

Exact Paired Permutation Testing Algorithms for NLP Systems

Anonymous ACL submission

Abstract

Significance testing has played a vital role in the development of NLP systems, providing confidence that one system is indeed better than another one. However, many significance tests involve hard computation problems, and so we rely on approximation methods such as Monte Carlo sampling. In this paper, we provide an exact dynamic programming algorithm that runs in quadratic time in the size of the dataset and performs the paired permutation test, a widely used test in comparing two systems, for the case of comparing accuracies between two classification systems. We show that Monte Carlo approximations are often too noisy to reliably determine whether we can reject the null hypothesis. We show that Monte Carlo approximations are often too noisy to reliably determine whether we can reject the null hypothesis with a significance level of $\alpha \approx 0.05$ for any number of sentence N . Additionally, we show that our exact algorithm is more efficient than the approximation algorithm for $N \leq 10K$.

1 Introduction

Statistical hypothesis testing (Lehmann and Romano, 2005) is a fundamental evaluation technique in the sciences. Thus, it should come as no surprise that statistical hypothesis testing is an often employed technique in natural language processing (e.g., Dietterich (1998); Koehn (2004); Ojala and Garriga (2010); Clark et al. (2011); Berg-Kirkpatrick et al. (2012)) for the comparison of competing system, e.g., can we claim that system A has lower error than system B with high confidence.

In this paper, we develop an efficient, exact algorithm for the paired permutation test (Good, 2000), a commonly used method for system comparison in NLP (Yeh, 2000; Dror et al., 2018, 2020; Deutsch et al., 2021). An exact paired permutation test involves a summation over a specific set of permutations. The naïve algorithm to perform this summation runs in exponential time as enumerates

all permutations. Thus, practitioners resort to running a Monte Carlo (MC) approximation to the permutation test. However, as with all stochastic simulation, this process introduces additional error when determining whether or not we may reject the null hypothesis. This paper introduces a new dynamic programming algorithm for the performing the computations required by a paired permutation test exactly. Furthermore, we show that the algorithm is efficient when comparing two systems based on accuracy, the standard metric for evaluation in many NLP tasks, e.g. part-of-speech tagging and dependency parsing. Furthermore, in the appendix (App. B), we give an exact quintic algorithm for comparing F_1 , but its runtime is too slow to be used on large datasets.

We compare the relative accuracy and efficiency of the exact test and the MC approximation for comparing part-of-speech taggers on the English Universal Dependency Dataset (Nivre et al., 2018). We experimented with the number of MC samples K and the number of sentences N . We find that in all settings, the estimation error of MC can be unreliable even when $\geq 20K$ samples are taken. Additionally, we find that our exact algorithm is faster than using MC with $20K$ and $40K$ samples for datasets with $\leq 6K$ and $\leq 10K$ sentences. Our Python implementation¹ runs the exact permutation test in under three seconds for a dataset with $10K$ sentences. Overall, our exact method appears to be a more practical alternative to MC.

2 Paired Permutation Test

A statistical hypothesis test attempts to reject a null hypothesis H_\emptyset at significance level α . It is common practice to set $\alpha = 0.05$, but there is a movement to lower the scientific standard (Ioannidis, 2018). We now turn to the paired permutation test (Good, 2000), the focus on this work. Suppose we want to

¹We will release our code publicly upon publication.

080 evaluate the performance of System A and System
081 B on an input \mathbf{x} with N entries that has a set of true
082 predictions \mathbf{y} . Suppose System A predicts $\hat{\mathbf{a}}$ and
083 System B predicts $\hat{\mathbf{b}}$. The type of the prediction
084 varies between tasks, e.g., $\mathbf{y} \in \{+1, -1\}^N$ in the
085 simplest case of positive and negative predictions
086 while $\mathbf{y} \in \{1, \dots, C\}^N$ for general classification
087 problems. In many NLP tasks, such as part-of-
088 speech tagging, dependency parsing, *inter alia*, a
089 sentence-level prediction may be decomposed into
090 word-level predictions. Then $\mathbf{y} \in \{1, \dots, C\}^{LN}$
091 where L is the maximum sentence length.

092 Now, suppose we wish to test the hypothesis
093 that $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ are different under a scoring function
094 f . We do this by performing a paired permutation
095 test on the null-hypothesis, H_\emptyset , that there is *no*
096 difference between $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ under f :

$$097 \quad H_\emptyset : \text{effect}_f(\hat{\mathbf{a}}, \hat{\mathbf{b}}) \stackrel{\text{def}}{=} |f(\hat{\mathbf{a}}) - f(\hat{\mathbf{b}})| \stackrel{?}{=} 0 \quad (1)$$

098 We attempt to reject H_\emptyset by considering the prob-
099 ability of having seen an effect size as small as
100 observed under the null distribution P_\emptyset .

101 **Definition 1.** Given two sets of output, $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$,
102 both of size N , a **paired permutation** (henceforth
103 **permutation**) $\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle$ is a pair of data that is
104 composed by swapping elements between $\hat{\mathbf{a}}$ and
105 $\hat{\mathbf{b}}$. Specifically, for entry $n \leq N$, we have that
106 $\hat{a}'_n = \hat{a}_n$ and $\hat{b}'_n = \hat{b}_n$ (no swap) or $\hat{a}'_n = \hat{b}_n$ and
107 $\hat{b}'_n = \hat{a}_n$ (swap). The set of all permutations of $\hat{\mathbf{a}}$
108 and $\hat{\mathbf{b}}$ is given by $S(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ and $|S(\hat{\mathbf{a}}, \hat{\mathbf{b}})| = 2^N$.

109 Under a paired permutation test, the null-
110 distribution, P_\emptyset , is defined to be the uniform
111 distribution over the paired permutations (i.e.,
112 $P_\emptyset(\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle) = 2^{-N}$). Then, if our observed
113 effect is $o = \text{effect}_f(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ and we have an
114 effect random variable drawn from P_\emptyset , $E =$
115 $\text{effect}_f(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \sim P_\emptyset$, we reject our null-
116 hypothesis if $p = \mathbb{P}(E \leq o) < \alpha$ where α is
117 pre-set. More formally, p can be computed as

$$118 \quad p = \frac{\sum_{\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle \in S(\hat{\mathbf{a}}, \hat{\mathbf{b}})} \mathbb{1}[\text{effect}_f(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq o]}{|S(\hat{\mathbf{a}}, \hat{\mathbf{b}})|} \quad (2)$$

119 The paired permutation test, thus, compares the
120 difference between each individual set of predic-
121 tions. Peyrard et al. (2021) show that this approach
122 is important for reliable significance testing in NLP.

```

1: def monte_carlo( $\hat{\mathbf{a}}, \hat{\mathbf{b}}$ ):
2:    $o \leftarrow \text{effect}(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ 
3:    $p \leftarrow 0$ 
4:   for  $i \in 1, \dots, K$ :
5:      $\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle \leftarrow \langle \mathbf{0}, \mathbf{0} \rangle$ 
6:     for  $n \in 1, \dots, N$ :
7:        $\langle \hat{a}'_n, \hat{b}'_n \rangle \leftarrow \text{RandomSwap}(\hat{a}_n, \hat{b}_n)$ 
8:        $p += \frac{1}{K} \mathbb{1}[\text{effect}(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq o]$ 
9:   return  $p$ 

```

Figure 1: Monte Carlo sampling approximation algo-
rithm for the paired permutation test.

Approximating Paired Permutation Test In
general, the paired permutation test requires us
to compute the sum in (2). The naïve computation
clearly runs in exponential time as it requires the
enumeration of all $|S(\hat{\mathbf{a}}, \hat{\mathbf{b}})| = 2^N$ paired per-
mutations. Therefore, most practical implemen-
tations of the paired permutation test use a MC
approximation, whereby one randomly samples
paired permutations to construct an approximate
null-distribution. We give this MC algorithm in
Fig. 1. Note that the algorithm works for *any* effect
function. In the following section, we provide an
exact algorithm for the case where effect measures
the absolute difference in accuracy.

3 Exact Test for Accuracy

We present a dynamic programming (DP) approach
for the paired permutation test. In the general case,
the runtime of our DP depends on the chosen effect.
However, in the case of accuracy, we are able to
derive an efficient exact algorithm.

$$143 \quad A(\hat{\mathbf{a}}) \stackrel{\text{def}}{=} \frac{t(\hat{\mathbf{a}})}{t(\hat{\mathbf{a}}) + f(\hat{\mathbf{a}})} = \frac{t(\hat{\mathbf{a}})}{N} \quad (3)$$

144 where $t(\hat{\mathbf{a}})$ and $f(\hat{\mathbf{a}})$ are the number of true predic-
145 tions and false predictions made in $\hat{\mathbf{a}}$ with regards
146 to \mathbf{y} , respectively. We can decompose $t(\hat{\mathbf{a}})$ among
147 each prediction such that $t(\hat{\mathbf{a}}) = \sum_{n=1}^N t(\hat{a}_n)$
148 where $t(\hat{a}_n) \in \{0, \dots, L\}$.

149 The aim of our test is to determine if the p -value
150

$$151 \quad p_A = \frac{\sum_{\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle \in S(\hat{\mathbf{a}}, \hat{\mathbf{b}})} \mathbb{1}[\text{effect}_A(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq o]}{|S(\hat{\mathbf{a}}, \hat{\mathbf{b}})|} \quad (4)$$

152 is less than the significance level α .

```

1: def perm_test_acc( $\hat{\mathbf{a}}, \hat{\mathbf{b}}$ ):
2:    $\mathbf{W} \leftarrow \mathbf{0}$ 
3:    $\mathbf{W}[0,0] \leftarrow 1$ 
4:   for  $n \in 1, \dots, N$ :
5:      $v_{\text{stay}} \leftarrow t(\hat{\mathbf{a}}'_n) - t(\hat{\mathbf{b}}'_n)$ 
6:      $v_{\text{swap}} \leftarrow t(\hat{\mathbf{b}}'_n) - t(\hat{\mathbf{a}}'_n)$ 
7:     for  $v \in \mathbf{W}[n]$ :
8:        $\mathbf{W}[n, v + v_{\text{stay}}] += \frac{1}{2} \mathbf{W}[n-1, v]$ 
9:        $\mathbf{W}[n, v + v_{\text{swap}}] += \frac{1}{2} \mathbf{W}[n-1, v]$ 
10:     $o \leftarrow |t(\hat{\mathbf{a}}) - t(\hat{\mathbf{b}})|$ 
11:     $p \leftarrow 0$ 
12:    for  $v \in \mathbf{W}[N]$ :
13:      if  $|v| \leq o$ :
14:         $p += \mathbf{W}[N, v]$ 
15:    return  $p$ 

```

Figure 2: Dynamic program to find exact p value for the paired-permutation test for accuracy.

Proposition 1. Given an input \mathbf{x} with predictions $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ and true predictions \mathbf{y} , for any paired permutation $\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle$, $\text{effect}_A(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq \text{effect}_A(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ iff $|t(\hat{\mathbf{a}}') - t(\hat{\mathbf{b}}')| \leq |t(\hat{\mathbf{a}}) - t(\hat{\mathbf{b}})|$

Proposition 1 indicates that we only need to care about the different in true predictions when performing the paired permutation test for accuracy. We build our DP by constructing a structure \mathbf{W} such that for any $n \in \{0, \dots, N\}$ and $v \in \mathcal{L}$ where $\mathcal{L} \stackrel{\text{def}}{=} \{-LN, \dots, LN\}$, $\mathbf{W}[n, l]$ is the probability that a paired permutation $\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle$ satisfies $t(\hat{\mathbf{a}}'_{:n}) - t(\hat{\mathbf{b}}'_{:n}) = v$ where we define $\hat{\mathbf{a}}'_{:n}$ to be the first n predictions of $\hat{\mathbf{a}}'$. Note that we do not consider the absolute value in the DP as we can not decompose an absolute difference into the difference of individual predictions. Once we have \mathbf{W} , then the row $\mathbf{W}[N]$ contains the distribution over $t(\hat{\mathbf{a}}') - t(\hat{\mathbf{b}}')$ and so we can find p_A . The algorithm is formalized as perm_test_acc in Fig. 2.²

Theorem 1. Given an input \mathbf{x} with predictions $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ and true predictions \mathbf{y} , perm_test_acc($\hat{\mathbf{a}}, \hat{\mathbf{b}}$) returns p_A in $\mathcal{O}(LN^2)$ time and $\mathcal{O}(LN)$ space.

²The correctness of perm_test_acc is given in Theorem 1. We provide proofs in in App. A.

Metric	Mean	Standard Dev.
Accuracy	0.9543	0.1116
Sentence length	12.08	10.60

Table 1: Distributions for of accuracy and sentence length for POS tagging using Stanza (Qi et al., 2020) on the English UD test dataset (Nivre et al., 2018).

3.1 Practical Implementation

Fig. 2 shows a \mathbf{W} structure that is a $\mathbb{R}^{N \times LN}$ matrix which suggests we need $\mathcal{O}(LN^2)$ space. However, we note that at any iteration n , we only ever need row $\mathbf{W}[n-1]$ and $\mathbf{W}[n]$. Therefore, we only need to maintain two rows of the matrix and so only require $\mathcal{O}(LN)$ space for the algorithm.

4 Experiments

We demonstrate the efficiency of our exact algorithms by simulating paired permutation tests between two systems. In order to have some control over the p -value, N , and L , we randomly generate our two system outputs from a measured distribution. Specifically, we will use the Stanza³ (Qi et al., 2020) part-of-speech tag statistics when evaluating on the English Universal Dependencies (UD) test set (Nivre et al., 2018). We sample our outputs from the normal distribution where the mean and standard deviation match the rates of Stanza. We further sample the length of each sample sentence according to the distribution of lengths in the test set. The distributions are provided in Tab. 1.

Error Rate of Monte Carlo. We first examine the error rate of the monte_carlo as we increase the number of samples used. We sample $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ using the distributions in Tab. 1, however, we multiply the mean of $\hat{\mathbf{b}}$ by a factoring in proportion to N . We do this to obtain p -values roughly between 0.001 and 0.1 which is a typical range for α and so a paired permutation test would be required. For each K , we sampled five pairs of systems and run the method five times for each pair, giving 25 data points for each K and N . The reported error rates given in Fig. 3 are the averages of these errors.⁴ The mean multiplication factor as well as the average p -value for each N is given in Tab. 2. We see

³The code and pre-trained model are both freely accessible at <https://github.com/stanfordnlp/stanza>.

⁴We discard relative errors > 10 . We do this to be able to see the noise more clearly in a specific range. However, we note that requiring this additional filter further shows the unreliability of the MC method.

N	Mean	Mean p	Standard Dev. p
250	0.95μ	0.0512	0.1351
500	0.96μ	0.0189	0.0656
1000	0.97μ	0.0153	0.0534
2000	0.98μ	0.0202	0.0536
4000	0.985μ	0.0170	0.0538
8000	0.99μ	0.0360	0.1146

Table 2: Means used for $\hat{\mathbf{b}}$ distributions in Fig. 3. We use μ to reference the mean 0.9543 given in Tab. 1. The mean and standard deviation of the p -values of the paired permutation test are also given.

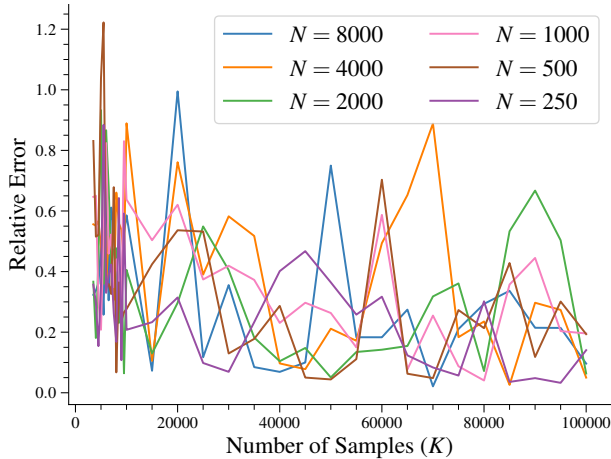


Figure 3: Relative errors of using monte_carlo for the paired permutation test. System $\hat{\mathbf{a}}$ is sampled according to Tab. 1 and system $\hat{\mathbf{b}}$ is sampled according to Tab. 2.

that the MC approximations are not reliable for all the values of N . While there is a downwards trend as we increase the number of cycles, we observe a lot of noise even when taking 25 attempts per N and K pair. The trend seems most clear until $K = 20,000$ at which point we see a lot of noise as K increases. We therefore suggest that 20,000 is the minimum number of samples required when performing a MC paired permutation test, though more is likely better.

Advantages of the Exact Test. When reporting system accuracies in the literature, an exact p -value avoids the estimation error associated with Monte Carlo, as the results above demonstrate. We now show that, empirically, the exact test is *more efficient* than the MC approximation when a large number of samples is taken; this is evinced in Fig. 4. We compare the runtime of perm_test_acc against monte_carlo for $K = 20,000$ and $K = 40,000$.⁵

⁵The experiment used an Apple M1 Max processor.

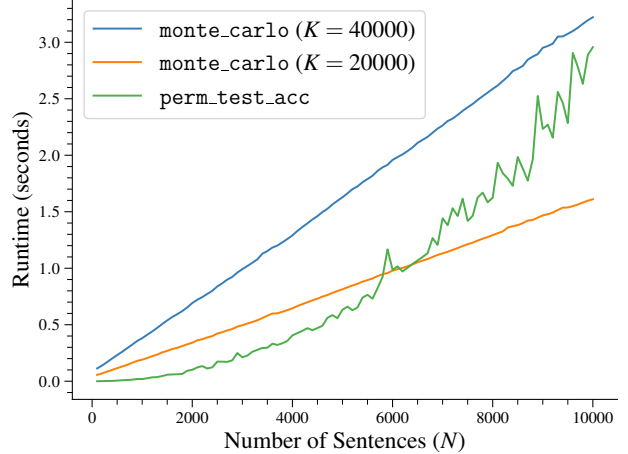


Figure 4: Runtime comparison of perm_test_acc and monte_carlo as a function of the number of sentences.

We can see that perm_test_acc is more efficient than monte_carlo with $K = 40,000$ and $K = 20,000$ for $N < 10,000$ and $N < 6,000$ (respectively). We note that the average test set size of the UD treebanks⁶ is just over 1,000 sentences, and only three treebanks had more than 6,000 sentences.⁷ Additionally, the standard split of the commonly used Penn treebank (PTB) (Marcus et al., 1993) provides a test set of about 5,500 sentences. Therefore, the perm_test_acc is more efficient than monte_carlo for most of the datasets that are used in NLP token-level classification problems.

5 Conclusion

We presented a dynamic programming algorithm to compute the exact p -value of a paired permutation test for the case of difference in accuracy. Our algorithm runs in $\mathcal{O}(LN^2)$ time and requires $\mathcal{O}(LN)$ space. We empirically show that when using MC approximation techniques, we often require $K > LN$ samples to obtain a “good enough” approximation. Therefore, not only is the MC method imprecise, it is also often slower than our exact algorithm for commonly used datasets. We also note that our dynamic program can be extended to compute exact p -values for the paired permutation test using other metrics such as the difference in F_1 scores (see App. B). However, these may be impractical for reasonably sized dataset.

⁶We examined a total of 129 treebanks as some languages have multiple treebanks.

⁷These were Czech, Japanese, and Russian which had test set sizes of roughly 10,000, 8,000, and 6,500 (respectively).

259
260
261
262
263
264
265
266

267
268
269
270
271
272
273
274

275
276
277
278
279

280
281
282

283
284
285
286
287
288
289

290
291
292
293
294

295
296
297

298
299
300

301
302
303
304
305

306
307

308
309
310
311

References

Taylor Berg-Kirkpatrick, David Burkett, and Dan Klein. 2012. [An empirical investigation of statistical significance in NLP](#). In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, EMNLP-CoNLL 2012, July 12-14, 2012, Jeju Island, Korea*, pages 995–1005. ACL.

Jonathan H. Clark, Chris Dyer, Alon Lavie, and Noah A. Smith. 2011. [Better hypothesis testing for statistical machine translation: Controlling for optimizer instability](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 176–181, Portland, Oregon, USA. Association for Computational Linguistics.

Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. [A statistical analysis of summarization evaluation metrics using resampling methods](#). *Transactions of the Association for Computational Linguistics*, 9:1132–1146.

Thomas G. Dietterich. 1998. [Approximate statistical tests for comparing supervised classification learning algorithms](#). *Neural computation*, 10(7):1895–1923.

Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. [The hitchhiker’s guide to testing statistical significance in natural language processing](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1383–1392, Melbourne, Australia. Association for Computational Linguistics.

Rotem Dror, Lotem Peled-Cohen, Segev Shlomov, and Roi Reichart. 2020. [Statistical Significance Testing for Natural Language Processing](#). Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers.

Phillip Good. 2000. [Permutation Tests A Practical Guide to Resampling Methods for Testing Hypotheses](#). Springer.

John P. A. Ioannidis. 2018. [The Proposal to Lower P Value Thresholds to .005](#). *JAMA*, 319(14):1429–1430.

Philipp Koehn. 2004. [Statistical significance tests for machine translation evaluation](#). In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.

Erich Leo Lehmann and Joseph P. Romano. 2005. [Testing Statistical Hypotheses](#). Springer.

Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. [Building a large annotated corpus of English: The Penn Treebank](#). *Comput. Linguistics*, 19(2):313–330.

Joakim Nivre, Mitchell Abrams, Željko Agić, Lars Ahrenberg, Lene Antonsen, Katya Aplonova, Maria Jesus Aranzabe, Gashaw Arutie, Masayuki Asahara, Luma Ateyah, Mohammed Attia, Aitziber Atutxa, Liesbeth Augustinus, Elena Badmaeva, Miguel Ballesteros, Esha Banerjee, Sebastian Bank, Verginica Barbu Mititelu, Victoria Basmov, John Bauer, Sandra Bellato, Kepa Bengoetxea, Yevgeni Berzak, Irshad Ahmad Bhat, Riyaz Ahmad Bhat, Erica Biagetti, Eckhard Bick, Rogier Blokland, Victoria Bobicev, Carl Börstell, Cristina Bosco, Gosse Bouma, Sam Bowman, Adriane Boyd, Aljoscha Burchardt, Marie Candito, Bernard Caron, Gauthier Caron, Gülşen Cebiroğlu Eryiğit, Flavio Massimiliano Cecchini, Giuseppe G. A. Celano, Slavomír Čěplö, Savas Cetin, Fabricio Chalub, Jinho Choi, Yongseok Cho, Jayeol Chun, Silvie Cinková, Aurélie Collomb, Çağrı Çöltekin, Miriam Connor, Marine Courtin, Elizabeth Davidson, Marie-Catherine de Marneffe, Valeria de Paiva, Arantza Diaz de Ilarraza, Carly Dickerson, Peter Dirix, Kaja Dobrovoljc, Timothy Dozat, Kira Droganova, Puneet Dwivedi, Marhaba Eli, Ali Elkahky, Binyam Ephrem, Tomaž Erjavec, Aline Etienne, Richárd Farkas, Hector Fernandez Alcalde, Jennifer Foster, Cláudia Freitas, Katarína Gajdošová, Daniel Galbraith, Marcos Garcia, Moa Gärdenfors, Sebastian Garza, Kim Gerdes, Filip Ginter, Iakes Goenaga, Koldo Gojenola, Memduh Gökırmak, Yoav Goldberg, Xavier Gómez Guinovart, Berta González Saavedra, Matias Grioni, Normunds Grūzītis, Bruno Guillaume, Céline Guillot-Barbance, Nizar Habash, Jan Hajič, Jan Hajič jr., Linh Hà Mý, Na-Rae Han, Kim Harris, Dag Haug, Barbora Hladká, Jaroslava Hlaváčová, Florinel Hociung, Petter Hohle, Jena Hwang, Radu Ion, Elena Irimia, Olájídé Ishola, Tomáš Jelínek, Anders Johannsen, Fredrik Jørgensen, Hüner Kaşıkara, Sylvain Kahane, Hiroshi Kanayama, Jenna Kanerva, Boris Katz, Tolga Kayadelen, Jessica Kenney, Václava Kettnerová, Jesse Kirchner, Kamil Kopacewicz, Natalia Kotsyba, Simon Krek, Sookyong Kwak, Veronika Laippala, Lorenzo Lambertino, Lucia Lam, Tatiana Lando, Septina Dian Larasati, Alexei Lavrentiev, John Lee, Phuong Lê Hồng, Alessandro Lenci, Saran Lertpradit, Herman Leung, Cheuk Ying Li, Josie Li, Keying Li, Kyung-Tae Lim, Nikola Ljubešić, Olga Loginova, Olga Lyahevskaya, Teresa Lynn, Vivien Macketanz, Aibek Makazhanov, Michael Mandl, Christopher Manning, Ruli Manurung, Cătălina Mărănduc, David Mareček, Katrin Marheinecke, Héctor Martínez Alonso, André Martins, Jan Mašek, Yuji Matsumoto, Ryan McDonald, Gustavo Mendonça, Niko Miekka, Margarita Misirpashayeva, Anna Missilä, Cătălin Mititelu, Yusuke Miyao, Simonetta Montemagni, Amir More, Laura Moreno Romero, Keiko Sophie Mori, Shinsuke Mori, Bjartur Mortensen, Bohdan Moskalevskyi, Kadri Muischnek, Yugo Murawaki, Kaili Müürisep, Pinkey Nainwani, Juan Ignacio Navarro Horñiacek, Anna Nedoluzhko, Gunta Nešpore-Bērzkalne, Luong Nguyễn Thị, Huyền Nguyễn Thị Minh, Vitaly Nikolaev, Rattima Nitisaroj, Hanna Nurmi, Stina Ojala, Adedayo Olúokun,

375 Mai Omura, Petya Osenova, Robert Östling, Lilja
376 Øvrelid, Niko Partanen, Elena Pascual, Marco
377 Passarotti, Agnieszka Patejuk, Guilherme Paulino-
378 Passos, Siyao Peng, Cemel-Augusto Perez, Guy Per-
379 rier, Slav Petrov, Jussi Piitulainen, Emily Pitler, Bar-
380 bara Plank, Thierry Poibeau, Martin Popel, Lauma
381 Pretkalniņa, Sophie Prévost, Prokopis Prokopidis,
382 Adam Przepiórkowski, Tiina Puolakainen, Sampo
383 Pyysalo, Andriela Rääbis, Alexandre Rademaker, Lo-
384 ganathan Ramasamy, Taraka Rama, Carlos Ramisch,
385 Vinit Ravishankar, Livy Real, Siva Reddy, Georg
386 Rehm, Michael Rießler, Larissa Rinaldi, Laura Rit-
387 uma, Luisa Rocha, Mykhailo Romanenko, Rudolf
388 Rosa, Davide Rovati, Valentin Roşca, Olga Rud-
389 ina, Jack Rueter, Shoval Sadde, Benoît Sagot, Shadi
390 Saleh, Tanja Samardžić, Stephanie Samson, Manuela
391 Sanguinetti, Baiba Saulīte, Yanin Sawanakunanon,
392 Nathan Schneider, Sebastian Schuster, Djamel Sed-
393 dah, Wolfgang Seeker, Mojgan Seraji, Mo Shen,
394 Atsuko Shimada, Muh Shohibussirri, Dmitry Sichi-
395 nava, Natalia Silveira, Maria Simi, Radu Simionescu,
396 Katalin Simkó, Mária Šimková, Kiril Simov, Aaron
397 Smith, Isabela Soares-Bastos, Carolyn Spadine, An-
398 tonio Stella, Milan Straka, Jana Strnadová, Alane
399 Suhr, Umut Sulubacak, Zsolt Szántó, Dima Taji,
400 Yuta Takahashi, Takaaki Tanaka, Isabelle Tellier,
401 Trond Trosterud, Anna Trukhina, Reut Tsarfaty, Fran-
402 cis Tyers, Sumire Uematsu, Zdeňka Urešová, Lar-
403 raitz Uria, Hans Uszkoreit, Sowmya Vajjala, Daniel
404 van Niekerk, Gertjan van Noord, Viktor Varga, Eric
405 Villemonte de la Clergerie, Veronika Vincze, Lars
406 Wallin, Jing Xian Wang, Jonathan North Washing-
407 ton, Seyi Williams, Mats Wirén, Tsegay Wolde-
408 mariam, Tak-sum Wong, Chunxiao Yan, Marat M.
409 Yavrumyan, Zhuoran Yu, Zdeněk Žabokrtský, Amir
410 Zeldes, Daniel Zeman, Manying Zhang, and Hanzhi
411 Zhu. 2018. [Universal dependencies 2.3](#). LIN-
412 DAT/CLARIN digital library at the Institute of For-
413 mal and Applied Linguistics (ÚFAL), Faculty of
414 Mathematics and Physics, Charles University.

415 Markus Ojala and Gemma C. Garriga. 2010. [Permuta-](#)
416 [tion tests for studying classifier performance](#). *The*
417 *Journal of Machine Learning Research*, 11:1833–
418 1863.

419 Maxime Peyrard, Wei Zhao, Steffen Eger, and Robert
420 West. 2021. [Better than average: Paired evaluation](#)
421 [of NLP systems](#). In *Proceedings of the 59th Annual*
422 *Meeting of the Association for Computational Lin-*
423 *guistics and the 11th International Joint Conference*
424 *on Natural Language Processing, ACL/IJCNLP 2021,*
425 *(Volume 1: Long Papers), Virtual Event, August 1-6,*
426 *2021*, pages 2301–2315. Association for Computa-
427 tional Linguistics.

428 Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and
429 Christopher D. Manning. 2020. [Stanza: A Python](#)
430 [natural language processing toolkit for many human](#)
431 [languages](#). In *Proceedings of the Association for*
432 *Computational Linguistics: System Demonstrations*.

433 Alexander Yeh. 2000. [More accurate tests for the statis-](#)
434 [tical significance of result differences](#). In *COLING*

2000 Volume 2: The 18th International Conference
on Computational Linguistics.

435
436

A Proofs for Section §3 (Exact Test for Accuracy)

Proposition 1. Given an input \mathbf{x} with predictions $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ and true predictions \mathbf{y} , for any paired permutation $\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle$, $\text{effect}_A(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq \text{effect}_A(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ iff $|t(\hat{\mathbf{a}}') - t(\hat{\mathbf{b}}')| \leq |t(\hat{\mathbf{a}}) - t(\hat{\mathbf{b}})|$

Proof.

$$\begin{aligned} \text{effect}_A(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq \text{effect}_A(\hat{\mathbf{a}}, \hat{\mathbf{b}}) &\iff |A(\hat{\mathbf{a}}') - A(\hat{\mathbf{b}}')| \leq |A(\hat{\mathbf{a}}) - A(\hat{\mathbf{b}})| \\ &\iff \left| \frac{t(\hat{\mathbf{a}}') - t(\hat{\mathbf{b}}')}{|\mathbf{y}|} \right| \leq \left| \frac{t(\hat{\mathbf{a}}) - t(\hat{\mathbf{b}})}{|\mathbf{y}|} \right| \\ &\iff |t(\hat{\mathbf{a}}') - t(\hat{\mathbf{b}}')| \leq |t(\hat{\mathbf{a}}) - t(\hat{\mathbf{b}})| \end{aligned}$$

Theorem 1. Given an input \mathbf{x} with predictions $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ and true predictions \mathbf{y} , $\text{perm_test_acc}(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ returns p_A in $\mathcal{O}(LN^2)$ time and $\mathcal{O}(LN)$ space.

Proof. We first prove that for all $n \in \{0, \dots, N\}$, $\mathbf{W}[n]$ is the probability distribution $\mathbb{P}(\text{effect}_A(\hat{\mathbf{a}}'_{:n}, \hat{\mathbf{b}}'_{:n}))$.

Base case: Then $n = 0$. We have that $\mathbf{W}[0, 0] = 1$ and $\mathbf{W}[0, v] = 0$ for all $v \in \mathcal{L} \setminus \{0\}$.

Inductive step: Assume that $\mathbf{W}[n-1]$ is the probability distribution $\mathbb{P}(\text{effect}_A(\hat{\mathbf{a}}'_{:(n-1)}, \hat{\mathbf{b}}'_{:(n-1)}))$. Let $v \in \mathcal{L}$ be a candidate difference and $v' = t(\hat{a}'_n) - t(\hat{b}'_n)$. We know that $\langle \hat{a}'_n, \hat{b}'_n \rangle$ is $\langle \hat{a}_n, \hat{b}_n \rangle$ with probability $\frac{1}{2}$ or $\langle \hat{b}_n, \hat{a}_n \rangle$ with probability $\frac{1}{2}$. Therefore, $v' = t(\hat{a}_n) - t(\hat{b}_n)$ with probability $\frac{1}{2}$ or $v' = t(\hat{b}_n) - t(\hat{a}_n)$ with probability $\frac{1}{2}$. Then

$$\mathbb{P}(\text{effect}_A(\hat{\mathbf{a}}'_{:n}, \hat{\mathbf{b}}'_{:n}) = v) = \frac{1}{2} \mathbb{P}(\text{effect}_A(\hat{\mathbf{a}}'_{:(n-1)}, \hat{\mathbf{b}}'_{:(n-1)}) = v - v') = \frac{1}{2} \mathbf{W}[n-1, v - v']$$

This is exactly what is done from Line 7 to Line 9 in Fig. 2. Therefore, $\mathbf{W}[n]$ is the probability distribution $\mathbb{P}(\text{effect}_A(\hat{\mathbf{a}}'_{:n}, \hat{\mathbf{b}}'_{:n}))$.

Line 11 to Line 14 in Fig. 2 construct the p -value using the following equation

$$\begin{aligned} \sum_{v \in \mathbf{W}[N]} \mathbb{1}[|v| \leq o] \mathbf{W}[N, v] &= \sum_{v \in \mathbf{W}[N]} \mathbb{1}[|v| \leq o] \mathbb{P}(\text{effect}_A(\hat{\mathbf{a}}', \hat{\mathbf{b}}') = v) \\ &= \mathbb{P}(\text{effect}_A(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq \text{effect}_A(\hat{\mathbf{a}}, \hat{\mathbf{b}})) = p_A \end{aligned}$$

where $o = |t(\hat{\mathbf{a}}) - t(\hat{\mathbf{b}})|$

The algorithm runs over two nested for-loops of sizes $\mathcal{O}(N)$ and $\mathcal{O}(LN)$ respectively. As the inner loop does constant amount of work per iteration, perm_test_acc runs in $\mathcal{O}(LN^2)$ time. The space complexity is discussed in §3.1. ■

B Exact Paired Permutation Test for F_1

We now derive a similar DP algorithm for the case of the F_1 score which we define as

$$F_1(\mathbf{x}) = \frac{t^+(\hat{\mathbf{a}})}{t^+(\hat{\mathbf{a}}) + \frac{1}{2}f(\hat{\mathbf{a}})} \quad (5)$$

where $t^+(\hat{\mathbf{a}})$ is the number of true positive predictions made in $\hat{\mathbf{a}}$ with regards to \mathbf{y} .

```

1: def perm_test_F1( $\hat{\mathbf{a}}, \hat{\mathbf{b}}, K$ ):
2:    $\mathbf{W} \leftarrow \mathbf{0}$ 
3:    $\mathbf{W}[0, \langle 0, 0, 0, 0 \rangle] \leftarrow 1$ 
4:   for  $n \in 1, \dots, N$ :
5:     for  $\langle t_a^+, f_a, t_b^+, f_b \rangle \in \mathbf{W}[n]$ :
6:        $v_{\text{stay}} \leftarrow \langle t_a^+ + t^+(\hat{\mathbf{a}}'_n), f_a + f(\hat{\mathbf{a}}'_n), t_b^+ + t^+(\hat{\mathbf{b}}'_n), f_b + f(\hat{\mathbf{b}}'_n) \rangle$ 
7:        $v_{\text{swap}} \leftarrow \langle t_a^+ + t^+(\hat{\mathbf{b}}'_n), f_a + f(\hat{\mathbf{b}}'_n), t_b^+ + t^+(\hat{\mathbf{a}}'_n), f_b + f(\hat{\mathbf{a}}'_n) \rangle$ 
8:        $\mathbf{W}[n, v_{\text{stay}}] += \frac{1}{2} \mathbf{W}[n-1, \langle t_a^+, f_a, t_b^+, f_b \rangle]$ 
9:        $\mathbf{W}[n, v_{\text{swap}}] += \frac{1}{2} \mathbf{W}[n-1, \langle t_a^+, f_a, t_b^+, f_b \rangle]$ 
10:     $o \leftarrow \left| \frac{t^+(\hat{\mathbf{a}})}{t^+(\hat{\mathbf{a}}) + \frac{1}{2}f(\hat{\mathbf{a}})} - \frac{t^+(\hat{\mathbf{b}})}{t^+(\hat{\mathbf{b}}) + \frac{1}{2}f(\hat{\mathbf{b}})} \right|$ 
11:     $p \leftarrow 0$ 
12:    for  $\langle t_a^+, f_a, t_b^+, f_b \rangle \in \mathbf{W}[N]$ :
13:      if  $\left| \frac{t_a^+}{t_a^+ + \frac{1}{2}f_a} - \frac{t_b^+}{t_b^+ + \frac{1}{2}f_b} \right| \leq o$ :
14:         $p += \mathbf{W}[N, \langle t_a^+, f_a, t_b^+, f_b \rangle]$ 
15:    return  $p$ 

```

Figure 5: Dynamic program to find exact p value for the paired-permutation test for F_1 .

The aim of our significance test is to decide whether the p -value

$$p_{F_1} = \frac{\sum_{\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle \in S(\hat{\mathbf{a}}, \hat{\mathbf{b}})} \mathbb{1}[\text{effect}_{F_1}(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq o]}{|S(\hat{\mathbf{a}}, \hat{\mathbf{b}})|} \quad (6)$$

is less than the significance level α where $\text{effect}_{F_1}(\hat{\mathbf{a}}, \hat{\mathbf{b}}) \stackrel{\text{def}}{=} |F_1(\hat{\mathbf{a}}) - F_1(\hat{\mathbf{b}})|$. Unfortunately, unlike accuracy, we cannot decompose the F_1 score into a single additive component. We can write $\text{effect}_{F_1}(\hat{\mathbf{a}}, \hat{\mathbf{b}})$ as

$$\text{effect}_{F_1}(\hat{\mathbf{a}}, \hat{\mathbf{b}}) = \left| \frac{t^+(\hat{\mathbf{a}})}{t^+(\hat{\mathbf{a}}) + \frac{1}{2}f(\hat{\mathbf{a}})} - \frac{t^+(\hat{\mathbf{b}})}{t^+(\hat{\mathbf{b}}) + \frac{1}{2}f(\hat{\mathbf{b}})} \right|$$

Therefore, we have four variables that we can decompose along the data points, $t^+(\hat{\mathbf{a}})$, $f(\hat{\mathbf{a}})$, $t^+(\hat{\mathbf{b}})$, and $f(\hat{\mathbf{b}})$. We construct a similar DP to `perm_test_acc`, however instead of maintaining the difference in true predictions, we maintain a tuple of the four aforementioned variables. We give this algorithm as `perm_test_F1` in Fig. 5 As each variable can be any of $\mathcal{O}(LN)$ values, this makes our DP have a runtime of $\mathcal{O}(L^4 N^5)$. Unfortunately, while the algorithm is polynomial in time, the quintic factor makes it impractical for common NLP datasets as described in §4

Theorem 2. *Given an input \mathbf{x} with predictions $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}$ and true predictions \mathbf{y} , `perm_test_F1`($\hat{\mathbf{a}}, \hat{\mathbf{b}}$) returns p_{F_1} in $\mathcal{O}(L^4 N^5)$ time and $\mathcal{O}(L^4 N^4)$ space.*

Proof. For any $n \in \{0, \dots, N\}$ and $\langle t_a^+, f_a, t_b^+, f_b \rangle \in \mathcal{L}^4$, we define $E_n(t_a^+, f_a, t_b^+, f_b)$ to be the event that $t(\hat{\mathbf{a}}'_n) = t_a^+$, $f(\hat{\mathbf{a}}'_n) = f_a$, $t(\hat{\mathbf{b}}'_n) = t_b^+$, and $f(\hat{\mathbf{b}}'_n) = f_b$. We first prove that $\mathbf{W}[n]$ is the probability distribution over the tuples such that

$$\mathbf{W}[n, \langle t_a^+, f_a, t_b^+, f_b \rangle] = \mathbb{P}(E_n(t_a^+, f_a, t_b^+, f_b)) \quad (7)$$

Base case: Then $n = 0$. We have that $\mathbf{W}[0, \langle 0, 0, 0, 0 \rangle] = 1$ and $\mathbf{W}[0, v] = 0$ for all $v \in \mathcal{L} \setminus \{\langle 0, 0, 0, 0 \rangle\}$.

Inductive step: Assume that $\mathbf{W}[n-1]$ is the probability distribution described in (7). We know that $\langle \hat{a}'_n, \hat{b}'_n \rangle$ is $\langle \hat{a}_n, \hat{b}_n \rangle$ with probability $\frac{1}{2}$ or $\langle \hat{b}_n, \hat{a}_n \rangle$ with probability $\frac{1}{2}$. Then, if we let $\langle t_a^+, f_a, t_b^+, f_b \rangle \in \mathcal{L}^4$, we can find the following probability

$$\begin{aligned} \mathbb{P}(E_n(t_a^+, f_a, t_b^+, f_b)) &= \frac{1}{2} \mathbb{P}\left(E_{n-1}\left(t_a^+ - t(\hat{a}'_n), f_a - f(\hat{a}'_n), t_b^+ - t(\hat{b}'_n), f_b - f(\hat{b}'_n)\right)\right) \\ &= \frac{1}{2} \mathbf{W}\left[n-1, \langle t_a^+ - t(\hat{a}'_n), f_a - f(\hat{a}'_n), t_b^+ - t(\hat{b}'_n), f_b - f(\hat{b}'_n) \rangle\right] \end{aligned}$$

This is exactly what is done from Line 8 to Line 9 in Fig. 5. Therefore, $\mathbf{W}[n]$ is the probability distribution $\mathbb{P}(E_n(t_a^+, f_a, t_b^+, f_b))$.

Line 11 to Line 14 in Fig. 5 construct the p -value using the following equation

$$\sum_{\langle t_a^+, f_a, t_b^+, f_b \rangle \in \mathbf{W}[N]} \mathbb{1}\left[\left|\frac{t_a^+}{t_a^+ + \frac{1}{2}f_a} - \frac{t_b^+}{t_b^+ + \frac{1}{2}f_b}\right| \leq o\right] \mathbf{W}[N, \langle t_a^+, f_a, t_b^+, f_b \rangle] \quad (8)$$

$$= \sum_{\langle t_a^+, f_a, t_b^+, f_b \rangle \in \mathbf{W}[N]} \mathbb{1}\left[\left|\frac{t_a^+}{t_a^+ + \frac{1}{2}f_a} - \frac{t_b^+}{t_b^+ + \frac{1}{2}f_b}\right| \leq o\right] \mathbb{P}(E_N(t_a^+, f_a, t_b^+, f_b)) \quad (9)$$

$$= \sum_{\langle \hat{\mathbf{a}}', \hat{\mathbf{b}}' \rangle \in S(\hat{\mathbf{a}}, \hat{\mathbf{b}})} \mathbb{1}\left[\text{effect}_{F_1}(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq o\right] \mathbb{P}\left(E_N\left(t(\hat{\mathbf{a}}'), f(\hat{\mathbf{a}}'), t(\hat{\mathbf{b}}'), f(\hat{\mathbf{b}}')\right)\right) \quad (10)$$

$$= \mathbb{P}\left(\text{effect}_{F_1}(\hat{\mathbf{a}}', \hat{\mathbf{b}}') \leq o\right) = p_{F_1} \quad (11)$$

The algorithm runs over two nested for-loops of sizes $\mathcal{O}(N)$ and $\mathcal{O}(L^4 N^4)$ respectively. As the inner loop does constant amount of work per iteration, `perm_test_F1` runs in $\mathcal{O}(L^4 N^5)$ time. We need to store two rows of \mathbf{W} at any given time. Therefore the space complexity is $\mathcal{O}(L^4 N^4)$ ■