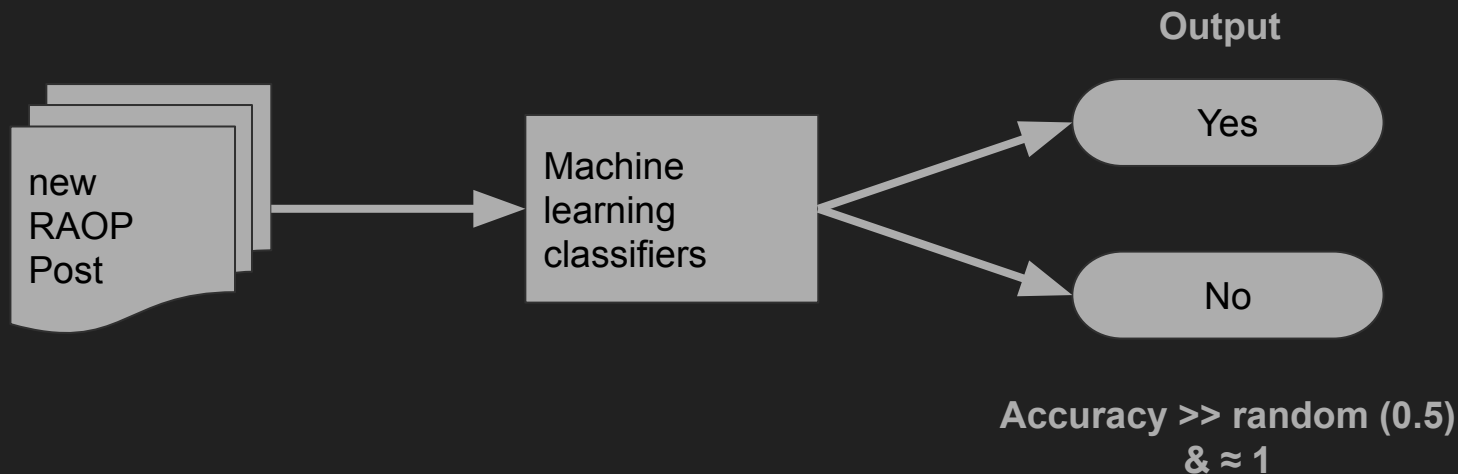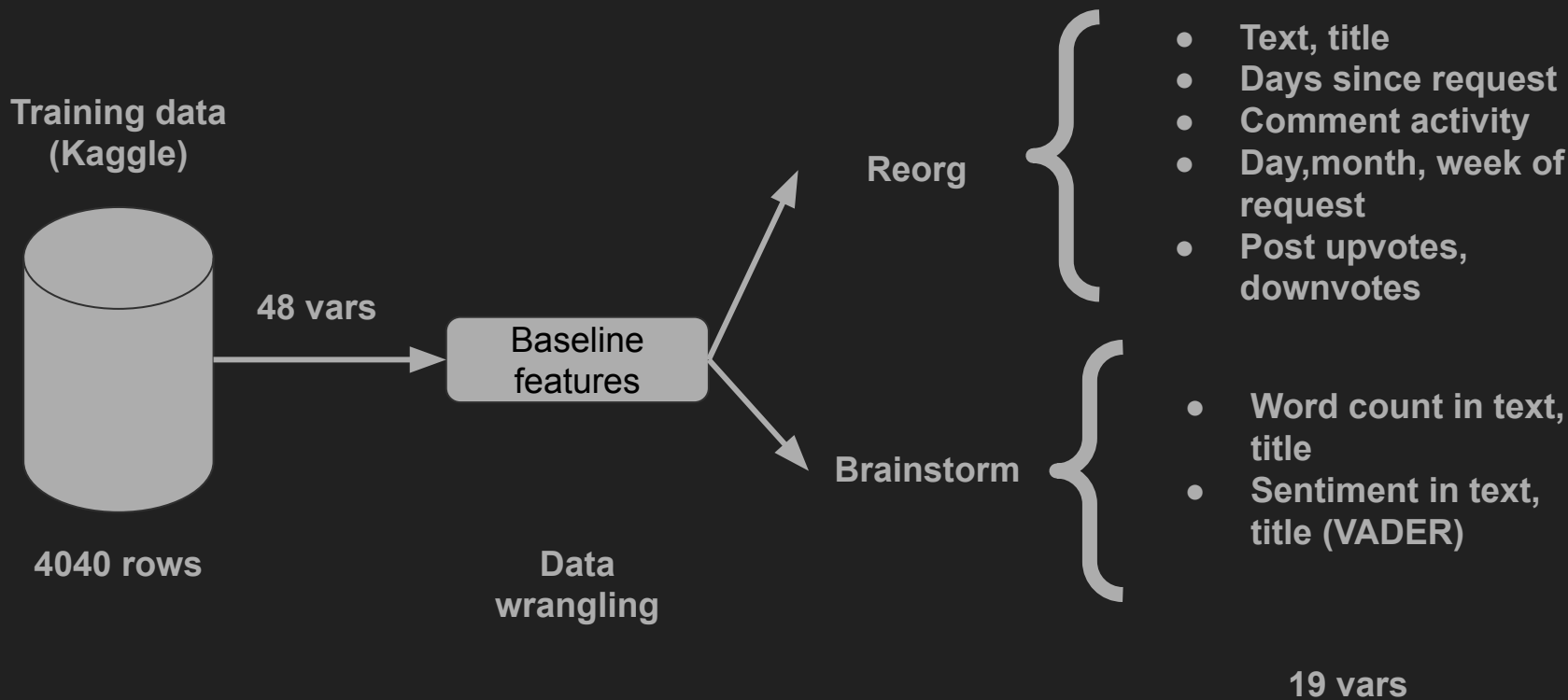# Random Acts of Pizza

Kanika Mahajan, Carolina Arriaga
W207 Applied Machine Learning

# Problem description

In the Random Acts of Pizza challenge (Kaggle competition) we are asked to **predict** if a message posted on Reddit group RAOP will get a pizza or not.

Output

new RAOP Post → Machine learning classifiers → Yes / No

**Accuracy >> random (0.5) & ≈ 1**
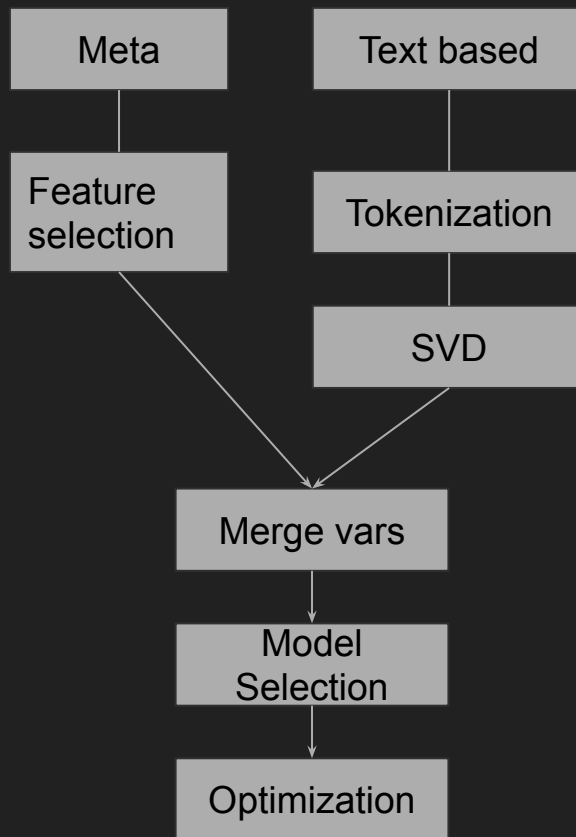
# Data description

# Challenges

- The initial data was unbalanced.
  - 24% of the posts received pizza
  - Classifiers might be biased towards predicting "False".
- We had a small data set (4,040 rows)
  - Test data not available (Kaggle competition)
  - Split in three ways -> even smaller dataset (75-15-10)
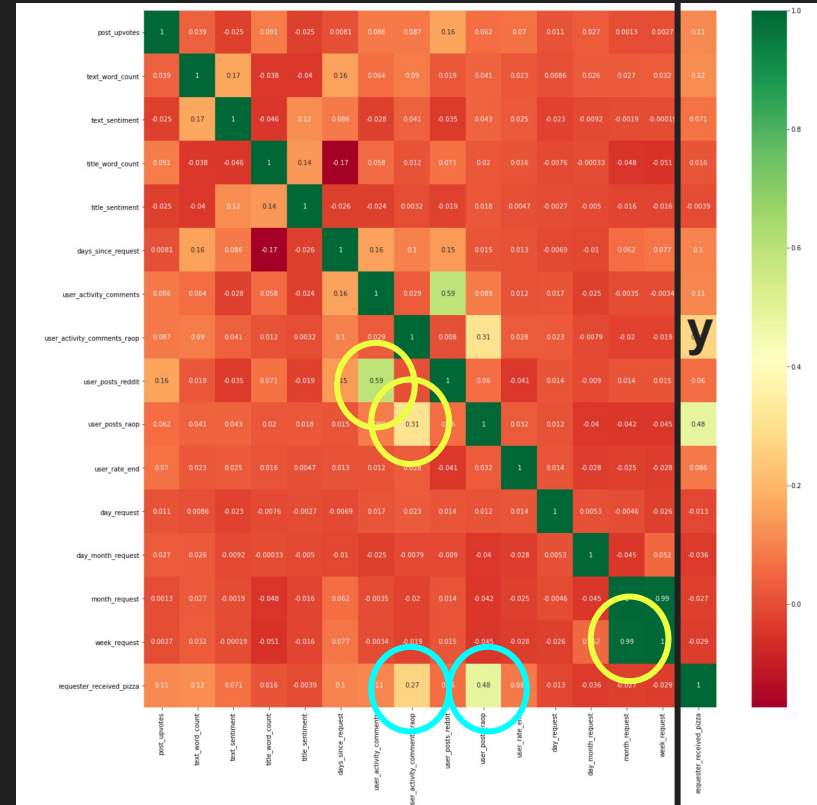- Team objective:

  get accuracy > 90%

# Our Approach

- Feature engineering
  - Meta - related to user or our understanding of language (e.g. post upvotes, activity)
  - Text based features - word counts, sentiment, TFIDF
  - Matrix correlation
  - Dimensionality reduction - SVD
- Model selection
  - Logistic regression
  - Random Forests (RF, Ada, XGB)
  - Multilayer Neural Network
  - Dense Neural Network
- Hyperparameter optimization

```
┌──────────┐        ┌──────────┐
│   Meta   │        │Text based│
└────┬─────┘        └────┬─────┘
     │                   │
┌────┴─────┐        ┌────┴──────┐
│ Feature  │        │Tokenization│
│selection │        └────┬──────┘
└────┬─────┘             │
     │              ┌────┴─────┐
     │              │   SVD    │
     │              └────┬─────┘
     └──────┐   ┌────────┘
         ┌──┴───┴───┐
         │Merge vars│
         └────┬─────┘
         ┌────┴─────┐
         │  Model   │
         │Selection │
         └────┬─────┘
         ┌────┴─────┐
         │Optimization│
         └──────────┘
```

# Feature Selection - Meta

Finding features that are independent
from each other to avoid multicollinearity.

- Highly correlated covariates:
  - Week - Month
  - Comments - Posts on Reddit
  - Comments raop - Posts on raop
- Removed week of request, posts
  on reddit, comments on raop.
- High correlation between output
  variable posts on raop (blue) - good
  variables.



Correlation matrix

# Feature engineering - TFIDF + SVD
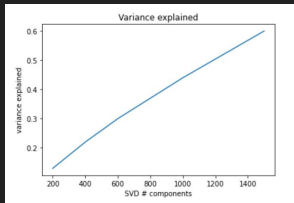
We tried several approaches:

- Only tf-idf, word count (using token pattern)
- Better processor:
  - Removes special characters
  - Lemmatization and stemming
  - Remove stop words
- Used ngrams
  - Range (1-4)
- SVD

Tokenization

Up to 300k ngrams

SVD

Tested with: 200, 400, 600, 1000, 1500



Model Eval

Optimize Baseline Logistic regression

Random forest (XGB, Ada, RF)

# Model Overview - Baseline

**Logistic Regression**

- Multiple componentes = 1500
- L2 Regularization

**Trees - (DT, RF, Ada)**

- Estimators=20
- Depth=4

**7-NN**

- Tried [1,3,5,7,9,11]

Accuracy (Dev set)

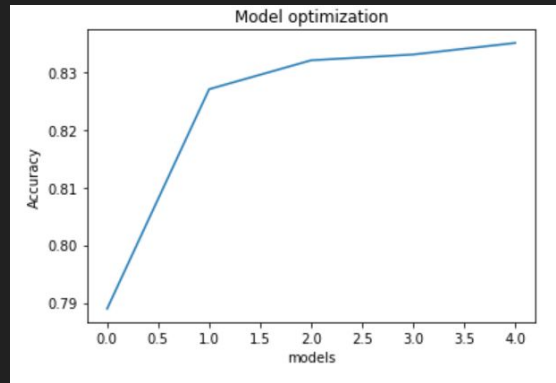| Model | Accuracy |
|-------|----------|
| Logit | 84% |
| DT | 75% |
| RF | 82% |
| Ada | 83% |
| 7-NN | 72% |

# Model Overview - Optimization

**Logistic Regression**

- **SVD componentes** = [**200**, 400, 600, 1000, 1500]
- Multiple smoothing values **C**=[0.01,0.1,**1**,10,100]
- L2 Regularization

**Trees - (XGB, Ada, RF)**

- Multiple estimators [15, 30, 80, 100, 120, **250**]
- Multiple depths [2,3,4,5,6,**12**]



Accuracy (Dev set)

| Logit | XGB |
| --- | --- |
| 84.5% | 83% |

1% improvement

# Model Overview - New models

Accuracy (Dev set)

**Multilayer Neural Network**

- Used MLPClassifier with dimension (216,216,216) - our best result of components.

| MLNN | DNN |
|------|-----|
| 82.17% | 92.17% |

+7% improvement

**Dense Neural Network**

- Included 3 sigmoid layers
  - 2 of size 216
  - Third layer of size 1
  - Sigmoid

```
Layer (type)                    Output Shape                Param #
=================================================================
dense (Dense)                   (None, 216)                 46872

dense_1 (Dense)                 (None, 216)                 46872

dense_2 (Dense)                 (None, 1)                   217
=================================================================
Total params: 93,961
Trainable params: 93,961
Non-trainable params: 0
```

# Model Overview - New model optimization

**Dense Neural Network**

- Considered 1512 SVD (~60% variance)
- Stochastic Gradient Descent- batch size = 1
- Mini_batches = [20,60,100,150]
- 2 Layers
- Units = [200 & 500] each layer
- Epochs = 100

Dev set

| Model | Accuracy |
|---|---|
| Stochastic GD | 86% |
| Mini batch (20) | 97% |
| Mini batch (60) | 98% |
| Mini batch (150) | 99% |

+7% improvement

# DNN - Best model

```
Model: "sequential"
_____
Layer (type)                 Output Shape              Param #
=================================================================
dense (Dense)                (None, 200)               302600
_____
dense_1 (Dense)              (None, 200)               40200
_____
dropout (Dropout)            (None, 200)               0
_____
dense_2 (Dense)              (None, 1)                 201
=================================================================
Total params: 343,001
Trainable params: 343,001
Non-trainable params: 0
_____
Accuracy (DNN): 0.994 - batch=150
```

Increased parameters

~ 3.6x

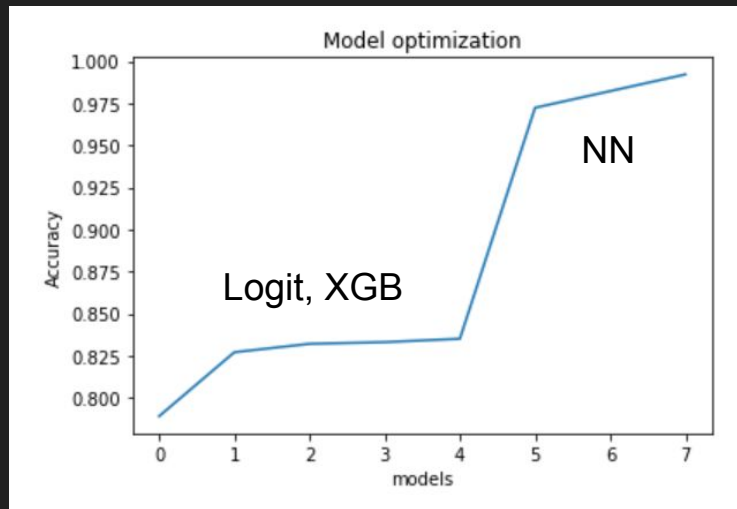Increased accuracy to 99%

Added Noise (dropout=0.5)

Mini batch GD

# Optimized parameters - Best models

| Model | Best hyperparameters | Accuracy |
|---|---|---|
| Logit | C = 1, L2, SVD=1512 | 85% |
| XGB | Depth=4<br>Estimators=120<br>SVD=1512 | 84% |
| DNN | SVD=1512 (input)<br>Sequential<br>L1= 200 comps & Sigmoid<br>L2= 200 comps & Sigmoid<br>Dropout = 0.5<br>L3= 1 comps & Sigmoid<br><br>Batch size=150<br>Epochs = 100<br><br>Optimizer=Adam | 99% |

# Results

We tested our optimized models with the validation data.



Accuracy

| Model | Dev | Test |
|---|---|---|
| Logit | 86% | 82% |
| XGB | 84% | 83% |
| DNN mini batch (150) | 98% | 99% |

# What worked

- Removing covariates with high correlation.
- SVD components ~ hundreds.
- Logistic regression was slightly better than forests, also faster.
- Dense Neural Network predicted the best results.

# What didn't work

- Baseline model
  - KNN model, accuracy was much lower compared to RF and Logistic Regression.
- Multilayer Neural Network
  - Didn't do better than trees or logit.
  - Didn't explore further.
- Dense Neural Network
  - We tried different activation layers (relu, softmax) - sigmoid was the best

# Lessons learned

- We could have created an ensemble
  - Include multiple model predictions
- Some models were running well in a computer but due to computing power not so well in other.
  - We could have created a container
- Compiling takes a long time
  - Need to find ways to optimize computing time.