# CONDESCENSION DETECTION *for*
## STACKOVERFLOW

By Amelia Lui

# TABLE of CONTENTS

# INTRODUCTION

- **What is StackOverflow?**
- **Motivation**: A kinder, more productive learning experience
- **Goal**: Classify a comment as condescending or not condescending

## 02

# METHODOLOGY

- **Data**: <u>Stack Overflow Data</u> from Kaggle
- **Tools**: Numpy, Pandas, Matplotlib, Seaborn, Sklearn, Vader, nltk, scipy, gensim
- **Topic Modeling**: LSA, NMF
- **Classification**: kNN, Logistic Regression, Random Forests
- **Model Evaluation**: emphasis on recall, but with context of confusion matrix

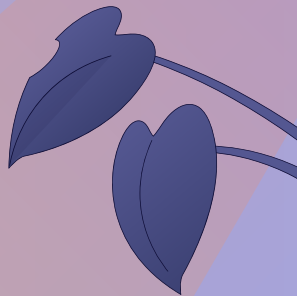# RESULTS: *topics*

**Latent Semantic Model**

**TOPIC ONE**
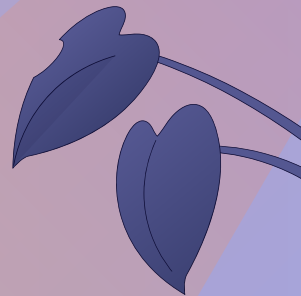*use, code, would, question, answer, like, think, c, also, need, want, using, time, get, work, example, could, mean, say, see, know, much*

**TOPIC TWO**
*jpeg, ocaml, postgresql, words, case-insensitive, -tiers, n-tiers, associative, foo, controlchars.quote, age, href=, cryptography*

**TOPIC THREE**
*asked, good, perfectly, accepted, answered, correct, subjective, post, voted, help, yes, wrong, upvote, google, original*

# RESULTS: *topics*

**Non-negative Matrix Factorization**

**TOPIC ONE**

mean, string, new, example, different, best, statement, name, read, work, say, get, words, foo, syntax, var, f, x, file, variable, upvote, phrase

**TOPIC TWO**

code, use, would, think, like, c, need, time, way, want, could, good, method, get, function, c++, type, object, say, example, better

**TOPIC THREE**

question, answer, asked, good, one, answers, would, valid, c++, ask, vote, perfectly, accepted, nice, different, wrong, better.

# RESULTS: *topics*
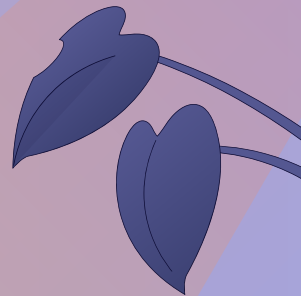
## Non-negative Matrix Factorization

**TOPIC ONE**

mean, string, new, example, different, best, statement, name, read, work, say, get, words, foo, syntax, var, f, x, file, variable, upvote, phrase

**TOPIC TWO**

code, use, would, think, like, c, need, time, way, want, could, good, method, get, function, c++, type, object, say, example, better

**TOPIC THREE**

question, answer, asked, good, one, answers, would, valid, c++, ask, vote, perfectly, accepted, nice, different, wrong, better.

# RESULTS: *topics*

## Final Model: Latent Semantic Model w/ TF-IDF

### 1: Further Questions

*use, code, would, question, answer, like, think, c, also, need, want, using, time, get, work, example, could, mean, say, see, know, much*

### 2: Technical

*jpeg, ocaml, postgresql, words, case-insensitive, -tiers, n-tiers, associative, foo, controlchars.quote, age, href=, cryptography*

### 3: Reviewing Comments

*asked, good, perfectly, accepted, answered, correct, subjective, post, voted, help, yes, wrong, upvote, google, original*

# RESULTS: *classification*
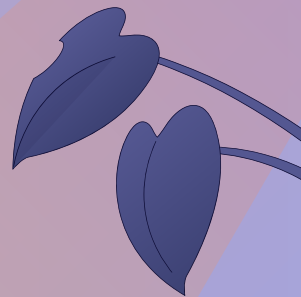
**Baseline Model**

Let this model predict the *majority class* every time
**accuracy score**: 89.07%
**precision score**: 0.00%
**recall score**:  0.00%
**f1 score**: 0.00%

# RESULTS: *classification*

## Final Model

kNN (k=3), with upsampled, scaled data and lower decision threshold
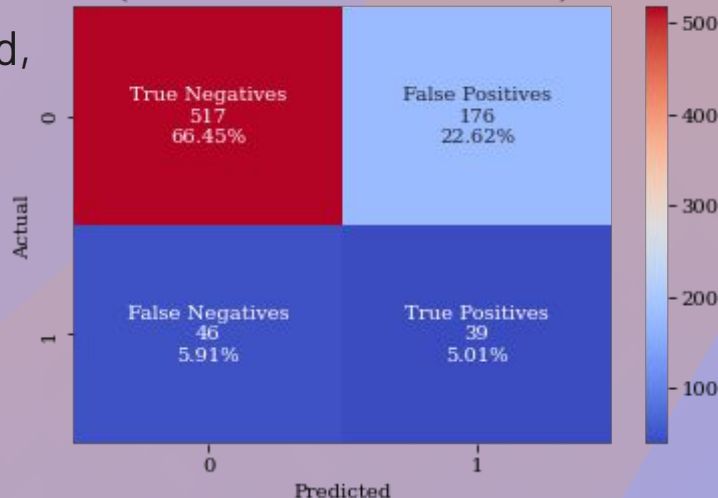
**accuracy score**: 71.47%

**precision score**: 18.14%

**recall score**: 45.88%

**f1 score**: 26.00%

### Confusion Matrix: kNN (lower decision threshold)

| | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | True Negatives 517 66.45% | False Positives 176 22.62% |
| Actual 1 | False Negatives 46 5.91% | True Positives 39 5.01% |

# CONCLUSION

- **Application**
  - Place warning to choose words more kindly if comment detected as condescending
- **Further Work**:
  - Get rid of technical terms
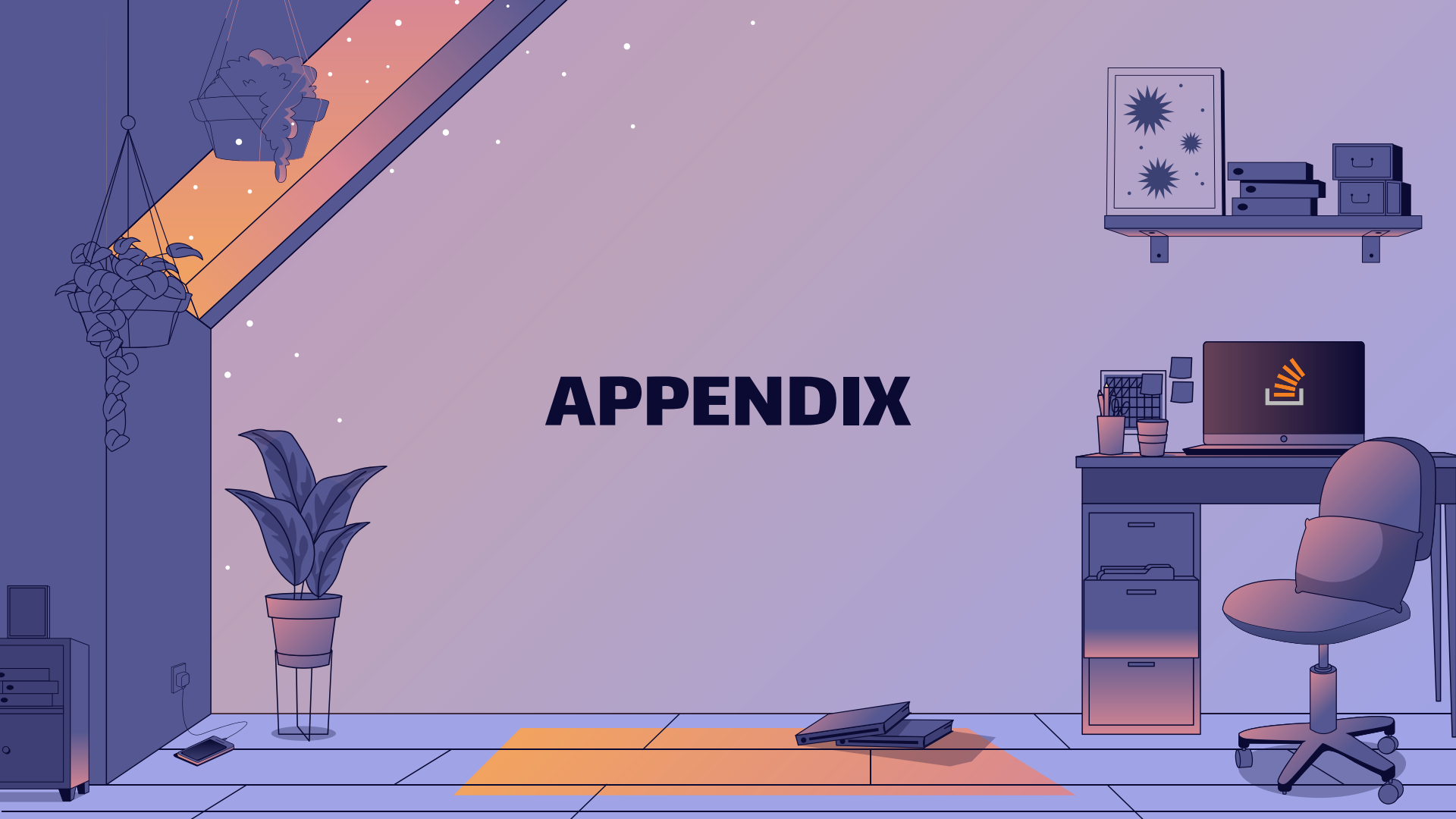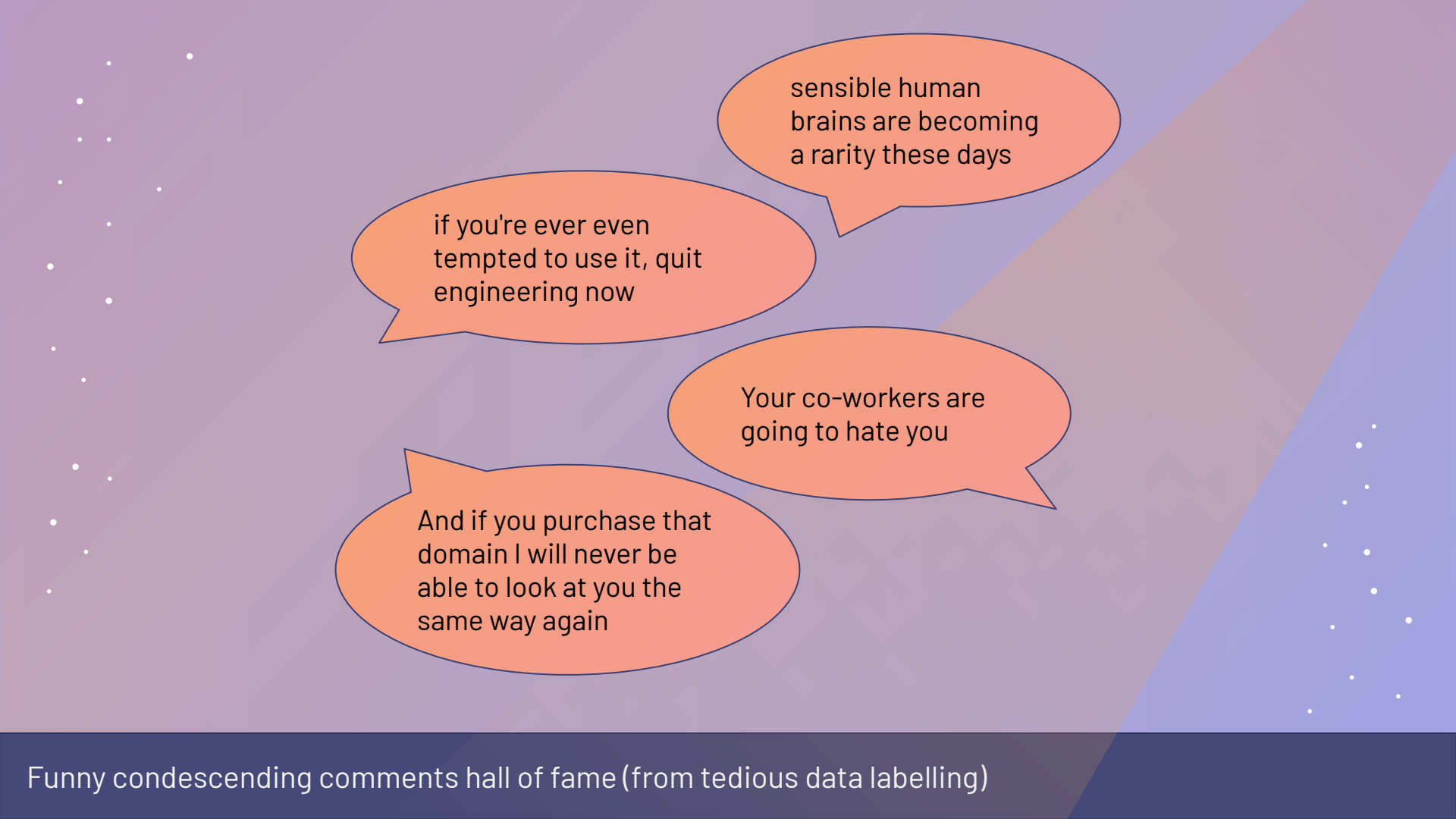  - Look at other classification model
  - More data

# Thank You!

Slides by Slidesgo

# APPENDIX

sensible human brains are becoming a rarity these days

if you're ever even tempted to use it, quit engineering now

Your co-workers are going to hate you

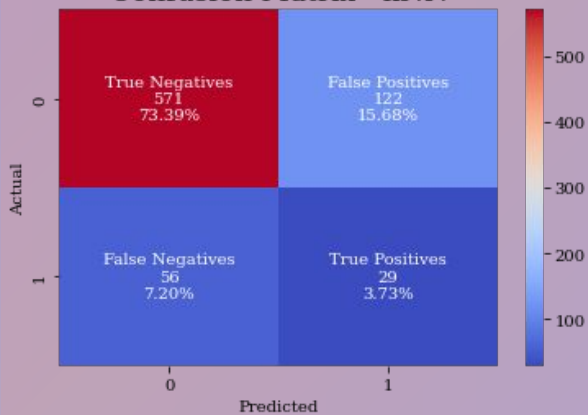And if you purchase that domain I will never be able to look at you the same way again

Funny condescending comments hall of fame (from tedious data labelling)

| kNN (k=3) | Logistic Regression | Random Forests |
|---|---|---|
| **accuracy score**: 74.55% | **accuracy score**: 10.93% | **accuracy score**: 86.25% |
| **precision score**: 14.01% | **precision score**: 10.93% | **precision score**: 15.62% |
| **recall score**: 25.88% | **recall score**: 100.00% | **recall score**: 5.88% |
| **f1 score**: 18.18% | **f1 score**: 19.70% | **f1 score**: 8.55% |

Confusion Matrix - kNN

| | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | True Negatives 571 73.39% | False Positives 122 15.68% |
| Actual 1 | False Negatives 56 7.20% | True Positives 29 3.73% |

Confusion Matrix - LogReg

| | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | True Negatives 0 0.00% | False Positives 693 89.07% |
| Actual 1 | False Negatives 0 0.00% | True Positives 85 10.93% |

Confusion Matrix - Random Forests

| | Predicted 0 | Predicted 1 |
|---|---|---|
| Actual 0 | True Negatives 666 85.60% | False Positives 27 3.47% |
| Actual 1 | False Negatives 80 10.28% | True Positives 5 0.64% |

For easy comparison of classification models performance (before lowering decision threshold)

**kNN (k=3)**

**accuracy score**: 71.47%

**precision score**: 18.14%

**recall score**: 45.88%

**f1 score**: 26.00%

Final model

Confusion Matrix: kNN
(lower decision threshold)

**Logistic Regression**
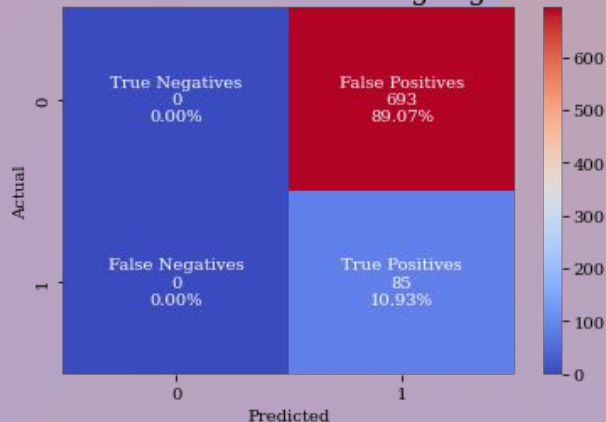
**accuracy score**: 10.93%

**precision score**: 10.93%

**recall score**: 100.00%

**f1 score**: 19.70%
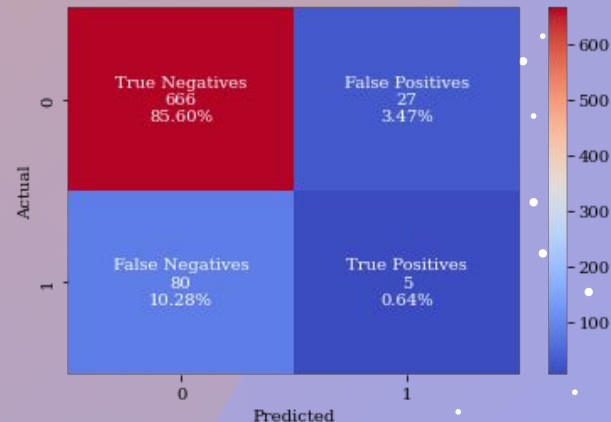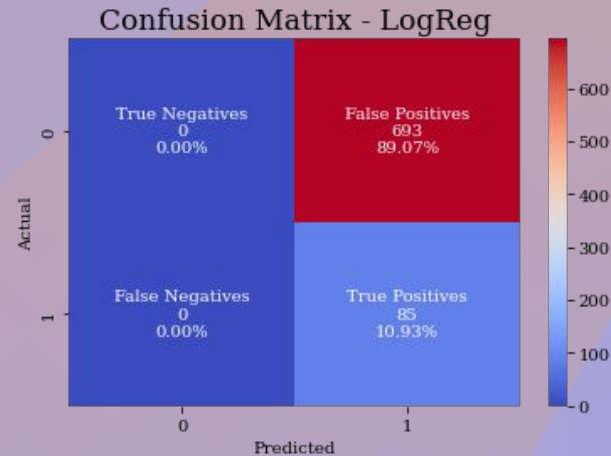
Same results

Confusion Matrix - LogReg

For easy comparison of classification models performance (after lowering decision threshold)