • i . was and the my to she a me
- 237 • . and the to they a for , service we
- 238 • . i to the ! a for it \$ and
- 239 • . the , and a was i of it to
- 240 • . the , i and was it a to of
- 241 • we . the to and was our , he us
- 242 • the . , and was to a in room i
- 243 • i , . the to and a of my that
- 244 • ! . and the is i , great this place

245 LDA Topics (preprocessed)

- 246 • sushi buffet place food fish good roll rice rolls eat
- 247 • pizza good sauce like ice cream place cheese sweet tea
- 248 • like time place store know new going car good work
- 249 • food great place service good bar time staff beer love
- 250 • room hotel like vegas place night nice stay people good
- 251 • minutes service said told order time asked got manager
- 252 went
- 253 • food good place burger fries like better sandwich pretty
- 254 meat
- 255 • ordered food like came table restaurant server got wait-
- 256 ress meal
- 257 • coffee love place amazing best delicious menu breakfast
- 258 great wine
- 259 • good chicken salad lunch food great service steak soup
- 260 nice

261 D-VAE Topics

- 262 • , of the to . and - i as with
- 263 • me i she they told to said my out do
- 264 • et un pour le de les la pas en est
- 265 • vegas most are has want its your those area mall
- 266 • worst horrible waited waiting walked attitude arrived paid
- 267 min left
- 268 • delicious sandwich cheese thai sauce rolls special tacos
- 269 salad fries
- 270 • awkward besides parmagiana themselves vodka reorder-
- 271 ing ditching succeeding closer bulging
- 272 • food restaurant was with (were our cheese us it
- 273 • und die das der m ist nicht war zu es
- 274 • highly beautiful thank color helped beyond wedding
- 275 thanks treat recommend

D-VAE Topics (preprocessed)

- 276 • dish try sauce ordered soup meat chicken fried better 277
- cream 278
- 279 • quick kitchen oh fan ordering mexican seating style gets 280
- kept 281
- 282 • hotel room vegas rooms told free called stay stayed said 283
- 284 • food minutes restaurant bar service order drinks ordered 285
- waiter waitress 286
- 287 • amazing great love awesome staff try lunch excellent 288
- attentive family 289
- 290 • bar place great night happy sushi good menu food drinks 291
- 292 • food place service good better great price sushi love buf- 293
- fet 294
- 295 • fries burger place sandwich try meat love ordered good 296
- food 297
- 298 • place bar find love music looking club store vegas people 299
- 300 • told said service order called car know store time manager 301

ETM-D Topics

- 302 • customer called never him guy money business car call 303
- his 304
- 305 • ja selbst dieser wurde meinem anderen bedienung abend 306
- stunden zimmer 307
- 308 • the . this on a and in that i not 309
- 310 • fait avec leur suis toujours leurs du ou en faire 311
- 312 • was were our so had but we table . at 313
- 314 • sips creators goodyear firday nibbling fake carrying elm 315
- rusty ginza 316
- 317 • fake creators elm recognize sips rusty firday bat girlies 318
- exacting 319
- 320 •) but all (as this like a \$ with 321
- 322 • fries plate salsa served soup fish sweet pho french deli- 323
- cious 324
- 325 • free always walk has parking floor looking most area 326
- enjoy 327

ETM-D (preprocessed)

- 311 • rice sushi chicken like food fried roll good buffet fish 312
- 313 • wine great time happy table menu dinner server bar night 314
- 315 • good place pizza food great bar selection love better 316
- prices 317
- 318 • room hotel rooms pool stay floor check strip night vegas 319
- 320 • car phone told called customer work time said problem 321
- business 322
- 323 • minutes food order said waitress asked table took waited 324
- drink 325
- 326 • love look fun kids store staff work new shop feel 327
- 328 • pizza cheese sandwich good sauce salad fries bread or- 329
- dered flavor 330

- friendly staff atmosphere great attentive excellent fantastic clean recommend highly
- expected totally quickly available read mediocre decor believe expensive oh

G.2 20ng:

Semantic Topics

- serious finding behavior basically riding sensitive definitely story yourself funny
- build standard & mark phi hardware mon using text full
- monitor setup ps resolution zero dealer fits feature price installed
- build just already service before us only without called because
- christian man religion us jewish because evil without yourself guns
- riding measure kept water funny sensitive definitely possibly obviously weight
- annual skylab nagorno ti badly inevitable orteg dizziness endowment ssp
- xbiff annual skylab nagorno ti badly inevitable orteg dizziness endowment
- richter annual skylab nagorno badly ti inevitable orteg dizziness ssp
- bhj richter ghj gallant ymon apda kodiaks pena xbiff claris

Syntactic Topics

- _ < = ' : > - ; *)
- point away believe were between them someone world time too
- push papers armed announced rockets contest stars genocide heads et
- this my as that would not but no or time
- published nuclear vote contest generation committee goals penalty powerful guidelines
- controller windows bought uses worked drives box anybody reply large
- george scientific authors contrib north shuttle round soviet ice protect
- entries programs important rates latest postscript solution couple contact digital
- and . that of this is in the to ,
- asking quality york won kings ii st sale byte toronto

LDA Topics

- * .) (the , - to and a
- . , the of - and a to : (
- # . the , of is it to and
- > ' / . x ; , : - the
- . the , to a i and is it of

- the , . of to and in a + 0.010*
- = - < : , (. >) +
- . the , of + 0.004* + 0.004* + 0.004* + 0.004* + 0.003*
- _ . the , - i to a and db
- . , the to that of + 0.020* + 0.018* + 0.018*

LDA Topics (preprocessed)

- game window better good play got win team like best
- people think way know good law government believe human life
- know think like said going mr lot people things time
- space new center university war number april list research press
- use like sure work word system people program available know
- want time people like know going good left new cost
- new card drive disk car video price problem like available
- file available key use data code bit files chip like
- years ago group like center year anybody weeks day believe
- people god think government time believe like year right said

D-VAE Topics

- ' _ > - = (: / <].txt
- never same see good very does really now look might
- senior master color details members press usa jobs league commercial
- was this it are who of " as in they
- increased tape vote hockey myers ahead satellite lock director papers
- never same look law effect let post case enough number
- authors bill sj south stars received flame postings nuclear picture
- . the of are this in (a - it
- sale shipping asking please sell excellent advance condition stereo cd
- division rocket papers premises magazine press mission document myers reported

D-VAE (preprocessed)

- thanks help bit memory problem anybody disk work drive windows
- thought took home subject taken problems place having agree response
- body force questions turn change contact effect longer buy able
- key bit probably chip think know work problem system trying
- good took second home year car great years best bad

418	• bit research law university data key graphics public soft-	
419	ware program	
420	• god law people rights second jesus new human gun issue	
421	• space new people know work years university said let list	
422	• version file graphics dos window set software windows	
423	bit run	
424	• god know jesus bible think people said believe body man	
425	ETM-D Topics	
426	• menial istanbul invoke cliff date metacard scars press	
427	billings rose	
428	• do better than must us might point its true make	
429	• do was we _ i he it you not they	
430	• > * - = : ' <) (;	
431	• menial measure istanbul papers quest trade metacard thus	
432	cliff zephyr	
433	• the . , a -) in and is on	
434	• date menial measure round devices flight attitude quest	
435	papers separate	
436	• getting better drive controller must chip help mac maybe	
437	problems	
438	• character menial washington istanbul english trade track	
439	scars cliff echo	
440	• menial scientists south cliff satellites release waive wing	
441	workspace quest	
442	ETM-D Topics (preprocessed)	
443	• people god jesus christian bible believe know right said	
444	like	
445	• games good bad team game come year went car years	
446	• key bit chip system data number public use access gov-	
447	ernment	
448	• dos video advance set thanks running write speed run	
449	windows	
450	• april david lost center article games city week deal inter-	
451	ested	
452	• drive dos advance thanks card windows anybody disk pc	
453	key	
454	• drive windows file use program pc like software window	
455	info	
456	• people government states law program public university	
457	years year state	
458	• people government public years key right new power law	
459	space	
460	• argument believe come jesus evidence certainly bible	
461	christian true body	

H Motivation from Cognitive Science and Linguistics

Cognitive science and linguistic research suggest that words appear in sentences for two reasons: to fulfill a syntactic function or to convey semantic content. Studies have shown that syntax and content words elicit different patterns of brain activity [Neville *et al.*, 1992] and have distinct developmental tendencies [Brown, 1973]. Additionally, it has been reported that humans acquiring language infer syntactical categories of words based on their short-term context within a document, while content words rely on long-range context [Redington *et al.*, 1998].

Theoretical linguistic considerations and psycholinguistic studies have shown that grammatical functions (*syntax*) and semantics (*content*) are distinct subprocesses within the language domain [Neville *et al.*, 1992]. Syntax and content - both terms refer to a set of words that help the writer create correct sentences that the reader understands. In natural language, syntax refers to the grammatical rules that determine the structure of the sentence. Semantics, on the other hand, deals with the content of a text. Sometimes grammatically correct words do not make sense even though they are grammatically sound. Semantics helps add a layer of meaning so that words together make sense [Culicover and Jackendoff, 2006].

A sentence in a corpus is a complex sequence of words, and the context determines the meaning of each word. For example, the word “sweet” can mean a flavor opposite of sour or something amazing, depending on the context. Similarly, syntax rules are also governed by the context in which they are used. Thus, syntax and semantics are rarely context-free because the surrounding words in a sentence or paragraph directly influence both grammatical rules and the meaning [Friederici and Weissenborn, 2007].

As for context, syntactic words have a short-range context in a document, while semantic words have a long-range context [Griffiths *et al.*, 2004]. Syntactic constraints lead to relatively short-range dependencies that span multiple words within a sentence. Semantic conditions lead to long-range dependencies: different sentences in the same document are likely to have similar content and use similar words. Hence the distinction between syntax and content words in a corpus is context-dependent. The context of words and interpretation is often subjective, leading to additional complications in defining syntax and content words in a corpus.

In tasks such as topic modeling and text classification that depend on the content of a corpus, syntax words are usually removed [Blei *et al.*, 2003]. And since the distinction between syntax and content words is not so clear, one is often limited to using manually curated word lists. Syntax word lists usually consist of stop words and punctuation. They are also task-specific, and sometimes a higher frequency of word occurrence is also considered as syntax [Griffiths *et al.*, 2004].

Curating a list of universal syntactic and semantic words is inefficient because their context is local and corpus dependent. At the local level, it is also possible for a word to have both features in a different context [Griffiths *et al.*, 2004]. Because these lists are not uniform, model performance and benchmarking are also inconsistent [Hoyle *et al.*, 2021].

520 Regarding linguistic composition, this work aims to auto-
521 mate the distinction of syntax and content words based on
522 context.

523 I Limitations

524 Few limitations of the proposed work are:

- 525 • In our framework, we focus on learning topics from syn-
526 tactic words and ignore syntactic classes. Since syntac-
527 tic classes are beneficial for determining part-of-speech
528 (POS) tags, in future developments, we intend to incor-
529 porate the learning of both syntactic classes and topics
530 from syntactic words.
- 531 • The scope of this research is centered on the integration
532 of syntax with neural topic models, and does not specifi-
533 cally address the investigation of syntactic topics or the
534 determination of an optimal number of syntactic topics
535 to be learned within our framework.

536 J Things we tried that did not work

537 In efforts to combine syntax and neural topic models, we
538 tried several architectures, but only three were found to be
539 effective, albeit not as superior as the proposed architecture.
540 The purpose of this section is to present these architectures
541 and provide a brief overview so as to aid future research in
542 this field. The architectures are presented in Figure 5, 6, and
543 7.

544 **Syntax aware attention GAN:** In this study, two adver-
545 sarial generators were implemented to generate the syntactic
546 and semantic components of a document corpus, which were
547 subsequently added to and compared against the original docu-
548 ment using a discriminator as shown in Figure 5. To supervise
549 the model with syntactic and semantic representations, a se-
550 quential network incorporating self-attention was employed.
551 The self-attention matrix and attention scores were utilized to
552 identify words with both long-range and short-range attention.
553 A significant challenge encountered with this model is that,
554 when utilizing a transformer within the sequential model, it
555 operates efficiently only with sub-word tokenizers rather than
556 word-level tokenizers, which are crucial for identifying topics.
557 Additionally, alternative sequential networks such as Recur-
558 rent Neural Networks (RNNs) require processing a significant
559 amount of text (in terms of length) to accurately determine
560 attention scores. Furthermore, determining an appropriate
561 threshold for identifying both long-range and short-range con-
562 text through attention scores remains a problem.

563 **A separate Decision Layer:** In this model, a decision layer
564 was introduced following the encoder in order to determine
565 the allocation of each word in a document to either syntax or
566 semantics as shown in Figure 6. As the decoders in this model
567 do not have separate supervision, the similarity of the outputs
568 from the decoders was maximized in order to enforce their
569 differentiation. The resulting outputs were then combined
570 to form the predicted representation of the document. One
571 limitation of this model is that the separation between syntax
572 and semantic words in the document corpus was not always
573 clear-cut.

A seq-net topic model: In this model, which is similar to 574
the final proposed model, a sequential network was utilized to 575
learn the class of each word as either syntax or semantics in 576
an unsupervised manner as shown in Figure 7. The objective 577
was the loss of the VAE model. While the model was able to 578
differentiate syntactic words to some extent, it was unreliable 579
and unclear why it worked for some datasets and not for others. 580

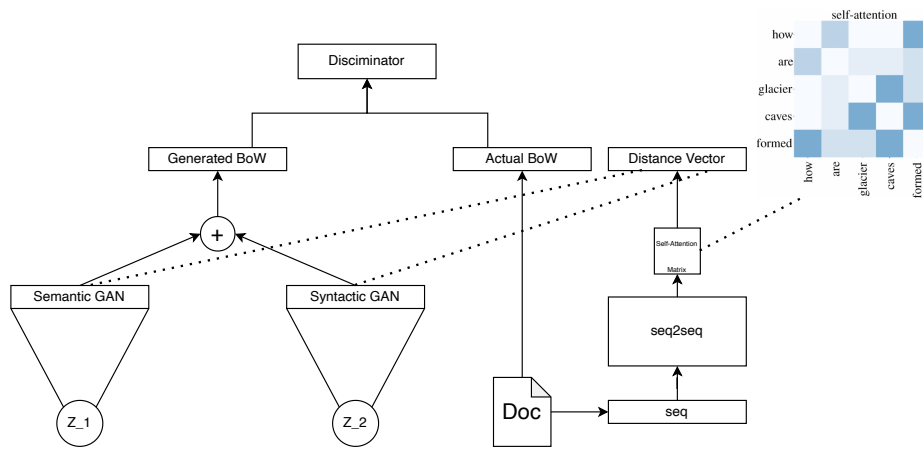


Figure 5: Syntax aware attention GAN - This figure illustrates the implementation of two adversarial generators to generate the syntactic and semantic components of a document corpus, which are then added to and compared against the original document using a discriminator. The separation is guided by the attention scores from a self-attention matrix.

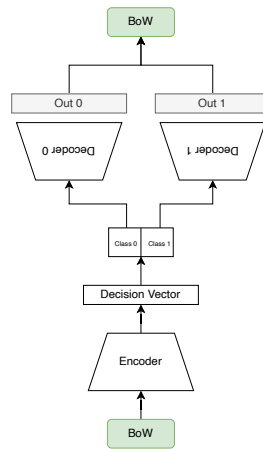


Figure 6: A separate Decision Layer - This figure illustrates the implementation of a decision layer following the encoder to determine the allocation of each word in a document to either syntax or semantics.

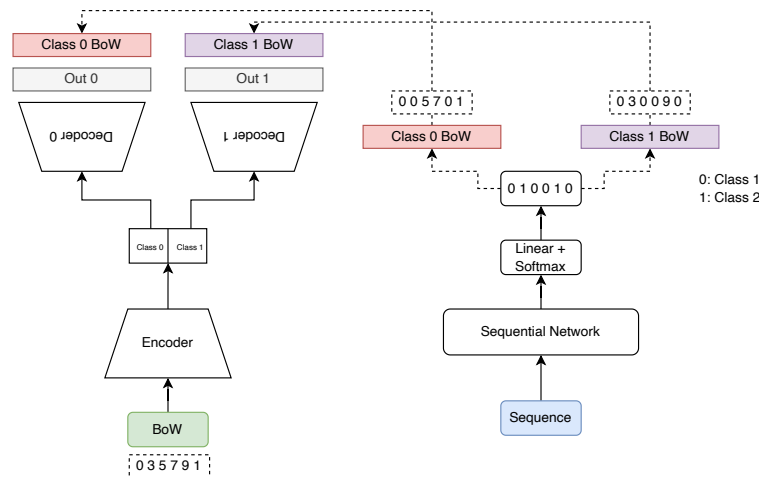


Figure 7: A seq-net topic model - This figure illustrates the implementation of a sequential network to classify each word as either syntax or semantics in an unsupervised manner using a sequential network.

References

- [Blei *et al.*, 2003] David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- [Brown, 1973] R Brown. A first language (pp. 54). Cambridge, MA: Harvard University Press, 1973.
- [Chelba *et al.*, 2017] Ciprian Chelba, Mohammad Norouzi, and Samy Bengio. N-gram language modeling using recurrent neural network estimation. *arXiv preprint arXiv:1703.10724*, 2017.
- [Culicover and Jackendoff, 2006] Peter W Culicover and Ray Jackendoff. The simpler syntax hypothesis. *Trends in cognitive sciences*, 10(9):413–418, 2006.
- [Friederici and Weissenborn, 2007] Angela D Friederici and Jürgen Weissenborn. Mapping sentence form onto meaning: The syntax–semantic interface. *Brain research*, 1146:50–58, 2007.
- [Griffiths *et al.*, 2004] Thomas Griffiths, Mark Steyvers, David Blei, and Joshua Tenenbaum. Integrating topics and syntax. *Advances in neural information processing systems*, 17, 2004.
- [Honnibal and Montani, 2017] Matthew Honnibal and Ines Montani. spaCy: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- [Hoyle *et al.*, 2021] Alexander Hoyle, Pranav Goel, Andrew Hian-Cheong, Denis Peskov, Jordan Boyd-Graber, and Philip Resnik. Is automated topic model evaluation broken? the incoherence of coherence. *Advances in Neural Information Processing Systems*, 34:2018–2033, 2021.
- [Ioffe and Szegedy, 2015] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, pages 448–456. PMLR, 2015.
- [Loper and Bird, 2002] Edward Loper and Steven Bird. Nltk: The natural language toolkit. *arXiv preprint cs/0205028*, 2002.
- [Neville *et al.*, 1992] Helen J Neville, Debra L Mills, and Donald S Lawson. Fractionating language: Different neural subsystems with different sensitive periods. *Cerebral Cortex*, 2(3):244–258, 1992.
- [Pennington *et al.*, 2014] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on EMNLP (EMNLP)*, pages 1532–1543, 2014.
- [Redington *et al.*, 1998] Martin Redington, Nick Chater, and Steven Finch. Distributional information: A powerful cue for acquiring syntactic categories. *Cognitive science*, 22(4):425–469, 1998.
- [Srivastava *et al.*, 2014] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014.