

Overview of the PAN@FIRE 2020 Task on the Authorship Identification of SOURCE CODE

Ali Fadel^a, Husam Musleh^a, Ibraheem Tuffaha^a, Mahmoud Al-Ayyoub^a, Yaser Jararweh^a, Elhadj Benkhelifa^b and Paolo Rosso^c

^aJordan University of Science and Technology, Jordan

^bStaffordshire University, UK

^cUniversitat Politècnica de València, Spain

Abstract

Authorship identification is essential to the detection of undesirable deception of others' content misuse or exposing the owners of some anonymous malicious content. While it is widely studied for natural languages, it is rarely considered for programming languages. Accordingly, a PAN@FIRE task, named Authorship Identification of SOURCE CODE (AI-SOCO), is proposed with the focus on the identification of source code authors. The dataset consists of crawled source codes submitted by the top 1,000 human users with 100 correct C++ submissions or more from the CodeForces online judge platform. The participating systems are asked to predict the author of a given source code from the predefined list of code authors. In total, 60 teams registered on the task's CodaLab page. Out of them, 14 teams submitted 94 runs. The results are surprisingly high with many teams and baselines breaking the 90% accuracy barrier. These systems used a wide range of models and techniques from pretrained word embeddings (especially, those that are tweaked to handle source code) to stylometric features.

Keywords

authorship-identification, source-code, datasets

1. Introduction

After the success of the previous tasks on source code such as SOURCE CODE re-use (SOCO) [1], Cross-Language SOURCE CODE re-use (CL-SOCO) [2], and Personality Recognition in SOURCE CODE (PR-SOCO) [3], a new task is proposed in this paper. Specifically, we describe the Authorship Identification of SOURCE CODE (AI-SOCO)¹ task, one of the tracks of the 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020).² In the following sections, we define the task, discuss the dataset we introduced to the research community, present the available baselines and results, and highlight the achieved accomplishments.

FIRE'20: Forum for Information Retrieval Evaluation, December 16–20, 2020, Hyderabad, India

✉ aliosm1997@gmail.com (A. Fadel); husam.sa3@gmail.com (H. Musleh); bro.t.1996@gmail.com (I. Tuffaha); malayyoub@gmail.com (M. Al-Ayyoub); yijararweh@just.edu.jo (Y. Jararweh); E.Benkhelifa@staffs.ac.uk (E. Benkhelifa); proso@dsic.upv.es (P. Rosso)

🌐 <https://github.com/AliOsm/> (A. Fadel)

🆔 0000-0002-3447-8331 (A. Fadel); 0000-0001-9372-9076 (M. Al-Ayyoub)

© 2020 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

CEUR Workshop Proceedings (CEUR-WS.org)

¹<https://sites.google.com/view/ai-soco-2020>

²<http://fire.irsi.res.in/fire/2020>

A. Word Capitalization

time limit per test: 2 seconds
memory limit per test: 256 megabytes
input: standard input
output: standard output

Capitalization is writing a word with its first letter as a capital letter. Your task is to capitalize the given word.

Note, that during capitalization all the letters except the first one remains unchanged.

Input

A single line contains a non-empty word. This word consists of lowercase and uppercase English letters. The length of the word will not exceed 10^3 .

Output

Output the given word after capitalization.

Figure 1: CodeForces Problem Sample.

2. Task Definition

AI-SOCO focuses on identifying the author who wrote a given piece of code out of a predefined set of authors. This facilitates solving issues related to cheating in academic, work and open source environments. The detection of cheating in academic communities is significant to properly assess the contribution of students or researchers. Moreover, in work environments, credit sometimes goes to people who did not deserve it. Such issues of plagiarism could arise in open source projects that are available on public platforms. Similarly, systems developed for the AI-SOCO task can be used in public or private online coding contests, whether done in coding interviews or in official coding training contests, to detect the cheating of applicants or contestants. Such systems can also play a significant role in detecting the source of anonymous malicious software.

3. Dataset and Evaluation

The dataset is composed of source codes collected from the open submissions in the CodeForces³ platform. CodeForces is an online judge for hosting competitive programming contests, where each contest consists of multiple problems. The problems' difficulty (aka index) ranges from A to Z, where A is the easiest and Z is the hardest. The majority of contests have 5 to 7 problems (A-G). So, the occurrence of problems with indexes bigger than G is very rare. Figure 1 shows an example of a CodeForces problem, while Figure 2 shows a sample of the possible solutions by three different CodeForces users.

A user can solve a problem by writing a solution for it using any of the available programming languages on the website, and then submitting the solution through the website. The solution is tested automatically against a set of test cases and its result can be correct (accepted) or incorrect (wrong answer, time limit exceeded, etc.).

For our dataset, we selected the top rated 1,000 users (according to the website's rating system) and collected 100 source codes per user. This gives a total number of 100,000 source codes. All collected source codes are correct and written using the C++ programming language. For each user, all collected source codes are from unique problems and were submitted within the

³<https://codeforces.com>

```

#include <string>
#include <iostream>
#include <ctype.h>
using namespace std;

int main() {
    string s;
    cin >> s;
    s[0] = toupper(s[0]);
    cout << s << endl;
    return 0;
}

#include <iostream>
#include <ctype>

using namespace std;

int main(){
    string s;
    cin >> s;
    cout << (char)toupper(s[0]);
    for(int i=1; i<s.size(); i++){
        cout << s[i];
    }
    return 0;
}

#include<iostream>
#include<algorithm>
#include<string>
#include<time.h>
using namespace std;
bool cmp(int a, int b)
{
    return a<b;
}
int main(){
int a[1001]={0},i,j,k=0,l,m,n,f=0;
string s,s1;
char c;
getline(cin,s1);
if(s1[0]<='z' && s1[0]>='a')
s1[0]='A'+(s1[0]-'a');
cout<<s1;

//system("pause");
return 0;
}

```

Figure 2: CodeForces Problem Solutions Sample.

Table 1
AI-SOCO Dataset Statistics

Users Count	1,000	Unique Problems	6,553
Solutions Count	100,000	Maximum Solutions/Problem	61
Tokens Count	22,795,141	Minimum Solutions/Problem	1
Whitespaces Count	46,944,461	Avg. Solutions/Problem	15.26
Unique Tokens	1,171,991	Avg. Solutions/Codeforces Index	2,439.02
Avg. Solutions/User	100	Median Solutions/Problem	12
Avg. Tokens/Solution	227.95	Unique User Countries	78
Avg. Whitespaces/Solution	469.45	Avg. Solutions/User Country	1,282.05
Maximum Tokens in a Solution	10,189	Minimum Tokens in a Solution	3

smallest possible time period to avoid the possibility of users changing their coding styles over time. The dataset was split randomly with 50% used for training, 25% used for development and the remaining 25% used for testing. The training, development and unlabeled testing datasets are available through our GitHub repository.⁴ The full version of the dataset (with testing dataset labels) is available through Zenodo.⁵ Table 1 shows some statistics about the dataset, while Figure 3 shows the solutions distribution over CodeForces indexes and Figure 4 shows the solutions distribution over C++ compiler versions.

For a multi-class classification problem like ours, many evaluation metrics are suggested in the literature. However, since our dataset is completely balanced, we are evaluating the performance of the systems with accuracy since it is a simple and easily interpreted metric.

4. Baseline Approaches

In this section, we present the baseline approaches we developed for this task along with their results. All implementations mentioned in this section are available on our GitHub repository (See Footnote 4). Table 2 shows all baselines results. The pretrained and fine-tuned models are

⁴<https://github.com/AliOsm/AI-SOCO>

⁵<https://zenodo.org/record/4059840>

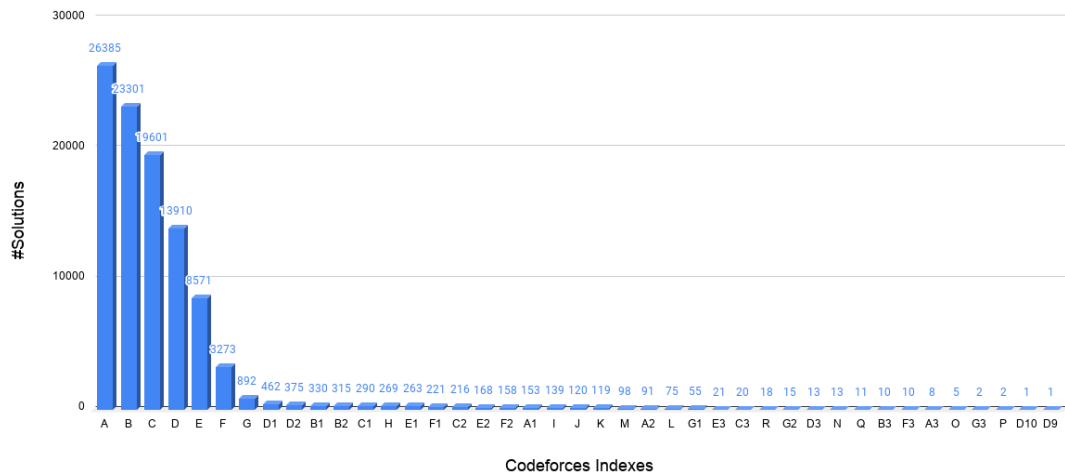


Figure 3: Number of Solutions per CodeForces Index.

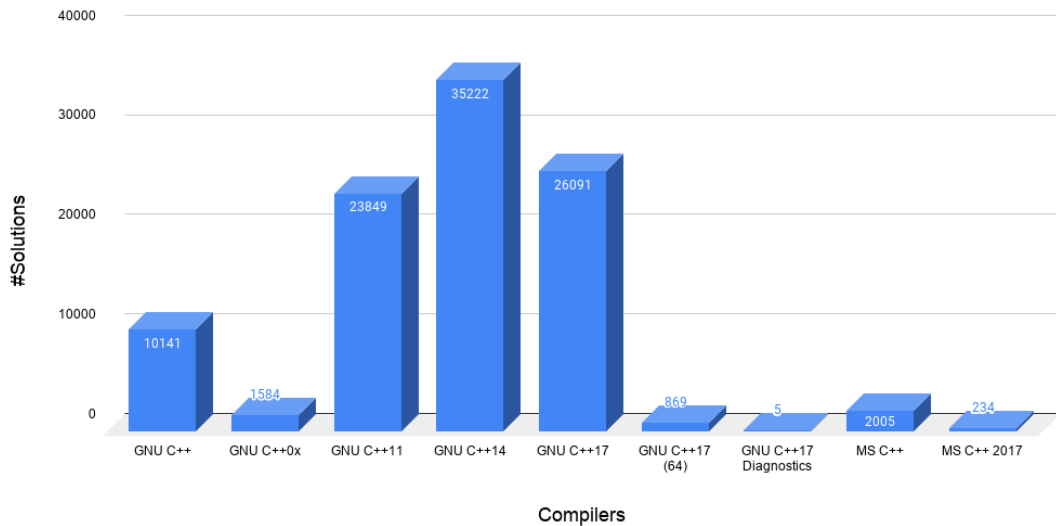


Figure 4: Number of Solutions per C++ Compiler.

available through HuggingFace models hub.⁶

4.1. Random Baseline

This baseline predicts a random author for each piece of code from the list of 1,000 authors. Its expected accuracy is 0.1%. We used this baseline to put a threshold for systems, whether they

⁶<https://huggingface.co/models?search=ai-soco>

Table 2

AI-SOCO Baselines Results on Development and Testing Datasets

Baseline Name	Development	Testing
Random	~0.1%	~0.1%
Characters Logistic	29.25%	29.92%
TF-IDF KNN	62.13%	62.78%
C++ RoBERTa Tiny	87.66%	87.46%
C++ RoBERTa Tiny-96	91.12%	91.02%
C++ RoBERTa Small	93.19%	92.88%

are better or worse than the randomness.

4.2. Characters Logistic Baseline

To build a code representation, this baseline converts each source code into a vector that represents the count of the 100 printable characters. Then, it builds a logistic regression [4] model on the vectorized representations. It achieved an accuracy of 29.25% and 29.92% on the development and testing datasets, respectively.

4.3. TF-IDF KNN Baseline

In this baseline, we tried to improve the code representation by vectorizing the source codes using term frequency-inverse document frequency (TF-IDF) [5] with 10K features. These features are fed into a K-Nearest Neighbors (KNN) [6] classifier with $K = 25$. Its accuracy is 62.13% and 62.78% on the development and testing datasets, respectively. This is a significant improvement over the previous baseline. One downside of this approach is efficiency. This lazy learner baseline is very slow in producing its predictions. It took about four hours to predict all examples in the development or testing datasets.

4.4. C++ RoBERTa Tiny

We pretrained a RoBERTa [7] model with a single Transformer [8] layer and 12 attention heads using the concatenation of the training dataset source codes. The model trained on Google Colab [9] platform⁷ with 8 TPU cores for 200 epochs, 32×8 batch size, 512 max sequence length and Masked Language Model (MLM) objective. Other parameters were set to their default values as mentioned in the `run_mlm.py` script provided by the HuggingFace Transformers [10] package.⁸ To tokenize the source codes, we trained a byte level BPE tokenizer [11] from the HuggingFace Tokenizers package⁹ with vocabulary size equals to 30K. Before tokenization, each four continuous spaces were converted into a single tab character (`\t`). After that, the model was fine-tuned on a classification task to predict the source code author out of the predefined list of the 1,000 authors. The fine-tuning was done on Google Colab platform using V100 GPU for 10

⁷<https://colab.research.google.com>

⁸<https://github.com/huggingface/transformers>

⁹<https://github.com/huggingface/tokenizers>

epochs, 32 batch size and 512 max sequence length. Source codes longer than 512 were truncated while fine-tuning. Finally, the model achieved 87.66% and 87.46% accuracy on development and testing datasets, respectively.

4.5. C++ RoBERTa Tiny-96

To do more experimentation, we tried to increase the number of attention heads in the model discussed in Section 4.4 from 12 to 96 attention heads. The model was trained and fine-tuned using the same pretraining and fine-tuning procedures as the previously mentioned one. The only change was in the batch size, where we set it to 16×8 instead of 32×8 for training and 16 instead of 32 for fine-tuning. The model achieved 91.12% and 91.02% accuracy on development and testing datasets, respectively.

4.6. C++ RoBERTa Small

To compare the effect of increasing the number of attention heads vs. increasing the number of Transformer layers, we trained another model with 6 Transformer layers and 12 attention heads. This model follows C++ RoBERTa Tiny pretraining and fine-tuning procedures as well. It achieved 93.19% and 92.88% accuracy on development and testing datasets, respectively.

5. Competition Results

The competition was hosted on the CodaLab¹⁰ platform. Surprisingly, 60 teams registered for the competition, 14 of which made 94 runs on the Evaluation and Post-Evaluation phases. The participating teams' accuracy results were higher than expected, ranging from 74.52% up to 95.11% on the testing dataset.

5.1. Overview of Participating Teams

In this section, we discuss the details of the teams' submissions that were disclosed by the time this paper is written. Table 3 shows the results of all submissions.

We start by discussing the submissions made by the deadline.

1. **Team UoB [12]:** The best approach in the competition was proposed by the Team UoB. The authors proposed several approaches and the best one was built using byte-level n-grams. Each source code is encoded into a vector representing the 20,000 most commonly occurring byte-level n-grams in the training dataset, where $n = 6$. Each n-gram in the vector was represented by a binary count, with 1 indicating the n-gram was present and 0 indicating the n-gram was absent in the given source code. The training dataset n-gram vector representations were then used to train a densely connected neural network classifier, which could then be used to predict authorship of the development and testing datasets. Their approach achieved 95.41% and 95.11% accuracy on the development and testing datasets, respectively.

¹⁰<https://competitions.codalab.org/competitions/25148>

2. **Team Yang1094 [13]:** According to the official results, the second best approach is Team Yang1094. Similar to some of the models of Team UoB [12], Team Yang1094 used an approach that relies on word-level and character-level n-gram features that are used to train logistic regression classifier. The character-level n-gram features (where $n = 2, \dots, 7$) filtered using term frequency-inverse document frequency (TF-IDF) gave the best results for this team with an accuracy of 94.28%.
3. **Team Alexa [14]:** Team Alexa proposed an ensemble of Naïve Bayes (NB) models and pretrained models. Specifically, their best model is a weighted combination of Multinomial NB (MultinomialNB), Bernoulli NB (BernoulliNB), and two versions of CodeBERTa. For the NB models, their features are the top 30,000 character-level n-grams (where $n = 1, \dots, 5$) with TF-IDF, whereas the CodeBERTa versions differ only in the learning rate. The resulting accuracy is 93.36%.
4. **Team LAST [15]:** This team's submissions are all based on the KNN procedure with $K = 1$. As for the preprocessing and feature extraction steps, they first used an indentation-aware tokenization followed by extracting token-level n-grams (where $n = 1, \dots, 4$) and skip-grams (1 for each 3-gram and 3 for each 4-gram). Then, they used a binary coding of the presence of the features and relevance frequency as a supervised term weighting scheme along with L2 normalization. The best accuracy reported for this team is 92.19%.
5. **Team FSU_HLJIT [16]:** Inspired by earlier work on plagiarism detection and microblog filtering, this team approached the task at hand as a ranking problem. The best performing system from this team ranked the source codes according to the number of occurrences of the character level n-grams (where $n = 15$). Despite its simplicity, this system got an accuracy of 91.57%.
6. **Team UMUTeam [17]:** This team's top system consisted of a combination of char-level n-grams and other author's traits features. The former were combinations between 1 and 8 length applying TF-IDF with Sublinear TF scaling. On the other hand, the author's traits included (1) the average length of code blocks and comments, (2) regular expressions to detect languages other than English, such as Arabic, Indian, or Russian; as well as (3) characters employed in the creation of ASCII-art. In addition, since some source codes included contact data such as nicknames, URLs, or Twitter accounts, the UMUTeam also compiled a Bag of Words (BOW) composed of the tokens that followed the words *name* and *author* in the block comments. Once all the features were compiled, they were filtered by applying a feature selection consisting of discarding low-variance features and the models were trained using a Random Forest (RF) classifier. The result is an accuracy of 91.16%.
7. **Team Abdalrhman:** This team used word-level and character-level n-grams with some preprocessing and cleaning steps for the dataset. For the preprocessing step; they normalized numbers into a special token then cleaned all lines with length greater than 50 or less than 4 characters. As for extracting the features, they used word-level and character-level (n-grams from text inside word boundaries) using TF-IDF features, (where $n = 2, \dots, 6$) for both word-level and character-level. Then the features are stacked from both vectorizers and fed into Linear Support Vector Classification (LinearSVC) which achieved an accuracy of 90.88%.

8. **Team bits_nlp_2020 [18]:** This team’s top system used Word2Vec embeddings along with a Convolutional Neural Network (CNN). The word2vec embeddings were trained from scratch. After concatenating vectors of both Continuous BOW (CBOW) and skip-gram for each word, the authors used the resulting embeddings for their CNN. The submitted system consisted of an ensemble of five such models, with different hyperparameters. They assigned a result for each test case based on whichever model had the highest confidence. The resulting accuracy was 90.64%.
9. **Team gaojiaming:** This team did not disclose the details of its submitted systems, which achieved an accuracy of 86.16%.
10. **Team Chanchal [19]:** This team used an abstract syntax tree (AST) for the tokenization of the code. Then, they applied TF-IDF to word-level and character-level n-grams to generate the code representation. Finally, different classifiers were used. Their best model used word bigrams on a stacked model. The stacked model is an ensemble of extra tree classifier, RF classifier, and XG-Boost classifier. The achieved accuracy is 82.95%.
11. **Team Kode_Stylers [20]:** This team used what they call “naturalness” of code as the key to their solution. They used different methods to obtain features such as tokenization, n-gram TF-IDF, warning messages, and coding styles. They used these features to train different classifiers, such as RF and Transformer. Then, they applied an ensemble approach to get better results. The best result obtained by this team is 82.08%.
12. **Team meghana:** This team did not disclose the details of its submitted systems, which achieved an accuracy of 81.2%.

After the competition’s deadline, two teams made interesting submissions. Their details are as follows.

- **Team Twist Bytes:** This team did not disclose the details of its submitted systems, which achieved an accuracy of 94.4%. This result is the second best result reported for this task.
- **Team SSNCSE_NLP [21]:** This team used character count vectorization and TF-IDF to extract features that are then fed into an RF classifier. The resulting accuracy is 85.73%.

5.2. Discussion and Error Analysis

There are some common characteristics among all participating teams. For example, they have all explored the use of different versions of n-grams. In most cases, n-gram features proved to be very powerful achieving the best (or near-best) results. Moreover, byte-level and character-level n-grams seem more promising than word-level n-grams. As for classification, different models showed strong performance individually and within ensembles. Actually, many teams have reported that ensembling only improves the results by a small margin. Finally, deep learning models and pretrained models showed good performance, but they do not emerge as clear winners like the case for many authorship analysis tasks [22, 23]. This applies to general models, such as RoBERTa, as well as models customized for handling source code, such as CodeBerta.

Looking at the predictions made by the participating systems, one can see certain patterns and trends. Appendix A shows five “easy” samples that each participating system managed to

Table 3

Summary of the results of the participating teams. The asterisk sign (*) is used to indicate baseline approaches while the pound sign (#) is used to indicate teams which made their submissions after the deadline.

Rank	Team	Accuracy
1	UoB	95.11%
2	Yang1094	94.28%
3	Alexa	93.36%
*	C++ RoBERTa Small Baseline	92.88%
4	LAST	92.19%
5	FSU_HLJIT	91.57%
6	UMUTeam	91.16%
*	C++ RoBERTa Tiny-96 Baseline	91.02%
7	Abdalrhman	90.88%
8	bits_nlp_2020	90.64%
*	C++ RoBERTa Tiny Baseline	87.46%
9	gaojiaming	86.16%
10	Chanchal	82.95%
11	Kode_Stylers	82.08%
12	meghana	81.20%
*	TF-IDF KNN Baseline	62.78%
*	Characters Logistic Baseline	29.92%
*	Random Baseline	0.08%
#	Twist Bytes	94.40%
#	SSNCSE_NLP	85.73%

predict correctly, whereas Appendix B shows five “hard” samples that none of the participating system managed to predict correctly. Easy samples tend to exhibit easily detectable unique patterns, such as the use of typedef and #define, whereas hard examples tend to be short or use patterns that are common among many authors. We performed a simple intersection between the participating systems’ predictions and discovered that 17,501 test samples out the 25,000 test samples (i.e., 70% of the test set) were correctly predicted by all systems. On the other hand, only 589 test samples (i.e., 2.36% of the test set) were never predicted correctly by any system. This indicates that the problem might not be as difficult as we initially expected.

To support this conjecture, we studied the participating systems’ predictions at the user level. We found out that our test set, has 133 “very easy” users out of 1,000 users. These users have such an easily identifiable coding styles that all participating systems managed to predict all of their testing cases correctly. One might think that having 13.3% very easy users in the test set is not that bad. However, if we consider slightly less easier users, i.e., ones with prediction rates of 90% or more across all participating systems, we end up with 603 users. Going to lower prediction rates (e.g., 80% rate of 70% rate) covers a vast majority of the users (83.1% and 93.2%, respectively).

On the other hand, very few users are difficult to identify. For example, the most difficult user, User 579, had a prediction rate of only 35.14% across all participating systems. Other difficult users with prediction rates lower than 50% include User 198 (prediction rate 36%), User 998

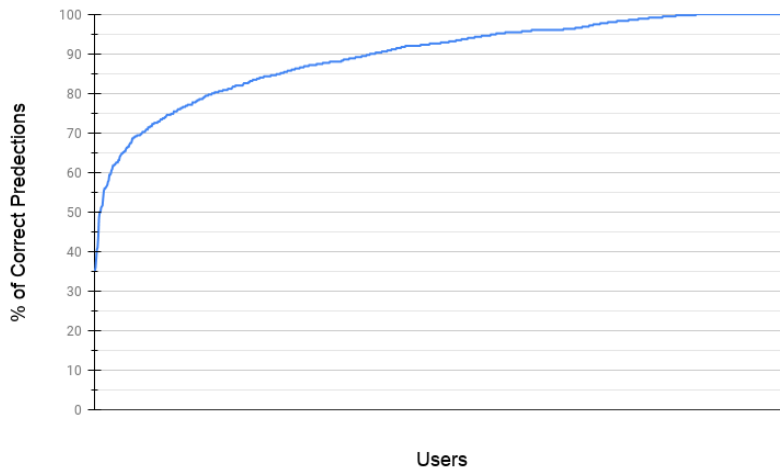


Figure 5: Percentage of Correctly Predicted Problems per User.

(prediction rate 38.29%), User 945 (prediction rate 40.57%), User 563 (prediction rate 41.14%), User 780 (prediction rate 44%), and User 205 (prediction rate 48.57%). Nonetheless, the number of hard users is rather low, which justifies the relatively high accuracy levels achieved by all participating systems. Figure 5 shows the percentage of correctly predicted samples per user.

Finally, instead of looking at this issue from the sample-level or the user-level, we look at it from the problem-level. It turns out that, for most problems, all participating systems managed to correctly identify the code authors. For very few problems, none of the participating systems managed to correctly identify any of the code authors for these problems. Figure 5 shows the percentage of systems that predicted samples of problems correctly.

Acknowledgments

The work of the last author was partially funded by the Spanish MICINN under the research project MISINFALKEHATE on MISinformation and MIScommunication in social media: FAKE news and HATE speech (PGC2018-096212-B-C31); as well as by the European Cooperation in Science and Technology under the COST Action 17124 DigForAsp.

References

- [1] E. Flores, P. Rosso, L. Moreno, E. Villatoro-Tello, PAN@FIRE: Overview of SOCO track on the detection of SOurce COde re-use, in: Notebook Papers of FIRE 2014, 2014.
- [2] E. Flores, P. Rosso, E. Villatoro-Tello, L. Moreno, R. Alcover, V. Chirivella, PAN@FIRE: Overview of CL-SOCO track on the detection of Cross-Language SOurce COde re-use., in: Post Proceedings of the Workshops at the 7th Forum for Information Retrieval Evaluation,

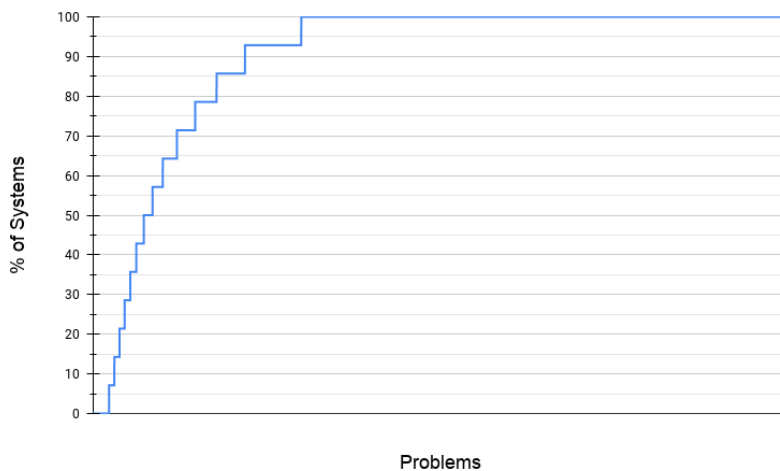


Figure 6: Percentage of Systems that Predicted Problems Correctly.

- Gandhinagar, India, December 4-6, 2015, volume 1587 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2015, pp. 1–5.
- [3] F. Rangel, F. González, F. Restrepo, M. Montes, P. Rosso, PAN@FIRE: Overview of the PR-SOCO track on Personality Recognition in SOURCE CODE, in: Working notes of FIRE 2016 - Forum for Information Retrieval Evaluation, Kolkata, India, December 7-10, 2016, volume 1737 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2016, pp. 1–15.
 - [4] C.-Y. J. Peng, K. L. Lee, G. M. Ingersoll, An introduction to logistic regression analysis and reporting, *The journal of educational research* 96 (2002) 3–14.
 - [5] G. Salton, C. Buckley, Term-weighting approaches in automatic text retrieval, *Information processing & management* 24 (1988) 513–523.
 - [6] G. Guo, H. Wang, D. Bell, Y. Bi, K. Greer, Knn model-based approach in classification, in: *OTM Confederated International Conferences "On the Move to Meaningful Internet Systems"*, Springer, 2003, pp. 986–996.
 - [7] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, V. Stoyanov, Roberta: A robustly optimized bert pretraining approach, *arXiv preprint arXiv:1907.11692* (2019).
 - [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, in: *Advances in neural information processing systems*, 2017, pp. 5998–6008.
 - [9] T. Carneiro, R. V. M. Da Nóbrega, T. Nepomuceno, G.-B. Bian, V. H. C. De Albuquerque, P. P. Reboucas Filho, Performance analysis of google colab as a tool for accelerating deep learning applications, *IEEE Access* 6 (2018) 61677–61685.
 - [10] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, et al., Huggingface’s transformers: State-of-the-art natural language processing, *ArXiv* (2019) arXiv–1910.
 - [11] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, Language models are

unsupervised multitask learners, OpenAI blog 1 (2019) 9.

- [12] A. Crosby, H. T. Madabushi, UoB at AI-SOCO 2020: Approaches to Source Code Classification and the Surprising Power of n-grams, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [13] Y. Yang, L. Kong, Z. Han, Y. Han, H. Qi, N-gram-based Authorship Identification of Source Code, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [14] M. Bni Younes, N. Al-Khdour, Team Alexa at Authorship Identification of SOURCE CODE (AI-SOCO), in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [15] Y. Bestgen, Boosting a KNN Classifier by Improving Feature Extraction for Authorship Identification of Source Code, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [16] Z. Han, T. Li, X. Wang, Y. Xu, M. Wu, Z. Li, Z. Wu, Y. Han, Ranking-based and Classification-based Approaches for Code Author Identification, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [17] J. A. García-Díaz, R. Valencia-García, UMUTeam at AI-SOCO'2020: Source Code Authorship Identification based on Character N-Grams and Author's Traits, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [18] A. V. Mandalam, Abhishek, Embedding-based Authorship Identification of Source Code, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [19] C. Suman, A. Raj, S. Saha, P. Bhattacharyya, Source Code Authorship Attribution Using Stacked Classifier, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [20] P. Sriiesaranusorn, S. Wattanakriengkrai, T. Son, T. Tanaka, C. Wiraatmaja, T. Ishio, R. G. Kula, Kode_Stylers: Author Identification through Naturalness of Code: An Ensemble Approach, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [21] N. N. A. Balaji, B. Bharathi, SSNCSE_NLP@Authorship Identification of SOURCE CODE (AI-SOCO) 2020, in: Proceedings of The 12th meeting of the Forum for Information Retrieval Evaluation (FIRE 2020), CEUR Workshop Proceedings, CEUR-WS.org, 2020.
- [22] W. Daelemans, M. Kestemont, E. Manjavacas, M. Potthast, F. Rangel, P. Rosso, G. Specht, E. Stamatatos, B. Stein, M. Tschuggnall, et al., Overview of pan 2019: Bots and gender profiling, celebrity profiling, cross-domain authorship attribution and style change detection, in: International Conference of the Cross-Language Evaluation Forum for European Languages, Springer, 2019, pp. 402–416.
- [23] M. N. AlRashdan, M. Abdullah, M. Al-Ayyoub, Y. Jararweh, Authorship analysis of english and spanish tweets, Proceedings of the Association for Information Science and Technology 57 (2020) e261.

A. Easy Samples

Listing 1: Easy Sample 1 | Problem ID: 99992 | User ID: 814

```
1 #include <bits/stdc++.h>
2 typedef long long ll;
3 typedef long double ld;
4 #define pb push_back
5 using namespace std;
6
7
8 int n, s;
9 int last[1010];
10 int main(){
11     ios::sync_with_stdio(0);
12     cin.tie(0);
13
14     cin >> n >> s;
15     for(int i = 0; i < n; i++){
16         int f, t;
17         cin >> f >> t;
18         last[f] = max(last[f], t);
19     }
20     int t = 0;
21     for(int i = s; i >= 1; i--){
22         if(t < last[i]) t = last[i];
23         t++;
24     }
25
26     cout << t << endl;
27
28     return 0;
29 }
```

Listing 2: Easy Sample 2 | Problem ID: 99452 | User ID: 742

```
1 #include <bits/stdc++.h>
2
3 using namespace std;
4
5 //using ll = int64_t;
6 using ll = long long;
7 using ull = uint64_t;
8 using i32 = int32_t;
9 using u32 = uint32_t;
10 using i64 = int64_t;
11 using u64 = uint64_t;
12 using pii = pair<int, int>;
13 using pll = pair<ll, ll>;
14 using ld = double;
15
16 #define X first
17 #define Y second
```

```

18
19 #ifndef ONLINE_JUDGE
20 #define FWRITE
21 #endif
22
23 namespace io
24 {
25
26 #ifndef FWRITE
27 #include <unistd.h>
28 #endif
29     const int BUFSIZE = 1 << 20;
30
31     int isize, osize;
32     char ibuf[BUFSIZE + 10], obuf[BUFSIZE + 10];
33     char* is, * it, * os = obuf, * ot = obuf + BUFSIZE;
34
35     char getchar()
36     {
37         if (is == it)
38         {
39             is = ibuf;
40 #ifdef FWRITE
41             it = ibuf + fread(ibuf, 1, BUFSIZE, stdin);
42 #else
43             it = ibuf + read(STDIN_FILENO, ibuf, BUFSIZE);
44 #endif
45             if (is == it) return EOF;
46         }
47         return *is++;
48     }
49
50     char getalpha()
51     {
52         char c = getchar();
53         while (!isalpha(c)) c = getchar();
54         return c;
55     }
56
57     void putchar(char c)
58     {
59         *os++ = c;
60         if (os == ot)
61         {
62 #ifdef FWRITE
63             fwrite(obuf, 1, BUFSIZE, stdout);
64 #else
65             write(STDOUT_FILENO, obuf, BUFSIZE);
66 #endif
67             os = obuf;
68         }
69     }
70
71     int inp() {

```

```

72     int x = 0, f = 0; char ch;
73     for (ch = getchar(); !isdigit(ch); ch = getchar())
74     {
75         if (ch == EOF) return -1;
76         if (ch == '-') f = 1;
77     }
78     for (; isdigit(ch); x = x * 10 + ch - '0', ch = getchar());
79     return f ? -x : x;
80 }
81
82 ll inp_ll() {
83     ll x = 0; int f = 0; char ch;
84     for (ch = getchar(); !isdigit(ch); ch = getchar())
85         if (ch == '-') f = 1;
86     for (; isdigit(ch); x = x * 10 + ch - '0', ch = getchar());
87     return f ? -x : x;
88 }
89
90 template<class T>
91 bool read(T& x)
92 {
93     x = 0;
94     char ch = getchar();
95     if (ch == EOF) return 0;
96     for (; !isdigit(ch); )
97     {
98         ch = getchar();
99         if (ch == EOF) return 0;
100    }
101    for (; isdigit(ch); x = x * 10 + ch - '0', ch = getchar());
102    return 1;
103 }
104
105 template<class T>
106 void write(T x)
107 {
108     static char s[22];
109     static char* it = s + 20;
110     static char* end = s + 20;
111
112     if (x < 0)
113     {
114         putchar('-');
115         x = -x;
116     }
117
118     do
119     {
120         *--it = x % 10 + '0';
121         x /= 10;
122     } while (x);
123     /*
124     if (!x)
125         *--it = '0';

```

```

126     while (x)
127     {
128         *-- it = x%10+'0';
129         x /= 10;
130     }
131     /*
132     for (; it < end; ++it)
133         putchar(*it);
134     }
135
136     template<>
137     void write(const char* s)
138     {
139         for (; *s; ++s) putchar(*s);
140     }
141
142     template<>
143     void write(char c)
144     {
145         putchar(c);
146     }
147
148     template<class T, class V>
149     void write(T x, V y)
150     {
151         write(x);
152         write(y);
153     }
154
155     template<class T>
156     void writeln(T x)
157     {
158         write(x);
159         putchar('\n');
160     }
161
162     struct ender
163     {
164         ~ender()
165         {
166             if (os != obuf)
167 #ifdef FWRITE
168                 fwrite(obuf, 1, os - obuf, stdout);
169 #else
170                 write(STDOUT_FILENO, obuf, os - obuf);
171 #endif
172         }
173     } __ender;
174
175 }
176
177 template<class T>
178 void print(const T& a)
179 {

```



```

180     for (auto x : a) printf("%d ", x); puts("");
181 }
182
183 int64_t power(int64_t a, int64_t b, int64_t p)
184 {
185     if (!b) return 1;
186     int64_t t = power(a, b >> 1, p);
187     t = t * t % p;
188     if (b & 1) t = t * a % p;
189     return t;
190 }
191
192 //mt19937_64 rd(chrono::steady_clock::now().time_since_epoch().count());
193 mt19937 rd(chrono::steady_clock::now().time_since_epoch().count());
194
195 using namespace io;
196
197 template<class T>
198 inline void freshmin(T& a, const T& b)
199 {
200     if (a > b) a = b;
201 }
202
203 template<class T>
204 inline void freshmax(T& a, const T& b)
205 {
206     if (a < b) a = b;
207 }
208
209 const ll B = 31;
210 //const ll MOD = 998244353;
211 const int INF = 1000000010;
212 //const ll INFll = 1000000000000000000LL;
213 const int MAXN = 400010;
214
215 int dx[] = { -1, 1, 0, 0, -1, -1, 1, 1 };
216 int dy[] = { 0, 0, -1, 1, -1, 1, -1, 1 };
217
218 ld det(ld x1, ld y1, ld x2, ld y2, ld x3, ld y3)
219 {
220     return x1 * y2 - x2 * y1 + x2 * y3 - x3 * y2 + x3 * y1 - x1 * y3;
221 }
222
223 int n;
224 vector<int> v[MAXN];
225 ll a[MAXN];
226
227 void solve()
228 {
229     n = inp();
230     for (int i = 1; i <= n; ++i) a[i] = 1LL << i;
231     ll sum = a[n];
232     for (int i = 1; i < n / 2; ++i) sum += a[i];
233     for (int i = n / 2; i < n; ++i) sum -= a[i];

```

```

234     sum = abs(sum);
235     writeln(sum);
236 }
237
238 int main()
239 {
240
241     for (int T = inp(); T --; )
242         solve();
243
244     return 0;
245 }

```

Listing 3: Easy Sample 3 | Problem ID: 98902 | User ID: 721

```

1  #include <bits/stdc++.h>
2  #include <ext/pb_ds/assoc_container.hpp>
3  #include <ext/pb_ds/tree_policy.hpp>
4
5  #pragma GCC optimize(-O3)
6  #pragma GCC optimize(0 fast)
7  #pragma GCC optimize("unroll-loops")
8
9  #define fi first
10 #define se second
11 #define sqr(x) (x) * (x)
12 #define p_b push_back
13 #define m_p make_pair
14 #define pll pair<ll, ll>
15 #define all(v) v.begin(), v.end()
16 #define pw(x) (1ll << x)
17
18 using namespace std;
19 using namespace __gnu_pbds;
20 typedef long long ll;
21 typedef long double ld;
22 const ll MAXN = 1123456;
23 const ll N = 1e6;
24 const ll MOD = 1e9 + 7;
25
26 template <typename T> using ordered_set = tree<T, null_type, less<T>,
    rb_tree_tag, tree_order_statistics_node_update>;
27
28 template <typename T> void vout(T s){cout << s << endl; exit(0);}
29
30
31
32 int main(){
33     ios_base :: sync_with_stdio(0);
34     cin.tie(0);
35
36     vector <pll> v;
37
38     ll n, m;

```

```

39     cin >> n >> m;
40
41     ll l = 1, r = n;
42
43     ll le = 1, ri = m;
44
45     ll st = 0;
46
47     while(l <= r){
48         if(!st){
49             v.p_b({l, le});
50             le++;
51         }else{
52             v.p_b({r, ri});
53             ri--;
54             if(ri == 0){
55                 l++;
56                 le = 1;
57                 r--;
58                 ri = m;
59             }
60         }
61         st = 1 - st;
62         if(l == r && le > ri)break;
63     }
64
65     for(auto i : v)cout << i.fi << " " << i.se << "\n";
66
67     return 0;
68 }

```

Listing 4: Easy Sample 4 | Problem ID: 98374 | User ID: 930

```

1 #include <bits/stdc++.h>
2 #define pb push_back
3 #define fast ios_base::sync_with_stdio(0);cin.tie(0);cout.tie(0);
4 #define int long long
5 #define pii pair<int,int>
6 #define all(x) x.begin(),x.end()
7 #define lb lower_bound
8 #define ld long double
9 using namespace std;
10 const int N = 3e5+1;
11 vector<pii> adj[N];
12 int dp[N];
13 int a[N];
14 int ans = 0;
15 void dfs (int src,int par){
16     dp[src] = a[src];
17     multiset<int,greater<int>> st;
18     for(auto it:adj[src]){
19         if(it.first!=par){
20             dfs(it.first,src);
21             dp[src] = max(dp[src],a[src] - it.second + dp[it.first]);

```

```

22     st.insert(dp[it.first] - it.second);
23 }
24 }
25 int nn = st.size();
26 if(nn>=2){
27     int fmx = *st.begin();
28     st.erase(st.begin());
29     int smx = *st.begin();
30     ans = max(ans, fmx + smx + a[src]);
31 }
32 ans = max(ans, dp[src]);
33 }
34 signed main() {
35     ios_base::sync_with_stdio(0);
36     cin.tie(0);
37     cout.tie(0);
38     int n;
39     cin>>n;
40     for(int i = 1; i<=n; i++)
41         cin>>a[i];
42     for(int i = 1; i<=n-1; i++){
43         int u, v, w;
44         cin>>u>>v>>w;
45         adj[u].pb({v, w});
46         adj[v].pb({u, w});
47     }
48     dfs(1, -1);
49     cout<<ans;
50 }

```

Listing 5: Easy Sample 5 | Problem ID: 97797 | User ID: 665

```

1 #include <bits/stdc++.h>
2 #include <ext/pb_ds/assoc_container.hpp>
3 #include <ext/pb_ds/tree_policy.hpp>
4 using namespace std;
5 using namespace __gnu_pbds;
6
7 template<typename T>
8 using ordered_set = tree<T, null_type, less<T>, rb_tree_tag,
9     tree_order_statistics_node_update>;
10 #define ll long long int
11 #define vi vector<int>
12 #define vl vector<ll>
13 #define pii pair<int, int>
14 #define pll pair<ll, ll>
15 #define pdd pair<double, double>
16 #define vii vector<pii>
17 #define vll vector<pll>
18 #define vd vector<double>
19 #define vb vector<bool>
20 #define el "\n"
21 #define ull unsigned long long int
22 #define ld long double

```

```

22 #define fi first
23 #define se second
24 #define pb push_back
25 #define eb emplace_back
26 #define mp make_pair
27 #define sqr(x) ((x) * (x))
28 #define FOR(i, a, b) for(ll i=a; i<=b; i++)
29 #define RFOR(i, b, a) for(ll i=b; i>=a; i--)
30 #define all(vec) (vec).begin(),(vec).end()
31 template <typename T> void maxi(T &x, T y) {if (y > x) x = y; }
32 template <typename T> void mini(T &x, T y) {if (y < x) x = y; }
33 #define error(args ...) { string _s = #args; replace(_s.begin(), _s.end(), ',',
    , ' '); stringstream _ss(_s); istream_iterator<string> _it(_ss); err(_it,
    args); }
34 void err(istream_iterator<string> it) {}
35 template<typename T, typename... Args>
36 void err(istream_iterator<string> it, T a, Args... args) {
37     cerr << *it << " =: " << a << "\n";
38     // cerr << "\n";
39     err(++it, args...);
40 }
41 template<class T1, class T2, class T3>
42 struct triple {
43     T1 a; T2 b; T3 c;
44     triple() : a(T1()), b(T2()), c(T3()) {};
45     triple(T1 _a, T2 _b, T3 _c) :a(_a), b(_b), c(_c) {}
46 };
47 template<class T1, class T2, class T3>
48 bool operator <(const triple <T1,T2,T3>&t1, const triple <T1,T2,T3>&t2) {
49     if(t1.a!=t2.a) return t1.a<t2.a;
50     else if(t1.b!=t2.b) return t1.b<t2.b;
51     else return t1.c<t2.c;
52 }
53 template <typename T1, typename T2>
54 inline std::ostream& operator << (std::ostream& os, const std::pair<T1, T2>&
    p) {
55     return os << "(" << p.first << ": " << p.second << ")";
56 }
57 template<typename T>
58 inline std::ostream &operator << (std::ostream & os, const std::vector<T>& v)
    {
59     bool first = true;
60     os << "[";
61     for(unsigned int i = 0; i < v.size(); i++) {
62         if(!first) os << ", ";
63         os << v[i];
64         first = false;
65     }
66     return os << "]";
67 }
68 template<typename T>
69 inline std::ostream &operator << (std::ostream & os, const std::set<T>& v) {
70     bool first = true;
71     os << "{";

```

```

72     for (typename std::set<T>::const_iterator ii = v.begin(); ii != v.end();
73         ++ii) {
74         if(!first) os << ", ";
75         os << *ii;
76         first = false;
77     }
78     return os << "}";
79 }
80 template<typename T1, typename T2>
81 inline std::ostream &operator << (std::ostream & os, const std::map<T1, T2>& v
82 ) {
83     bool first = true;
84     os << "{";
85     for (typename std::map<T1, T2>::const_iterator ii = v.begin(); ii != v.
86         end(); ++ii) {
87         if(!first) os << ", ";
88         os << *ii ;
89         first = false;
90     }
91     return os << "}";
92 }
93 template<typename T>
94 inline std::ostream &operator << (std::ostream & os, const std::unordered_set<
95     T>& v) {
96     return os << std::set<T>(v.begin(), v.end());
97 }
98 template<typename T1, typename T2>
99 inline std::ostream &operator << (std::ostream & os, const std::unordered_map<
100     T1, T2>& v) {
101     return os << std::map<T1, T2>(v.begin(), v.end());
102 }
103
104 const ll MOD = 1e9+7;
105 const ll INF = 1e18;
106 const double EPS = 1e-6;
107 /* ***** */
108
109 int main()
110 {
111     ios_base::sync_with_stdio(false);
112     cin.tie(NULL);
113
114     int n;
115     cin >> n;
116     vector<int> a(n);
117     for (int i = 0; i < n; ++i) {
118         ll x;
119         cin >> x;
120         a[i] = (int) __builtin_popcountll(x);
121     }
122     ll ans = 0, sum = 0;
123     vector<int> cnt(2, 0);
124     cnt[0]++;
125     for (int e : a) sum += e, cnt[sum & 1]++;

```

```

121     for (int e : cnt) ans += (1ll) e * (e - 1) / 2LL;
122     for (int i = 0; i < n; ++i) {
123         sum = 0;
124         int mx = -1;
125         for (int j = i; j < i + 65 && j < n; ++j) {
126             sum += a[j];
127             maxi(mx, a[j]);
128             if (sum % 2 == 0 && mx > sum - mx) ans--;
129         }
130     }
131     cout << ans << el;
132
133     return 0;
134 }

```

B. Hard Samples

Listing 6: Hard Sample 1 | Problem ID: 43 | User ID: 428

```

1 #include <bits/stdc++.h>
2 using namespace std;
3
4 int main(){
5     int a, b, c;
6     cin >> a >> b >> c;
7     cout << min(a + 2, min(b + 1, c)) * 3 - 3;
8     return 0;
9 }

```

Listing 7: Hard Sample 2 | Problem ID: 15564 | User ID: 289

```

1 #include <iostream>
2
3 #include <vector>
4 #define pb push_back
5
6
7 using namespace std;
8 int n, st = 0;
9
10 vector<vector<int>> v;
11
12 int main()
13 {
14     cin >> n;
15     v.resize(n+1);
16     v[st].pb(0);
17     st++;
18     for (int i = 0; i < n; i++)
19     {
20         int x;

```

```

21     cin >> x;
22     int l = 0;
23     int r = st;
24     while (l < r)
25     {
26         int c = (l+r) / 2;
27         if (v[c].back() >= x) l = c+1; else r = c;
28     }
29     if (l == st || v[l].back() > x)
30     {
31         v[st].pb(x);
32         st++;
33         continue;
34     }
35     v[l].pb(x);
36 }
37 for (int i = 0; i < st; i++)
38 {
39     for (int j = 0; j < v[i].size(); j++) if (v[i][j] != 0) cout << v[i][
40         j] << ' ';
41     cout << "\n";
42 }
43 return 0;
44 }

```

Listing 8: Hard Sample 3 | Problem ID: 31849 | User ID: 148

```

1 #include <iostream>
2 #include <vector>
3 #include <algorithm>
4 #include <cmath>
5 #include <map>
6 #include <set>
7 #include <sstream>
8
9 using namespace std;
10
11 map<pair<int ,int >,int> lol;
12
13 int gcd(int a, int b)
14 {
15     if (!b) return a;
16     return gcd(b, a%b);
17 }
18
19 pair<int ,int > kek[200500];
20
21 int main()
22 {
23     int n;
24     cin >> n;
25     for (int i=0; i<n; i++)
26     {
27         string s;

```



```

28     cin >> s;
29     for (auto &c:s)
30     {
31         if (c<'0' || c>'9')
32         {
33             c=' ';
34         }
35     }
36     stringstream z(s);
37     int a,b,c;
38     z >> a >> b >> c;
39     a+=b;
40     b=gcd(a,c);
41     a/=b;
42     c/=b;
43     ++lol[{a,c}];
44     kek[i]={a,c};
45 }
46 for (int i=0;i<n;i++)
47 {
48     cout << lol[kek[i]] << ' ';
49 }
50 return 0;
51 }

```

Listing 9: Hard Sample 4 | Problem ID: 49881 | User ID: 71

```

1 #include <bits/stdc++.h>
2 #define pb push_back
3 using namespace std;
4 const int N=1e5+5;
5
6
7 int main()
8 {
9     ios_base::sync_with_stdio(0);cin.tie(0);cout.tie(0);
10
11     cout<<2000<<endl;
12     for(int i=1000;i>=1;i--)
13         cout<<i<<" "<<1<<" "<<i<<" "<<2<<endl;
14
15     for(int i=1;i<=1000;i++)
16         cout<<i<<" "<<1<<" "<<i<<" "<<2<<endl;
17     return 0;
18 }

```

Listing 10: Hard Sample 5 | Problem ID: 65320 | User ID: 21

```

1 #include <bits/stdc++.h>
2 using namespace std;
3
4 int input() {
5     char c; cin >> c;
6     return c == '1';

```

```

7 }
8
9 voProblem id print_ans(vector<int> r, vector<int> c) { | User ID:
10     int ans = 0;
11     for (int i: r) ans += i;
12     for (int i: c) ans += i;
13
14     cout << ans << '\n';
15     for (int i = 0; i < r.size(); i++) {
16         if (r[i]) cout << "row " << i << '\n';
17         if (c[i]) cout << "col " << i << '\n';
18     }
19 }
20
21 int main() {
22     ios::sync_with_stdio(false); cin.tie(0);
23
24     int n; cin >> n;
25     int a[n][n], c[n];
26     for (int i = 0; i < n; i++) {
27         for (int j = 0; j < n; j++) {
28             a[i][j] = input();
29         }
30     }
31     for (int i = 0; i < n; i++) {
32         for (int j = 0; j < n; j++) {
33             a[i][j] ^= input();
34         }
35     }
36     for (int i = 0; i < n; i++) {
37         c[i] = input();
38     }
39     vector<int> tog_row(n), tog_col(n);
40     for (int i = 0; i < n; i++) {
41         for (int j = 0; j < n; j++) {
42             if (c[i] + c[j] == 0) {
43                 if (a[i][j]) {
44                     cout << -1; exit(0);
45                 }
46             } else if (c[i] + c[j] == 1) {
47                 if (a[i][j]) {
48                     if (c[j] == 1) {
49                         if (!tog_row[i]) {
50                             tog_row[i] = 1;
51                         }
52                     } else {
53                         if (!tog_col[j]) {
54                             tog_col[j] = 1;
55                         }
56                     }
57                 }
58             }
59         }
60     }

```

```

61
62     for (int i = 0; i < n; i++) {
63         for (int j = 0; j < n; j++) {
64             if (tog_row[i] || tog_col[j]) {
65                 if (tog_row[i] && c[j]) a[i][j] ^= 1;
66                 if (tog_col[j] && c[i]) a[i][j] ^= 1;
67
68                 if (a[i][j]) {
69                     cout << -1; exit(0);
70                 }
71             }
72         }
73     }
74
75     string sample_space;
76     for (int i = 0; i < n; i++) {
77         if (!tog_row[i] && c[i]) {
78             for (int j = 0; j < n; j++) {
79                 if (!tog_col[j] && c[j]) {
80                     sample_space.push_back(a[i][j] + '0');
81                 } else {
82                     sample_space.push_back('X');
83                 }
84             }
85             break;
86         }
87     }
88     if (sample_space.empty()) {
89         print_ans(tog_row, tog_col);
90     }
91
92     for (int i = 0; i < n; i++) {
93         if (!tog_row[i] && c[i]) {
94             string space1, space2;
95             for (int j = 0; j < n; j++) {
96                 if (tog_col[j] == 0 && c[j]) {
97                     space1.push_back(a[i][j] + '0');
98                     space2.push_back(1 - a[i][j] + '0');
99                 } else {
100                     space1.push_back('X');
101                     space2.push_back('X');
102                 }
103             }
104
105             if (space1 == sample_space) continue;
106             if (space2 == sample_space) {
107                 tog_row[i] = 1;
108                 for (int j = 0; j < n; j++) {
109                     if (c[j]) a[i][j] ^= 1;
110                 }
111             } else {
112                 cout << -1;
113                 exit(0);
114             }

```

```
115     }
116   }
117   for (int i = 0; i < n; i++) {
118     if (sample_space[i] == '1') {
119       tog_col[i] = 1;
120     }
121   }
122
123   print_ans(tog_row, tog_col);
124
125   return 0;
126 }
```
