



Suiside Prediction in Workplaces

Seng Fong Chang, Shuoxin Dai, Yiyang Duanmu, Derek He, Yunning Qiu

Table of Contents



Business Goals

01



Data Preprocessing &
Feature Engineering



Modeling &
Performance Evaluation



Business
Recommendations

04



Return on Investment



Conclusion



Future Improvement

05



Conclusion

06



Future Improvement

07

Business Goals

What business problems are we trying to solve, and what kind of solutions are we providing

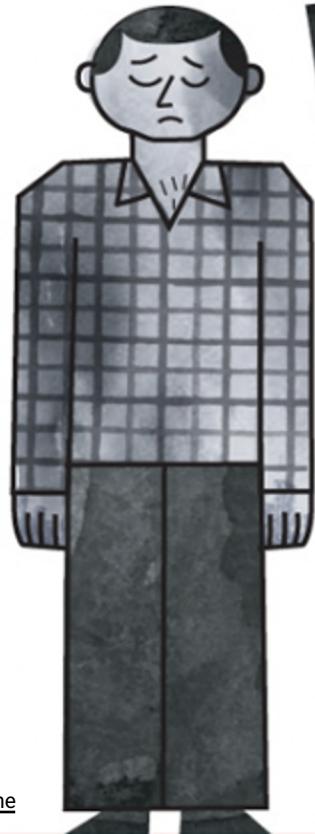


Image Source:
[hazards magazine](#)

What's the point?

Work-related suicides are not reportable. They are not counted. They are not a workplace prevention priority.

They should be. Suicides caused by bad jobs are a real and growing problem at work.

Problems & Solutions

Business Problem:

When employees attempt suicide, they cost companies money in the following ways:

1. Medical Costs (Hospitalization) and/or Compensation(Death)
2. Post-traumatic Stress Syndrome Training for other employees
3. Hurt Brand Images (Potential loss of existing employees, talents who could be onboard or customers)

Business Solution:

We want to save companies costs and improve employees' productivity, leading to increased profits.

We build a suicide prediction tool that predicts employees' intentions to suicide pre-emptively based on text data. Text classification identifies whether or not a person would commit suicide. Topic modeling identifies workers' flag behaviors or emotions related to mental health concerns, and companies can implement interventions to ensure corporate efficiency.

Suicides cases raise at work between 2018 and 2020, according to a report by the Bureau of Labor Statistics. Suicides occur most often among people in their working years of 24 to 64 years old.

Our Data

Source	The dataset is a collection of posts from "SuicideWatch" and "depression" subreddits of the Reddit platform. The posts are collected using Pushshift API. All posts that were made to "SuicideWatch" from Dec 16, 2008(creation) till Jan 2, 2021, were collected while "depression" posts were collected from Jan 1, 2009, to Jan 2, 2021.																		
Features	<p>The dataset displays a text column and a class column (suicide, non-suicide)</p> <table><thead><tr><th></th><th>text</th><th>class</th></tr></thead><tbody><tr><td>2</td><td>Ex Wife Threatening SuicideRecently I left my ...</td><td>suicide</td></tr><tr><td>3</td><td>Am I weird I don't get affected by compliments...</td><td>non-suicide</td></tr><tr><td>4</td><td>Finally 2020 is almost over... So I can never ...</td><td>non-suicide</td></tr><tr><td>8</td><td>i need helpjust help me im crying so hard</td><td>suicide</td></tr><tr><td>9</td><td>I'm so lostHello, my name is Adam (16) and I've...</td><td>suicide</td></tr></tbody></table>		text	class	2	Ex Wife Threatening SuicideRecently I left my ...	suicide	3	Am I weird I don't get affected by compliments...	non-suicide	4	Finally 2020 is almost over... So I can never ...	non-suicide	8	i need helpjust help me im crying so hard	suicide	9	I'm so lostHello, my name is Adam (16) and I've...	suicide
	text	class																	
2	Ex Wife Threatening SuicideRecently I left my ...	suicide																	
3	Am I weird I don't get affected by compliments...	non-suicide																	
4	Finally 2020 is almost over... So I can never ...	non-suicide																	
8	i need helpjust help me im crying so hard	suicide																	
9	I'm so lostHello, my name is Adam (16) and I've...	suicide																	
Number of Observations	232,074																		

Data Preprocessing & Feature Engineering

Image Source:
[SolveXia](#)



Data Preprocessing

Translate emoji to English description

For better interpretation, we replaced emojis with their English names (ex. 😭 → “loudly crying face”)

Remove stopwords

We removed NLTK (Natural Language Toolkit)'s default list of English stopwords and some unmeaningful punctuations like “.” “,”

Regex cleaning

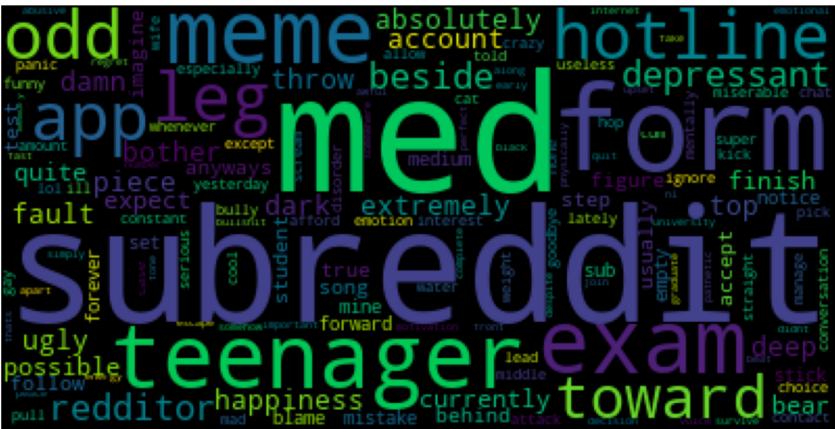
‘_WORK_’: jobs?|career|intern(ship)?|position
‘_INTERPERSONAL_’: co(\-)?worker|interpersonal|managers?|boss|supervisors|colleague|employees?|staffs?|network

Lemmatization

We used lemmatization over stemming because it converts words to their meaningful base forms. (ex. ‘Caring’ -> ‘Care’ rather than ‘Car’)

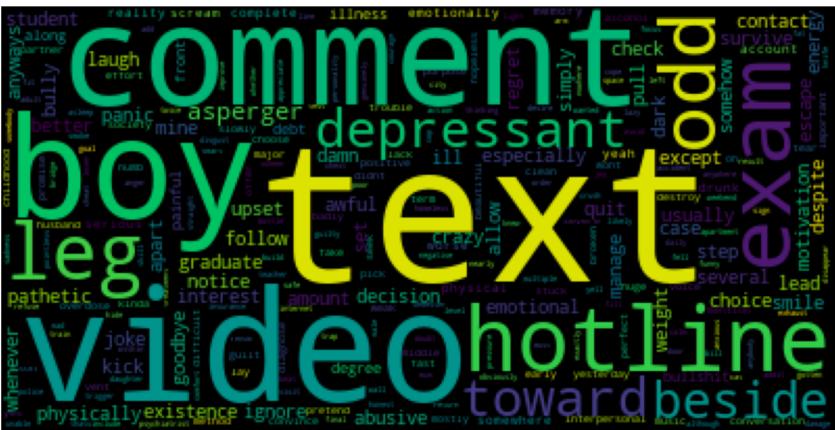
Word Cloud by all people after Pre-processing

- subreddit
 - med
 - teenager
 - form
 - meme



Word Cloud by Suicide People after Pre-processing

- text
 - comment
 - video
 - boy
 - hotline



Topic Modeling

Bigram

non-suicide	suicide
high school	want die
want die	like shit
suicidal thoughts	want end
want live	want kill
feel better	need help
...	...

Non-suicide: some negative expression but want live
Suicide: nearly all negative words

Trigram

non-suicide	suicide
smile face sunglass	wan na talk
cool smiling face	wan na chat
loudly crying face	bored wan na
face rolling eyes	loudly crying face
face tears joy	wan na die
...	...

Non-suicide: lots of positive emojis
Suicide: negative emojis, desire of communication

Modeling & Performance Evaluation

Image Source:
[GoodData](#)



Model Exploration & Architecture

1. Bert Tokenizer
2. Train/Test Sets 50-50 Split
3. Model Exploration
 - a. Logistic (baseline)
 - b. Random Forest
 - c. LightGBM
 - d. XGBoost
 - e. CatBoost
 - f. Bert



Logistic Regression

In this algorithm, the probabilities describing the possible outcomes of a single trial are modelled using a logistic function.



Random Forest

A meta-estimator that fits a number of decision trees on various sub-samples of datasets and uses average to improve the predictive accuracy of the model.



LightGBM

A fast, distributed, high-performance gradient boosting framework based on decision tree algorithm, used for ranking, classification, etc.



XGBoost

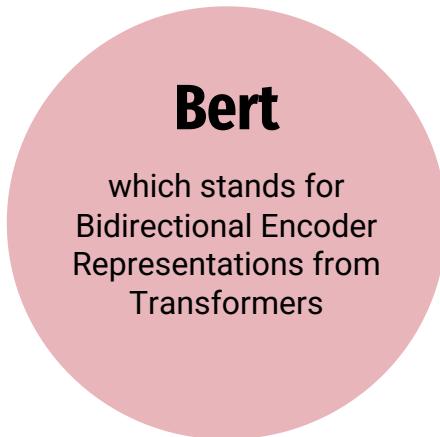
A scalable, distributed gradient-boosted decision tree (GBDT) machine learning library that provides parallel tree boosting.



CatBoost

It provides a gradient boosting framework which among other features attempts to solve for Categorical features using a permutation driven alternative compared to the classical algorithm.

Model Exploration - Bert



- 01
- 02
- 03
- 04

BertTokenizer

Use BertTokenizer from pre-trained and set truncation, MAX_LEN = 50, and batch_size = 32

Bert Classification

Use BertForSequenceClassification from pre-trained

Optimizer & Parameters

Optimizer: AdamW(lr=2e-5, correct_bias=False)

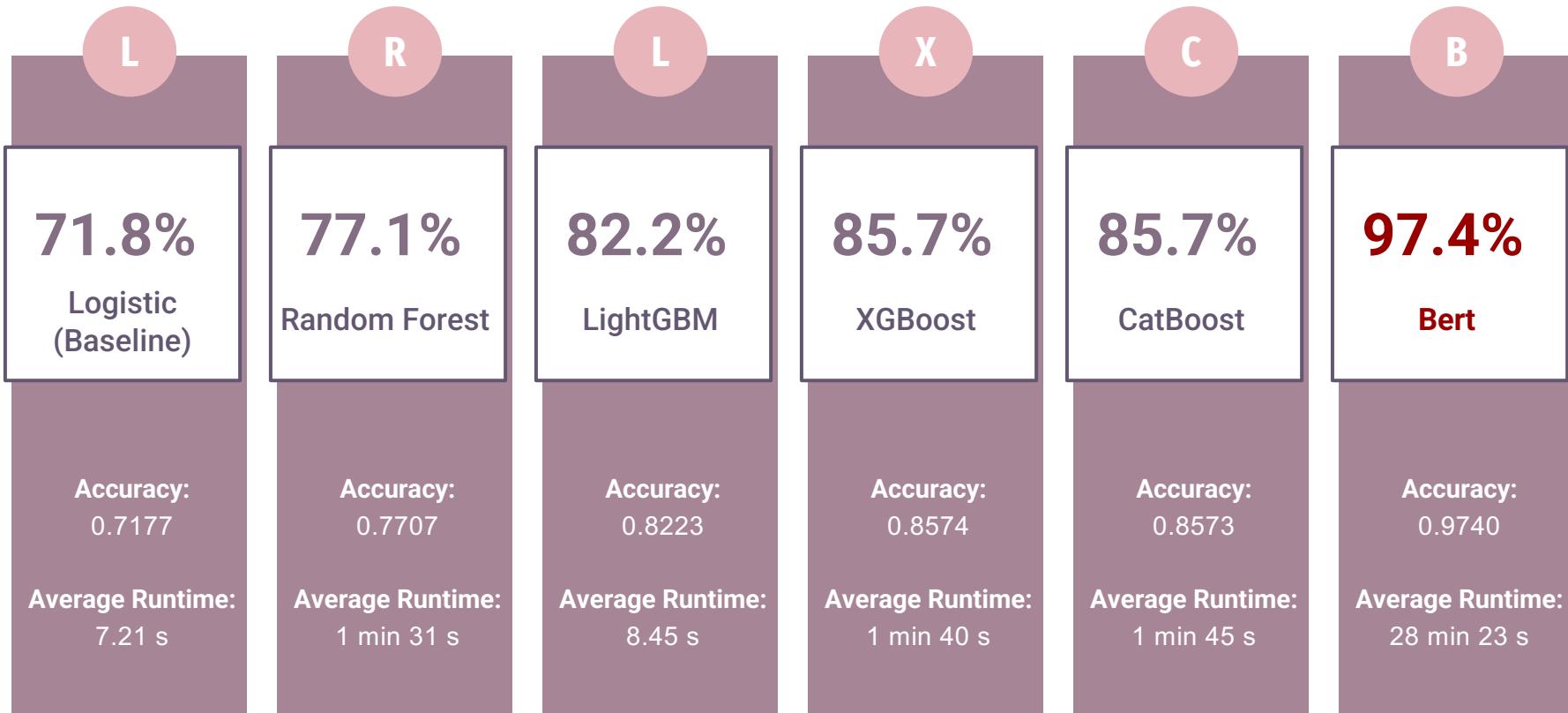
AdamW: a variant of the optimizer Adam that has an improved implementation of weight decay which is a form of regularization to lower the chance of overfitting

Parameters (109,483,778 trainable parameters)
weight_decay_rate

Evaluation Function

Accuracy & Loss

Model Performance



Business Recommendations



Image Source:
[Make a Difference](#)

Model Implementation Overview

ML Model

Our model gives high accuracy predicting potential suicide (XGBoost & CatBoost: 85.7% R²)



Suicide Detection

Implement an intelligent suicide detection product using NLP and ML classification methods



Business Mental Health Program



Financial Benefits

Increases working efficiencies of employees and saves medical costs for the company



Safe Workplaces

Boosting the happiness of employees establishes positive company image in the long run

Implementation Roadmap

01

Suicide Prediction Tool Use

As Facebook is using AI in the realm of suicide prevention, many tech companies are also looking for novel tools for predicting suicide. Our Suicide Prediction Tool can mine online posts(social media, blogs, etc.) or messages from workplace communication tools for words linked to suicide and identify flag behaviors or emotions to alert the companies of their employees' mental health issues.

02

Workplace Suicide Prevention Program

Suicide prediction tool can be incorporated to the Workplace Suicide Prevention Program in companies. Once a flagging suicidal reason is identified by the tool, the program will promote help-seeking and help-giving to the targeted employee at the workplace.

03

Improvement & Concern

Since we build our tool based on posts on Reddit that are not only posted by employed people, we would like to collect more employee data in the future and rebuild our tool. Also, using online data for suicide prediction is still subject to oversight and review to ensure its effectiveness, safety, and ethical permissibility.

Return on Investment

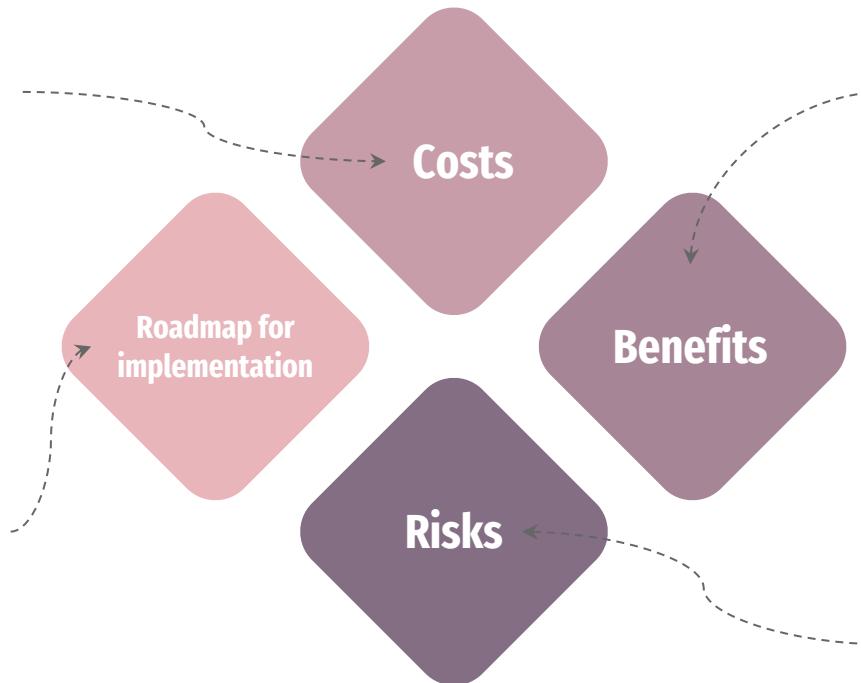
Image Source:
[Talk space](#)



Return on Investment

Enough data on the cost of inpatient psychiatric care for suicidality, police and EMS response, crisis institutes, and so on. Allocation of human resources that includes data cleaning and time for model making.

Implement suicide prediction tool in the company and incorporate it to Workplace Suicide Prevention Program. Once high risk people are identified, provide follow-ups or cognitive behavioral therapy. The researchers found that both interventions could be cost-effective as long as the models used to predict suicide risk have a high degree of accuracy.



Based on existing research literature, we estimate that each year of investment in the suicide prediction tool will at least prevent 3,600 suicide attempts and 140 deaths over the next 28 years. This leads to a reduction of 0.13% in the number of suicide attempts and deaths. We project the financial benefits of stopping these attempts and deaths to be \$1,100 per dollar invested in the tool. These benefits include savings in medical costs and increased profits.

The results may not be applicable to other geographies or populations, and the correlations described by the results are only applicable to the particular data sets used within the specified location.

Conclusion

Image Source:
Tesis Master



Conclusion

- Advanced NLP techniques can create business opportunities and identify potential suicidal actions at workplaces.
- Bert Model gives the highest classification accuracy of 97.4%, but it also has the longest runtime of 28 min 23 s. This could be implemented into businesses' mental health program or a suicide detection product that can be widely used.
- By looking at the result of model performance, we can promote help-seeking and reduce the incidence of suicide and suggest suicide-prevention campaigns.
- Suicide prediction can result in significant potential cost savings as a result of fewer suicide deaths and reduced life years lost.

Future Improvement

Image Source:
[accounting web](http://accountingweb)



Future Improvement

More Model Explorations

Tune Parameters

Automatic Pipelines

1

For **efficiency**, we would like to try fasttext, DistilBert, FastBert and so on.

For **accuracy**, we would like to try more advanced models which are improved based on Bert such as RoBerta, DeBerta-V3. Finally, we would like to find the balance (Trade-offs) between efficiency and accuracy.

2

We did not have enough time to tune parameters. Hence, we would like to tune parameters to control running time, limit overfitting, and improve accuracy in the future.

3

Functionalize things, automate workflows and create pipelines to deal with streaming data.

Reference

<https://analyticsindiamag.com/7-types-classification-algorithms/>

<https://towardsdatascience.com/lightgbm-vs-xgboost-which-algorithm-win-the-race-1ff7dd4917d>

<https://www.nvidia.com/en-us/glossary/data-science/xgboost/>

<https://en.wikipedia.org/wiki/Catboost>

<https://healthitanalytics.com/news/suicide-risk-prediction-models-prove-cost-effective-in-healthcare>

<https://www.rand.org/pubs/periodicals/health-quarterly/issues/v5/n2/09.html>