

Sparsity Regularization for Chinese Spelling Check

Anonymous ACL submission

Abstract

The Chinese Spelling Check (CSC) research objective is to detect and correct the spelling errors in the input. Generally, the number of incorrect characters in the input is far less than the correct, so the error probability sequence of input sentence predicted by the detection module should be sparse and sharp. However, all existing work has ignored this problem. In this paper, we add a sparsity regularization item to the objective function to make the output of the detection module close to sparse and sharp. We study two kinds of regularization: L_1 regularization and minimum entropy regularization. Extensive experiments on the SIGHAN show that the sparsity regularization proposed in this paper can effectively improve the performance of the CSC model while without increasing the computational complexity. In addition, the robustness experiment results show that our method is robust.

1 Introduction

Chinese Spelling Check (CSC) task is an important task in Natural Language Processing (NLP) community, which aims to correct spelling errors in Chinese input (Yu & Li, 2014). There are usually two types of errors in Chinese text editing and text Recognition: errors of visually similar characters caused by Optical Character Recognition (OCR) or Wubi input method, and errors caused by the misuse of phonologically similar characters (Duan Jianyong, 2021). Spelling errors will affect the semantics of the sentence and then negatively impact downstream text processing tasks. Therefore, it is necessary and practical to study the CSC.

The CSC task consists of two subtasks: error detection and error correction. For the CSC task, complete Chinese sentences with/without spelling errors will be given as the input. For the detection subtask, the detection module should return the locations of the incorrect characters. For the correction subtask, the correct module should point out the correct characters based on the output of detection subtask. The error correction problem is a follow-up problem of error detection for checking spelling errors (Hládek et al., 2020; Tseng et al., 2015).

The input of the detection module is Chinese sentence $X = (x_1, x_2, \dots, x_n)$, and x_i represents the i -th character in the input. The output is a probability vector $P = (p_1, p_2, \dots, p_n)$, p_i represents the probability that x_i is an incorrect character. In the input of the CSC, the number of incorrect characters is usually far less than correct characters. Therefore, the probability vector output by the detection module should be sparse and sharp. For example, as shown in Figure 1, for the input “我以前想要高诉你，可是忘了，真户秃。” (Correct sentence: “我以前想要告诉你，可是忘了，真糊涂。”), meaning is “I wanted to tell you before, but I forgot, so confused.”), the best output of detection module is a probability vector which with 1 at the location of “高” (告, tell), and “户秃” (糊涂, confused), and 0 at the other positions. The best output of the detection module is sparse and sharp, but this is a tremendous challenge for the detection module.

Based on the above problem, and inspired by the attention with sparsity regularization (Zhang et al., 2018), this paper uses sparsity regularization item to constrain the probability vector output of the detection module. We study two kinds of sparsity regularization: L_1 regularization and minimum entropy regularization. L_1 regularization will constrain the sum of the absolute value of the probability vector output by

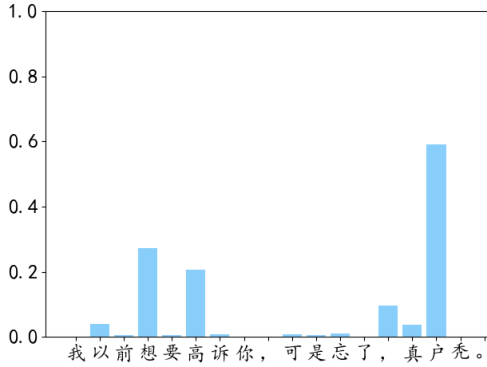


Figure 1: The example of the output of spelling detection.

82 detection module, so as to promote the sparsity of
 83 the probability vector. Minimum entropy
 84 regularization minimizes the entropy of the
 85 probability vector output by the detection module.
 86 As we all know, the more certain, the smaller
 87 entropy. So, when p_i is 0 or 1, the entropy of the
 88 probability vector P is the smallest. Therefore,
 89 minimum entropy regularization can make the
 90 probability vector output by the detection module
 91 close to sparse and sharp.

92 In order to test the effectiveness of our
 93 approach, we evaluated our method on the
 94 SIGHAN dataset. On the premise of using the
 95 same training set, validation set, and test set and
 96 keeping the same experimental setup, the two
 97 kinds of sparsity regularization used in this paper
 98 can effectively improve the performance of the
 99 CSC model. In addition, we use five new test sets
 100 to test the robustness of our approach. In these test
 101 sets, the proportion of without spelling errors
 102 sentences increases sequentially. The main
 103 contributions of this paper are as follows:

- 104 • We firstly propose use sparsity
 105 regularization to improve the performance
 106 of CSC model while without increasing the
 107 computational complexity of the model, and
 108 analyze why the sparsity regularization can
 109 be effective in CSC task from theoretical
 110 level.
- 111 • We study two kinds of sparsity
 112 regularization: L_1 regularization and
 113 minimum entropy regularization, and test
 114 the effectiveness and robustness of these
 115 two kinds of sparsity regularization through
 116 experiments.

117 2 Related Work

118 Early CSC methods were based on rules such as
 119 chunk, syntax, and grammar to determine spelling
 120 errors in the input (Hirst & Budanitsky,
 121 2005). This kind of approach has high correction
 122 accuracy but poor generalization performance.
 123 Because the rules need to be defined by experts
 124 based on knowledge and experience, the rules can
 125 only cover limited situations, and the formulation
 126 of these rules takes a lot of time. Then the
 127 researchers tried to use the statistical language
 128 model to complete the CSC task. First, replace the
 129 characters in the input sentence sequentially
 130 according to the confusion set, then use the
 131 statistical language model to score these replaced
 132 sentences, and give correction suggestions based
 133 on the score (Huang et al., 2014; Yeh et al., 2017).
 134 But the method based on the statistical language
 135 model has shortcomings, that is, the expressive
 136 and learning capability of the model are poor.

137 Deep learning technology can contribute to this
 138 problem. The literature (Duan et al., 2019)
 139 employs Bi-LSTM combined with conditional
 140 random field (CRF) to achieve the error detection,
 141 and CRF was used to predict the best annotation
 142 sequence. The literature (Wang et al., 2021)
 143 employs Lattice-LSTM to dynamically integrate
 144 characters, words, and pinyin information, and
 145 then the CRF layer detected errors according to
 146 the integrated information. The focus of these
 147 works is to integrate multiple features to improve
 148 the performance of the model. Different from
 149 these works, the literature (Wang et al., 2019)
 150 applies the pointer-generation model, which
 151 excellent performance in automatic
 152 summarization tasks, to the error correction task.
 153 For each character in the input sentence, the
 154 pointer-generation model decides whether to copy
 155 the character or replace by the character generated
 156 by the Generation network according to the
 157 probability output by the detection module. As far
 158 as we know, this is the first CSC model using the
 159 soft strategy. However, these methods are difficult
 160 to adopt when the corpus is limited, because these
 161 models need enough labeled corpus for training
 162 (Tan et al., 2020).

163 In recent years, BERT has been applied to
 164 many natural language processing tasks as a pre-
 165 trained representation model and has shown
 166 strong performance (Devlin et al., 2019). Some
 167 researchers applied the BERT to the CSC task and
 168 achieved state-of-the-art (SOTA) results. The

169 literature (Cheng et al., 2020) employs BERT to
 170 obtain the distributed representation of the input
 171 and employs the graph convolutional neural
 172 network (GCN) to learn the similarity knowledge
 173 in phonological and visual pronunciation to
 174 complete CSC task. The literature (Zhang et al.,
 175 2020) employs the Bi-GRU network as the
 176 detection module and then employs BERT to
 177 complete the error correction task according to the
 178 output of the detection module. The literature (Li
 179 et al., 2021) proposes a cloze-style detector-
 180 corrector framework (DCSpell) that firstly detects
 181 whether a character is erroneous before correcting
 182 it. DCSpell employs the discriminator of
 183 ELECTRA (Clark et al., 2016) as the Detector to
 184 detect the positions of incorrect characters and
 185 then employs BERT to correcting it.

186 The model based on BERT has achieved the
 187 SOTA on the CSC task, but some researchers
 188 pointed out that the error detection capability of
 189 BERT is poor (Zhang et al., 2020; Hong et al.,
 190 2019). Therefore, how to improve the detection
 191 capability of the CSC model base on BERT is a
 192 critical issue, which is also the focus of this paper.

193 3 Method

194 3.1 Soft-Masked BERT

195 The sparsity regularization method proposed in
 196 this paper is an improvement method that can
 197 improve the performance of the CSC model. The
 198 sparsity regularization method can be widely used
 199 in CSC tasks at a theoretical level. In order to
 200 show the superiority of our method, we use the
 201 strong baseline model Soft-Masked BERT (SM)
 202 (Zhang et al., 2020). The model structure is shown
 203 in Figure 2.

204 The SM model is composed of a detection
 205 module and correction module. The detection
 206 module uses a bidirectional GRU (Bi-GRU)
 207 model and the correction module uses BERT. The
 208 input of the detection module is the word
 209 embedding sequence $E = (e_1, e_2, \dots, e_n)$, e_i
 210 represents the embedding representation of the
 211 character x_i . The e_i consists of word embedding,
 212 position embedding, and segment embedding like
 213 the input of BERT. The output of the detection
 214 module is a probability vector $P = (p_1, p_2, \dots, p_n)$,
 215 p_i represents the probability that the character x_i
 216 is incorrect.

217 In the SM model, for each character of the
 218 sequence, the p_i is calculated as

$$219 \quad p_i = P_d(g_i = 1|X) = \sigma(W_d h_i^d + b_d) \quad (1)$$

220 where $P_d(g_i = 1|X)$ denotes the conditional
 221 probability predicted by the detection module, σ
 222 denotes the sigmoid function, h_i^d is the hidden
 223 state of the Bi-GRU, W_d and b_d are parameters of
 224 the model.

225 According to the probability output by the
 226 detection module, the soft-masked embedding e'_i
 227 for the x_i is calculated as **Formula (2)**.

$$228 \quad e'_i = p_i \cdot e_{\text{mask}} + (1 - p_i) \cdot e_i \quad (2)$$

229 where e_i is the embedding of x_i outputted by
 230 BERT and e_{mask} is the mask embedding of
 231 character [MASK] outputted by BERT. When the
 232 p_i close to 0, the e'_i is close to the e_i ; otherwise, it
 233 is close to the e_{mask} , which can reduce the impact
 234 of incorrect characters on the semantics of the
 235 input sentence.

236 The correction network is a sequential multi-
 237 class labeling task, which is based on the BERT
 238 model. The input of the correction network is the
 239 sequence $E' = (e'_1, e'_2, \dots, e'_n)$ that consists of
 240 soft-masked embedding. The output is a character
 241 sequence corrected by correction network, we
 242 define it as $Y = (y_1, y_2, \dots, y_n)$.

243 BERT is based on the transformer model and
 244 the hidden states of the output of BERT are
 245 denoted as $H^c = (h_1^c, h_2^c, \dots, h_n^c)$. For each x_i , the
 246 corresponding output of the SM model is
 247 calculated as follows

$$248 \quad P_c(y_i = j | X) = \text{softmax}(W h_i^c + b) [j] \quad (3)$$

249 Where $P_c(y_i = j | X)$ define the probability that
 250 x_i is corrected as character j in the output of the
 251 correct module, softmax denote the softmax
 252 function, W and b are parameters, h_i^c is the hidden
 253 state which is obtained by residual connection,
 254 defined as

$$255 \quad h_i^c = h_i^e + e_i \quad (4)$$

256 The training of the SM model is conducted end-
 257 to-end. For the dataset $D = \{(X^j, Y^j)\}_{j=1}^N$, X^j is
 258 the j -th input sequence, Y^j is the correct sequence
 259 of the X^j , N denotes that there are N with\without
 260 error sentences in the dataset. The goal of the
 261 training process is to optimize the two objective
 262 functions, corresponding to the detection module
 263 and error module.

$$264 \quad \mathcal{L}_d = \sum_{j=1}^N \sum_{i=1}^n \log P_d(g_i^j | X_j, \theta_d) \quad (5)$$

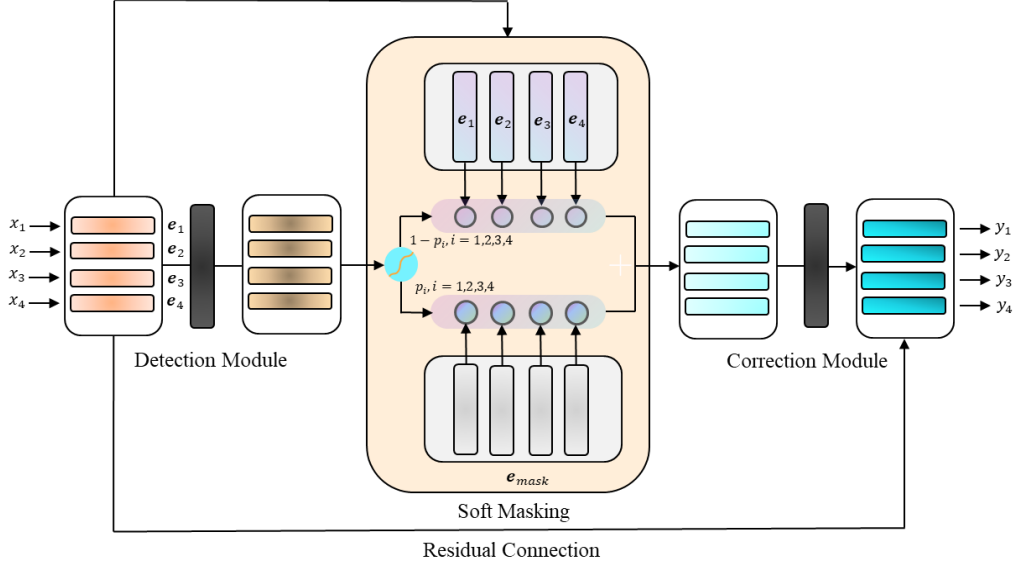


Figure 2: Structure of Soft-Masked BERT.

$$\mathcal{L}_c = \sum_{j=1}^N \sum_{i=1}^n \log P_c(y_i^j | X_j, \theta_c) \quad (6)$$

Where \mathcal{L}_d and \mathcal{L}_c denotes the objective of the detection module and the correction module, respectively. θ_d and θ_c denote the parameters of the detection module and correction module. The linear combination of \mathcal{L}_d and \mathcal{L}_c is the overall training objective of the SM model.

$$\mathcal{L} = \lambda \cdot \mathcal{L}_c + (1 - \lambda) \cdot \mathcal{L}_d \quad (7)$$

Where λ is a hyper-parameter and $\lambda \in (0,1)$.

3.2 Sparsity Regularization

The output of the detection network is a probability vector $P = (p_1, p_2, \dots, p_n)$. In the input of the CSC, the number of incorrect characters is usually far less than correct characters. Therefore, the probability vector P should be sparse and sharp. Accordingly, we attempt to study the sparsity regularization method for the CSC task. Specifically, we add a sparsity regularization item to the objective function as a penalty term to constraint the distribution of probability vector P .

$$\mathcal{L}' = \lambda \cdot \mathcal{L}_c + (1 - \lambda) \cdot \mathcal{L}_d + \beta \cdot \sum_{j=1}^N R(P^j) \quad (8)$$

in which $R(P^j)$ is the sparsity regularization item, β is the hyper-parameter which balances the sparsity regularization term and the log-likelihood function. In this paper, we study two sparsity regularization methods: L_1 regularization and minimum entropy regularization.

292

L_1 regularization L_0 regularization and L_1 regularization are the most commonly used regularization method to control the sparsity of parameters. The return of L_0 regularization is positively correlated with the number of non-zero elements in the input, so minimizing the L_0 regularization term can make the parameters sparse. But L_0 regularization is a non-convex function, which means it is difficult to optimize. Therefore, researchers usually use L_1 regularization instead of L_0 regularization. For the probability vector P^j , L_1 regularization item is:

$$L_1(P^j) = \sum_{i=1}^n |p_i^j| \quad (9)$$

Where p_i^j is a probability and $p_i^j \in (0,1)$. Since the optimization strategy in this paper is to maximize the objective function \mathcal{L}' , $R(P^j)$ based on the L_1 regularization in this paper is defined as:

$$R(P^j) = - \sum_{i=1}^n p_i^j \quad (10)$$

In the CSC task, maximizing the $R(P^j)$ based on L_1 regularization will cause all p_i^j to approach 0. For j -th sentence in the dataset, If the character at the k -th position is incorrect, \mathcal{L}_d will constrain p_k^j to approach 1. So L_1 regularization can control the distribution of P^j close to sparse theoretically. However, since the L_1 regularization item expects all p_i^j close to 0, even if the character is incorrect. This means that L_1 regularization cannot contribute to the probability vector P^j to be sharp.

321 We aim to contribute to the probability vectors P^j
 322 to be sparse and sharp. Accordingly, we introduce
 323 the minimum entropy approach.

324
 325 **Minimum entropy regularization** In a statistical
 326 learning community, Minimum entropy is used
 327 less frequently than maximum entropy and L_1
 328 regularization, there is also has been used as a
 329 regularization item. The literature (Grandvalet &
 330 Bengio, 2005) uses minimum entropy
 331 regularization to help the classifier based on semi-
 332 supervised learning make full use of beneficial
 333 unlabeled data, aiming to improve the robustness
 334 of the classifier. The literature (Zhang et al., 2018)
 335 uses the minimum entropy regularization item to
 336 contribute the distribution of attention vectors to
 337 sparse and sharp.

338 In the CSC task, we argue that P^j should be
 339 sparse and sharp. That is, the probability vector P^j
 340 should have low entropy. We use the minimum
 341 entropy regularization to meet this prior as follows.

$$\text{Ent}(P^j) = - \sum_{i=1}^n p_i^j \log(p_i^j) \quad (11)$$

$$R(P^j) = - \text{Ent}(P^j) = \sum_{i=1}^n p_i^j \log(p_i^j) \quad (12)$$

344 The derivative of the minimum entropy
 345 function $f(x) = x \log(x)$ is

$$f'(x) = 1 + \log(x) \quad (13)$$

347 When $x \in (0, 1/e)$, $f'(x) < 0$, and when $x \in$
 348 $(1/e, 1)$, $f'(x) > 0$, and $f(0) = f(1) = 0$.
 349 Therefore, when $x \in (0,1)$, maximizing the
 350 minimum entropy function will make the value of
 351 x to 0 or 1. For the CSC task, the value of the
 352 probability vector output by the detection module
 353 is preferably 0 or 1, corresponding to the location
 354 of the correct character and incorrect character. In
 355 summary, minimum entropy regularization can
 356 contribute to the CSC task at the theoretical level.

357 4 Experiment

358 In this section, we analyzed the experimental
 359 results of our method on the CSC task in detail
 360 and compared it with the results of baselines
 361 based on different methods. We also discussed
 362 how to choose model parameters and verified the
 363 robustness of our method.

"Id":	"A2-0023-1"
"Original_text":	"下个星期, 我跟我朋晴打算去法国玩儿。"
"Wrong_ids":	" [9] "
"Correct_text":	"下个星期, 我跟我朋友打算去法国玩儿。"

Table 1: Data sample.

Train and Val Data	Line
(Wang et al., 2018)	271,329
SIGHAN13	700
SIGHAN14	1301
SIGHAN15	970
Total	274,300
Test Data	Line (Line with error character)
SIGHAN13	1000(971)
SIGHAN14	1062(529)
SIGHAN15	1100(550)
Total	3162(2074)

Table 2: The statistic of the data.

364 4.1 Dataset

365 **Train and Validation set** In this paper, the
 366 training data is composed of (Wu et al., 2013), (Yu
 367 et al., 2014), and (Tseng et al., 2015). Following
 368 the literature (Wang et al., 2019), we add an
 369 additional 271K samples to our training data,
 370 which are provided by the literature (Wang et al.,
 371 2018). To observe the training process in real-time,
 372 we randomly selected 10% of these training data
 373 as the validation set and the other sample as the
 374 training set. We show one data sample in Table 1.

375 As shown in Table 1, a sample consists of four
 376 parts: Id, Original_text, Wrong_ids, and
 377 Correct_text. The Id is the id of the sample. The
 378 original_text denotes the input of the CSC task
 379 with/without incorrect characters. The Wrong_ids
 380 point out the location of the incorrect character,
 381 and it is an empty list when there are without
 382 incorrect characters in the input. The Correct_text
 383 denotes the corrected sentence, and it is the same
 384 as the Original_text when there are without
 385 incorrect characters in the input.

386 **Test set** Following the literature (Wang et al.,
 387 2019), we used the test dataset from the
 388 SIGHAN13 (Wu et al., 2013), SIGHAN14 (Yu et
 389 al., 2014), and SIGHAN15 (Tseng et al., 2015)
 390 benchmarks. Like related work, we also use
 391 OpenCC1 to convert the characters in the test set
 392 from traditional Chinese to simplified Chinese.
 393 The statistic of the data is listed in Table 2.

¹ <https://github.com/BYVoid/OpenCC>

394 4.2 Evaluation Metrics

395 For this paper, we report precision (P), recall (R),
396 and F1 scores as evaluation metrics, which are
397 used by almost all CSC tasks. We report the
398 sentence-level metrics on the detection sub-task
399 and correction sub-task, i.e., we consider an input
400 sentence has been correctly only when the output
401 sentences of the model are completely consistent
402 with our expected. As shown in **Formula 24**
403 through 26, the metrics are measured with the
404 help of the confusion matrix which is shown in
405 Table 3 following the literature (Tseng et al.,
406 2015).

$$407 \quad \text{Precision} = \text{TP} / (\text{TP} + \text{FP}) \quad (14)$$

$$408 \quad \text{Recall} = \text{TP} / (\text{TP} + \text{FN}) \quad (15)$$

$$409 \quad \text{F1} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} \quad (16)$$

410 For the CSC task, the better precision means
411 that the fewer correct characters are recognized as
412 incorrect characters, and the better recall means
413 that the more incorrect characters can be
414 detected/corrected.

415 4.3 Baselines

416 We compare our method to the following
417 baselines.

418 **FASpell**² (Hong et al., 2019): The FASpell
419 model uses the pre-trained model BERT as a
420 denoising autoencoder (DAE) and decoder and
421 uses confidence to filter candidate modifications
422 instead of confusion set.

423 **DCSpell** (Li et al., 2021): The DCSpell model
424 proposes a cloze-style detector-corrector
425 framework that firstly detects whether a character
426 is erroneous before correcting it. DCSpell uses the
427 pre-trained discriminator ELECTRA (Clark et al.,
428 2016) as a detection module and uses BERT to
429 correct the incorrect characters.

430 **BERT-Finetune** (BERT-FT) (Devlin et al.,
431 2019): Add a Softmax layer after the last layer of
432 the BERT model, and use the training set to fine-
433 tune the BERT model so that it can complete the
434 CSC task.

435 **Soft-Masked BERT**³ (Zhang et al., 2020): The
436 SM model uses the Bi-GRU as the detection
437 module, and then rewrites the input embedding
438 according to the results of the detection module,

² <https://github.com/iqiyi/FASpell>

³ <https://github.com/gitabtion/SoftMaskedBert-PyTorch>

Confusion matrix		System Result	
		Positive (Erroneous)	Negative (Correct)
Gold Standard	Positive	TP	FN
	Negative	FP	TN

Table 3: Confusion matrix.

439 and then enters it into the BERT model to
440 complete the error correction.

441 4.4 Experimental Setting

442 In order to show the improvement of our method
443 on the performance of the SM model, we use the
444 same experimental setting as the literature (Zhang
445 et al., 2020). Our code is based on BERT⁴.
446 Following the literature (Zhang et al., 2020), when
447 the fine-tuning process, we kept the default hyper-
448 parameters and only used Adam. We did not use
449 the dynamic learning rate strategy. The learning
450 rate is set to $2e^{-5}$, the hidden dimension of Bi-
451 GRU is set to 256, the batch size is set to 320, and
452 the training epochs are set to 100. In order to save
453 time, we use early stop strategy. We stop the
454 training process early when the validation loss
455 does not decrease for ten consecutive epochs. The
456 value of λ is set to 0.8, which is the best value
457 reported by (Zhang et al., 2020)

458 4.5 Result analysis

459 We list the experimental results of our method and
460 four baseline models on the test set. It can be seen
461 from Table 4, the precision and F1 score of the
462 SM model with sparsity regularization method are
463 better than baselines including SM model. Such
464 experimental results show the effectiveness of the
465 sparsity regularization method. In addition, we
466 found that both the detection and correction recall
467 of the SM model are better than that of the BERT-
468 FT model. Such experimental results show that
469 the Bi-GRU model is more effective than the
470 BERT model as the detection module, which
471 further supports the conclusion of literature
472 (Zhang et al., 2020) and (Hong et al., 2019) that
473 the error detection capability of BERT is poor.

474 Compared with the L_1 regularization, the
475 minimum entropy regularization is better. This
476 result is consistent with the analysis in Section 3.2.
477 The L_1 regularization only can contribute to the
478 sparsity of the probability vector output by the
479 detection module, while minimum entropy
480 regularization can contribute to sparsity and

⁴ <https://github.com/google-research/bert>

Test set	Method	Detection			Correction		
		P	R	F1	P	R	F1
SIGHAN	FASpell	36.4	43.3	39.6	30.1	35.7	32.7
	DCSpell	62.0	54.8	58.2	57.9	51.1	54.3
	BERT-FT	84.9	60.4	70.6	84.3	57.8	68.4
	SM	84.0	61.5	71.0	83.4	58.8	68.9
	SM + L ₁	85.5	61.4	71.5	84.9	58.8	69.5
	SM + Entropy	86.1	62.6	72.5	85.6	60.1	70.6

Table 4: Performance of different methods on test set.

Method	Pe%	Detection			Correction		
		P	R	F1	P	R	F1
SM	20%	77.6	61.5	68.6	76.8	58.8	66.6
	40%	72.9	61.5	66.7	72.0	58.8	64.7
	60%	67.9	61.5	64.5	66.9	58.8	62.6
	80%	63.7	61.5	62.5	62.6	58.8	60.6
	100%	59.7	61.5	60.6	58.6	58.7	58.7
SM + L ₁	20%	78.8	61.4	69.0	78.0	58.8	67.1
	40%	73.5	61.4	66.9	72.7	58.8	65.0
	60%	68.6	61.4	64.8	67.7	58.8	63.0
	80%	64.4	61.4	62.9	63.4	58.8	61.0
	100%	60.5	61.4	60.9	59.4	58.8	59.1
SM + Entropy	20%	80.4	62.6	70.4	79.8	60.1	68.5
	40%	75.0	62.6	68.2	74.2	60.1	66.4
	60%	70.3	62.6	66.2	69.4	60.1	64.4
	80%	66.1	62.6	64.3	65.2	60.1	62.5
	100%	62.5	62.6	62.5	61.5	60.1	60.8

Table 5: The result of robustness study.

sharpness. So, it is reasonable that the minimum entropy regularization method is better than L₁ regularization.

Further analysis of the experimental results shows that the detection precision and recall of the SM model with the minimum entropy regularization are better than the original SM model. This finding suggests that the minimum entropy regularization method can help the SM model fewer recognize correct characters as incorrect characters, and help it detect more incorrect characters. This result means that the minimum entropy regularization method that we have proposed therefore assists in improving the detection ability of the CSC model thoroughly. It should be noted that our method does not need any change in model architecture, that is, our method is theoretically applicable to all Chinese spelling error detection (CSED) tasks and the CSC tasks that include error detection sub-task.

In addition, the minimum entropy regularization improves the F1 score of error detection more than that of error correction. This result shows that the SM model with the minimum entropy regularization can detect and

correct some sentences that the original SM model detects the incorrect characters but failed to correct. This result may be explained by the fact that the improvement of minimum entropy regularization for error correction comes from the advance of error detection. It can also make the probability vector output by the detection module sparse and sharp. According to this result, we can infer that our method is more suitable for the CSC model using soft strategy than the CSC model using fixed threshold strategy.

4.6 Robustness Study

Few input sentences contain incorrect characters in practical applications, such as text editing and correction. So, the robustness of the CSC model is very important when there are more and more without incorrect characters sentences in the test set. Therefore, we tested the robustness of our method.

We use the following methods to increase the number of sentences in the test set that without incorrect characters:

- 1) Randomly select Pe% of having incorrect characters samples from the test set;

Pe%	0%	20%	40%	60%	80%	100%
L_1	0.6	0.5	0.3	0.4	0.4	0.4
Entropy	1.7	1.9	1.7	1.8	1.9	2.1

Table 6: Effect of regularization item.

Method	β	Detection	Correction
		F1	F1
L_1	0.02	70.5	68.7
	0.04	71.5	69.5
	0.06	70.8	69.0
	0.08	70.9	69.1

Table 7: Effect of β (L_1 regularization).

Method	β	Detection	Correction
		F1	F1
L_1	0.02	72.1	70.1
	0.04	70.5	68.7
	0.06	72.5	70.6
	0.08	71.2	69.2

Table 8: Effect of β (minimum entropy regularization).

- 2) For each selected sample, use the Correct_text of the data to replace its Original_text, thereby generating a new test without incorrect characters sample;
- 3) Add all the new samples generated in step 2) to the test set.

We tested the performance of the SM model and our method when Pe% is 20%/40%/60%/80%/100%, and the results are reported in Table 5. During the robustness study, we will not fine-tune.

When Pe% takes different values, the recall of the three models always remains static. Because we just added some new without incorrect character samples to the test set. Therefore, the containing incorrect characters sample that the CSC model can detect and correct will not change while without fine-tune. So, the recall will not change.

As shown in Table 5, when the number of without incorrect characters samples in the test set increases, the precision of the three models decreases. However, minimum entropy regularization and L_1 regularization can still improve the performance of the CSC model. We study the improvement of the F1 score of the two kinds of sparsity regularization methods when Pe% takes different values.

Table 6 reports the result, it can be seen that the L_1 regularization can continually improve the performance of the CSC model when Pe% takes different values, but the improvement is slight.

The minimum entropy regularization is better than the L_1 regularization, and as the value of Pe% increases, the minimum entropy regularization also shows a gradual improvement in the model performance. Therefore, it can be concluded that the two kinds of sparsity regularization studied in this paper are robust when the number of without incorrect characters samples in the test set increases.

4.7 Effect of Hyper Parameter β

We use the hyperparameter β to balance the objective function and the sparsity regularization item. The value of β is determined by the correcting F1 score on the test set. According to the experimental results, the best value of β is 0.04 when we use L_1 regularization, and the best value of β is 0.06 when we use minimum entropy regularization. The detailed experimental results are reported in Table 7 and Table 8.

5 Conclusion

This paper proposes to use the sparse regularization method to improve the performance of the CSC model. Specifically, we studied two kinds of sparsity regularization methods: L_1 regularization and minimum entropy regularization. By adding a sparsity regularization item to the objective function, the output of the detection module is close to sparse and sharp, so as to improve the performance of the CSC model. Although we only conducted experiments on the SM model, it should be noted that our method does not need any change in model architecture. Therefore, our method can apply to all Chinese spelling error detection (CSED) tasks and the CSC tasks that include the sub-task theoretical level. Experiments on the SIGHAN test set show that our method can effectively improve the performance of the CSC model, especially the minimum entropy regularization method.

References

- Yu, J., & Li, Z. (2014). Chinese spelling error detection and correction based on language model, pronunciation, and shape. Proceedings of The Third CIPS-SIGHAN Joint Conference on Chinese Language Processing.
- Duan Jianyong, Y. Y., Wang Hao. (2021). Chinese Spelling Correction Method Based on Transformer Local Information and Syntax Enhancement Architecture. Acta Scientiarum Naturalium

- 611 Universitatis Pekinensis, 57(1), 61-67. 658
612 <https://doi.org/10.13209/j.0479-8023.2020.081>. 659
- 613 Hládek, D., Staš, J., & Pleva, M. (2020). Survey of 660
614 automatic spelling correction. *Electronics*, 9(10), 661
615 1670.
- 616 Tseng, Y.-H., Lee, L.-H., Chang, L.-P., & Chen, H.-H. 662
617 (2015). Introduction to SIGHAN 2015 bake-off for 663
618 Chinese spelling check. *Proceedings of the Eighth 664*
619 SIGHAN Workshop on Chinese Language 665
620 Processing. 666
- 621 Zhang, J., Zhao, Y., Li, H., & Zong, C. (2018). 667
622 Attention with sparsity regularization for neural 668
623 machine translation and summarization. 669
624 *IEEE/ACM Transactions on Audio, Speech, and 670*
625 *Language Processing*, 27(3), 507-518. 671
- 626 Hirst, G., & Budanitsky, A. (2005). Correcting real- 672
627 word spelling errors by restoring lexical cohesion. 673
628 *Natural Language Engineering*, 11(1), 87-111. 674
- 629 Huang, Q., Huang, P., Zhang, X., Xie, W., Hong, K., 675
630 Chen, B., & Huang, L. (2014). Chinese spelling 676
631 check system based on tri-gram model. 677
632 *Proceedings of the third CIPS-SIGHAN joint 678*
633 conference on Chinese language processing. 679
- 634 Yeh, J.-F., Chang, L.-T., Liu, C.-Y., & Hsu, T.-W. 680
635 (2017). Chinese spelling check based on N-gram 681
636 and string matching algorithm. *Proceedings of the 682*
637 4th Workshop on Natural Language Processing 683
638 Techniques for Educational Applications (NLPTEA 684
639 2017). 685
- 640 Duan, J., Wang, B., Tan, Z., Wei, X., & Wang, H. 686
641 (2019). Chinese spelling check via bidirectional 687
642 lstm-crf. 2019 IEEE 8th Joint International 688
643 Information Technology and Artificial Intelligence 689
644 Conference (ITAIC). 690
- 645 Wang, H., Wang, B., Duan, J., & Zhang, J. (2021). 691
646 Chinese Spelling Error Detection Using a Fusion 692
647 Lattice LSTM. *Transactions on Asian and Low- 693*
648 Resource Language Information Processing, 20(2), 694
649 1-11. 695
- 650 Wang, D., Tay, Y., & Zhong, L. (2019). Confusionset- 696
651 guided pointer networks for Chinese spelling check. 697
652 *Proceedings of the 57th Annual Meeting of the 698*
653 Association for Computational Linguistics. 699
- 654 Tan, M., Chen, D., Li, Z., & Wang, P. (2020). Spelling 700
655 Error Correction with BERT based on Character- 701
656 Phonetic. 2020 IEEE 6th International Conference 702
657 on Computer and Communications (ICCC). 703
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. 704
(2019). BERT: Pre-training of Deep Bidirectional 705
Transformers for Language Understanding. 706
NAACL-HLT (1).
- Cheng, X., Xu, W., Chen, K., Jiang, S., Wang, F., 707
Wang, T., Chu, W., & Qi, Y. (2020). SpellGCN: 708
Incorporating Phonological and Visual Similarities 709
into Language Models for Chinese Spelling Check. 710
Proceedings of the 58th Annual Meeting of the 711
Association for Computational Linguistics. 712
- Zhang, S., Huang, H., Liu, J., & Li, H. (2020). 713
Spelling Error Correction with Soft-Masked BERT. 714
Proceedings of the 58th Annual Meeting of the 715
Association for Computational Linguistics. 716
- Li, J., Yin, D., Wang, H., & Wang, Y. (2021). DCSpell: 717
A Detector-Corrector Framework for Chinese 718
Spelling Error Correction. *Proceedings of the 44th 719*
International ACM SIGIR Conference on Research 720
and Development in Information Retrieval. 721
- Clark, K., Luong, M.-T., Le, Q. V., & Manning, C. D. 722
(2016). ELECTRA: PRE-TRAINING TEXT 723
ENCODERS AS DISCRIMINATORS RATHER 724
THAN GENERATORS. *ELECTRA*, 85, 90. 725
- Hong, Y., Yu, X., He, N., Liu, N., & Liu, J. (2019). 726
Faspell: A fast, adaptable, simple, powerful chinese 727
spell checker based on dae-decoder paradigm. 728
Proceedings of the 5th Workshop on Noisy User- 729
generated Text (W-NUT 2019). 730
- Grandvalet, Y., & Bengio, Y. (2005). Semi-supervised 731
learning by entropy minimization. *CAP*, 367, 281- 732
296. 733
- Wu, S.-H., Liu, C.-L., & Lee, L.-H. (2013). Chinese 734
spelling check evaluation at SIGHAN bake-off 735
2013. *Proceedings of the Seventh SIGHAN 736*
Workshop on Chinese Language Processing. 737
- Yu, L.-C., Lee, L.-H., Tseng, Y.-H., & Chen, H.-H. 738
(2014). Overview of SIGHAN 2014 bake-off for 739
Chinese spelling check. *Proceedings of The Third 740*
CIPS-SIGHAN Joint Conference on Chinese 741
Language Processing. 742
- Wang, D., Song, Y., Li, J., Han, J., & Zhang, H. (2018, 743
oct) "nov). A Hybrid Approach to Automatic 744
Corpus Generation for Chinese Spelling 745
Check. *Proceedings of the 2018 Conference on 746*
Empirical Methods in Natural Language 747
Processing Brussels, Belgium. 748