

Московский авиационный институт (национальный исследовательский университет)  
Институт № 8 «Компьютерные науки и прикладная математика»  
Кафедра № 806 «Вычислительная математика и программирование»

## Интеллектуальная система обнаружения признаков компрометации корпоративной электронной почты

Выпускная квалификационная работа бакалавра

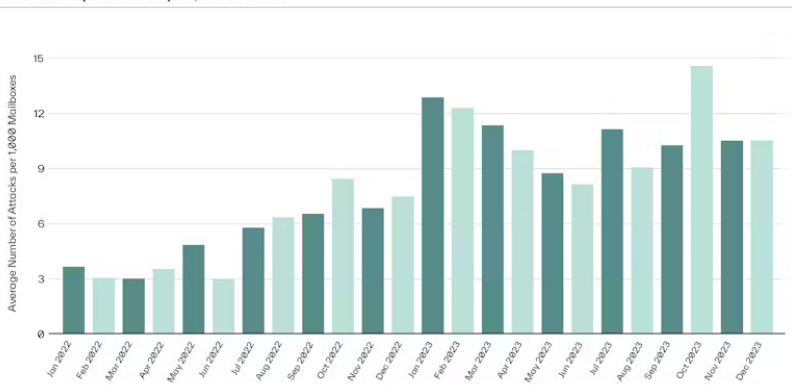
Студент группы М8О-406Б-20: Лисин Роман Сергеевич  
Научный руководитель: ст. преп. каф. 806 А. В. Борисов  
Консультант: канд. техн. наук, доц., доц. каф. 806 П. А. Ухов

Москва — 2024



В 2023 году количество BEC-атак выросло на 108% по сравнению с 2022 годом согласно источнику [1].

Median Monthly BEC Attacks per 1,000 Mailboxes



**Цель** — создание и обучение модели для автоматизации выявления компрометации корпоративной электронной почты.

**Задачи:**

- Подготовить тексты писем для обучения и тестирования моделей.
- Выполнить предобработку текстов для обучения и тестирования моделей.
- Реализовать различные модели для интеллектуального анализа деловой корреспонденции отправителя.
- Выявить лучшую модель для обнаружения компрометации корпоративной электронной почты.



Дано:

- Деловая корреспонденция компании Enron из открытого источника.

Требуется:

- Извлечь текст, написанный самим автором письма.
- Удалить подписи авторов, обращения.
- Обучить и протестировать модели на датасетах пяти выбранных сотрудников компании.



- Язык программирования Python
- Библиотеки для машинного обучения pandas, sklearn, numpy, nltk, tensorflow, pytorch, matplotlib, seaborn, tqdm, transformers
- Библиотека для работы с электронными письмами mail-parser



- Подготовка данных к обучению и тестированию
- Предобработка текстов писем
- Обучение и тестирование моделей
- Оценка результатов



## Пример электронного письма сотрудника компании Enron

```
Message-ID: <19790540.1075855679828.JavaMail.evans@thyme>  
Date: Tue, 12 Dec 2000 04:03:00 -0800 (PST)  
From: phillip.allen@enron.com  
To: christi.nicolay@enron.com  
Subject: Talking points about California Gas market  
Mime-Version: 1.0  
Content-Type: text/plain; charset=us-ascii  
Content-Transfer-Encoding: 7bit  
X-From: Phillip K Allen  
X-To: Christi L Nicolay  
X-cc:  
X-bcc:  
X-Folder: \Phillip_Allen_Dec2000\Notes Folders\Sent  
X-Origin: Allen-P  
X-FileName: pallen.nsf
```

Christy,

I read these points and they definitely need some touch up. I don't understand why we need to give our commentary on why prices are so high in California. This subject has already gotten so much press.

Phillip

----- Forwarded by Phillip K Allen/HOU/ECT on 12/12/2000  
12:01 PM -----



Пример датасета сотрудницы компании Enron Кейт Саймс с 2534 письмами

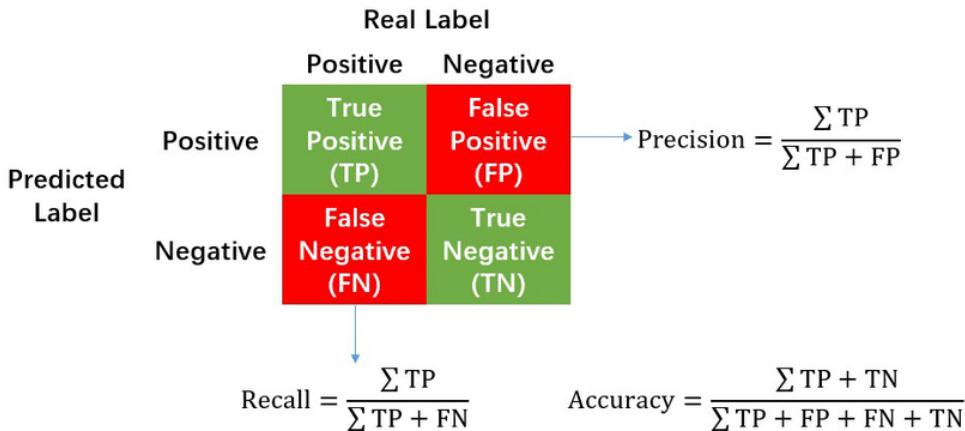
Unnamed: 0	text	label
0	I'm a little confused - 559066 is a Prebon dea...	1
1	Prebon is right on this. Both deals are 25 MW....	1
2	Hey there stranger!\n\nfor the pictures - they...	1
3	Two of these should have fees; two should not....	1
4	Mark's changing deal 581615 to APB - he had it...	1
...	...	...
2529	pending the sale of the Wilton Centre unit, I...	0
2530	Same to you! \nAnd I hope you and your family...	0
2531	\n I would appreciate your help in locating fi...	0
2532	\t2- SURVEY/INFORMATION EMAIL - 7/19/01\n\nCur...	0
2533	see if this works! If it does, see "13." w...	0





- Алгоритмы векторизации - Bag of Words, TF-IDF и токенизатор для нейронной сети BERT.
- Алгоритмы классификации - наивный классификатор Байеса, логистическая регрессия, метод опорных векторов, метод k-ближайших соседей, дерево решений, случайный лес, градиентный бустинг, перцептроны и рекуррентные, свёрточные, глубокие нейронные сети.



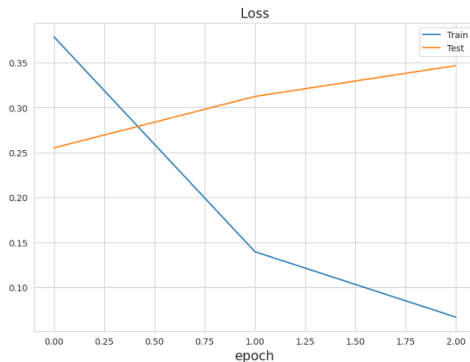
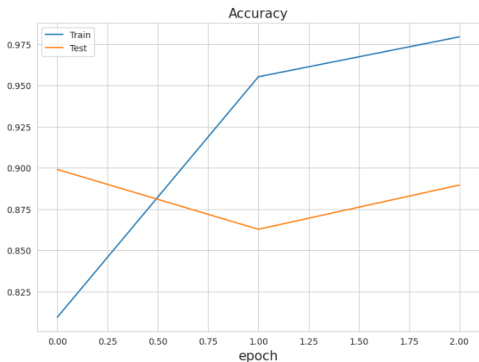


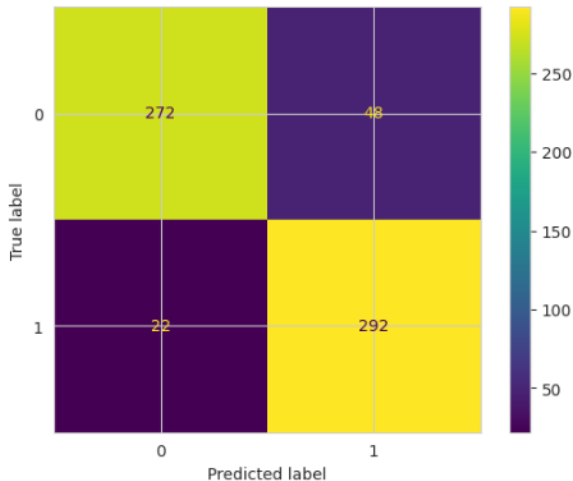
	Accuracy, %	Precision, %	Recall, %
BERT	88	86	90
SVM	87	92	81
Naive Bayes	86	82	92
Logistic Regression	86	90	82
CNN	85	85	86
LSTM	85	86	85
Random Forest	83	89	76
Gradient Boosting	81	86	73
Perceptron	80	90	69
KNN	78	75	86
RNN	77	75	83
Decision Tree	75	76	74



# Графики accuracy и loss по эпохам во время обучения и тестирования нейронной сети BERT

12





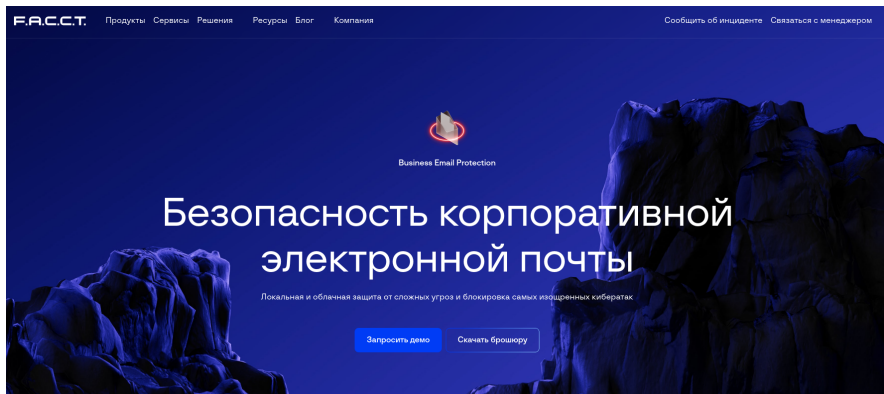
Репозиторий с исходным кодом расположен по ссылке  
[https://github.com/RomaMaster228/bec\\_detection](https://github.com/RomaMaster228/bec_detection)



1. Callie H. B. Report on BEC and VEC attacks. — URL: <https://abnormalsecurity.com/blog/bec-vec-attacks> (дата обращения 29.03.2024).



Данная работа будет использоваться в продукте по защите корпоративной почты Business Email Protection (BEP) российской компании F.A.C.C.T., занимающейся кибербезопасностью, для выявления компрометации корпоративной электронной почты.





- Классические этапы текстовой предобработки включают в себя перевод всех букв в тексте в нижний или верхний регистры, удаление цифр, чисел или замену на текстовый эквивалент, очистку от пунктуации, устранение стоп-слов, стемминг, лемматизацию.
- Также в предметной области корпоративной электронной почты из писем следует исключать HTML и CSS блоки, эмодзи, почтовые адреса, ссылки, пробельные символы с начала и конца текста, символы переноса строк.



В качестве набора данных для обучения будем использовать деловую корреспонденцию компании Enron, состоящую из 126841 электронных писем на английском языке.

```
enron_mail_20150507/  
└─ maildir  
    ├── allen-p  
    ├── arnold-j  
    ├── arora-h  
    ├── badeer-r  
    ├── bailey-s  
    ├── bass-e  
    ├── baughman-d  
    ├── beck-s  
    ├── benson-r  
    ├── blair-l  
    ├── brawner-s  
    ├── buy-r  
    ├── campbell-l  
    ├── carson-m  
    ├── cash-m  
    ├── causholli-m  
    ├── corman-s  
    ├── crandell-s  
    ├── cuilla-m  
    ├── dasovich-j  
    ├── davis-d  
    ├── dean-c  
    ├── delainey-d  
    └── derrick-j
```



- Компрометация корпоративной электронной почты является одним из самых разрушительных с финансовой точки зрения киберпреступлений, в результате которого только за предыдущий год убытки составили 2.7 миллиарда долларов.
- По данным ФБР средний ущерб успешной атаки на корпоративную электронную почту составляет более 125000 долларов.
- Одним из ярких примеров ВЕС является атака на американскую международную компанию по производству игрушек Mattel. В 2016 году хакеры взломали почтовый аккаунт недавно назначенного генерального директора компании и написали письмо сотруднику, в котором попросили перевести 3 миллиона долларов новому поставщику.

