# Predicting Mental Health Treatment at US Companies

*Supervised Machine Learning Report*
Esther Weon (esther.s.weon@gmail.com)

# OVERVIEW

# 1. INTRODUCTION

## OBJECTIVE:

Analyze mental health survey responses to find predictors for whether American employees will seek mental health treatment

## DATA (source):

- Global 2014 – 2016 mental health survey – majority in US
  - **1259** total responses (global)
  - **751** USA responses (~60% of total)
  - **27** features
- Explores employee & employer attitudes towards mental health in workplace
  - Employee demographic info; mental health history & treatment
  - Employer-offered mental health benefits & care options

# 2. DATASET

**FEATURES:**

- **Employee Info**
  - Demographics – Age, gender, location, type of employment
  - Write-in comments
  - Mental illness
    - Family history
    - Currently seeking treatment **(TARGET)**

- **Employer Attitudes & Work Environment Offerings**
  - Size & type of company
  - Availability of mental health benefits & care options
  - Value & prioritization of mental health in workplace

# 3. DATA CLEANING

## FEATURES TO CLEAN:

- **Age**
    - Invalid ages – below 0 & over 120
        - 5 rows dropped
    - Normalize using min-max scale

- **Gender**
    - Employees' write-in responses
    - Clean & Categorize
        - 'Male', 'Female', & 'Other' (i.e. Queer / Non-Binary) → 'Male' & 'Female'
        - 13 rows of 'Other' dropped

# FEATURES WITH MISSING VALUES:

- **state**
  - 515 → 11 rows, after focusing on just US employees
  - All 11 rows dropped

- **self_employed**
  - 18 NaN rows dropped

- **work_interfere**
  - 264 rows
  - Impute all NaNs to most common category ('Sometimes')

- **comments**
  - 1095 rows (significant percentage)
  - Save for later feature engineering

# 4. FEATURE ENGINEERING

## LOCATION & TIME FEATURES

- **Country → Continent**
  - Later removed to focus on US
- **US Regions**
  - North, South, East, West, Midwest

- **Years**
  - 2014, 2015, 2016
- **Seasons**
  - Spring, Summer, Fall, Winter

## COMPARISON FEATURES

- **mental_ vs. physical_consequence**
  - More worried about workplace consequences due to _____ health
- **mental_ vs. physical_interview**
  - More likely to mention _____ health in interview

## OTHER FEATURES

- **has_comment**
  - 0 (False) or 1 (True); majority 0
- **comment_len**
  - Char lengths of comments ; majority = 0

# 5. HANDLE SKEWED DATA

## SKEWED DATA

- Need 1:1 ratio for people seeking vs. not seeking mental health treatment
- **NOTE** – data not so skewed that skipping over-sampling would have given significantly worse scores

## OVER-SAMPLING VIA SMOTE

- Split into train & test datasets
- Over-sample the train dataset to achieve 1:1 ratio
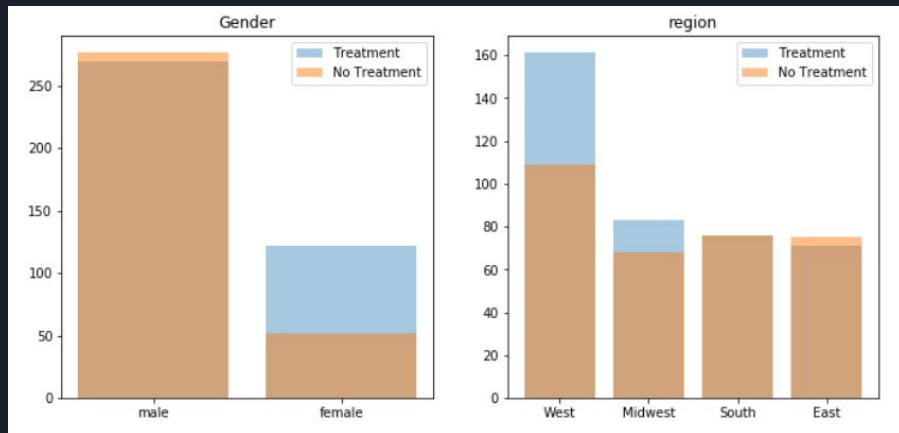    - Over-sampled train → Over-sampled X (train & test); Over-sampled Y (train & test)
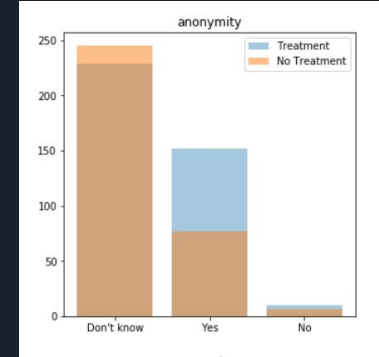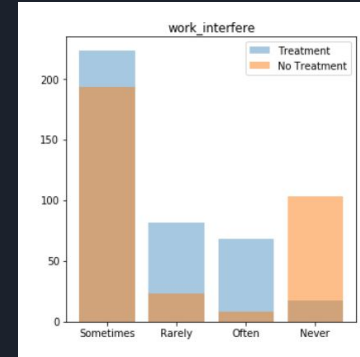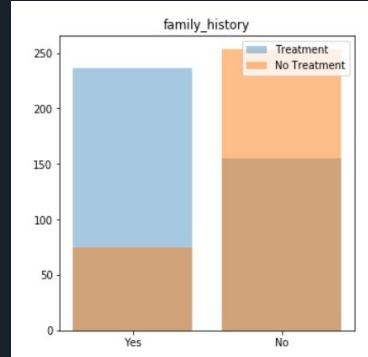
# 6. EXPLORE & ANALYZE

**NARROW FOCUS:**

- Just US employees
- Takes care of some missing rows (e.g. state)
- 2 groups to compare
  - Seeking vs. Not Seeking Treatment
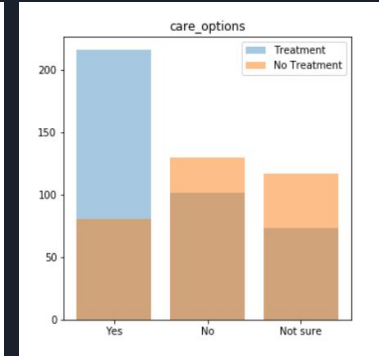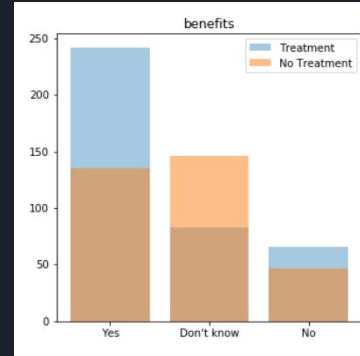
**EDA FINDINGS:**

- Gender
  - More women seeking treatment
- Region
  - More employees on West coast seeking treatment

## EDA FINDINGS, cont'd:

- **Employees seeking treatment more often have:**
  - Family history of mental illness
  - Mental health issues frequently affecting their work
  - Anonymous access to mental health help
  - Work-provided mental health benefits & care options

# 7. MODEL SELECTION

## TARGET:

Predict whether or not an employee is seeking mental health treatment

## STRATEGY:

- Create & evaluate several classifier models
    - Fit on over-sampled train
    - **TRAIN SCORE**: Score on over-sampled test
    - **TEST SCORE**: Score on test dataset
    - **WHOLE SCORE**: Score on whole dataset

## PREPARATION:

- **Normalization** of Continuous Variables
  - Age, comment_lens
- **Dummies** for Categorical Variables
  - Gender, self_employed, family_history, region, etc.

## MODELS TO BUILD:

- LASSO Logistic Regression
- Random Forest Classifier
- Naive Bayes Classifier
- Gradient Boosting Classifier
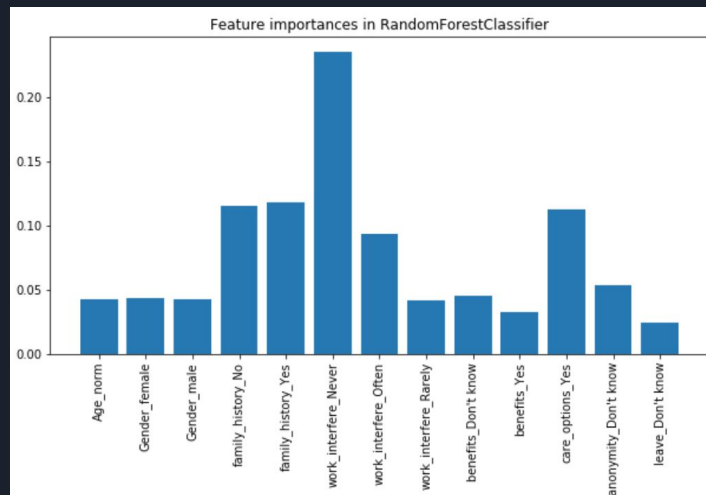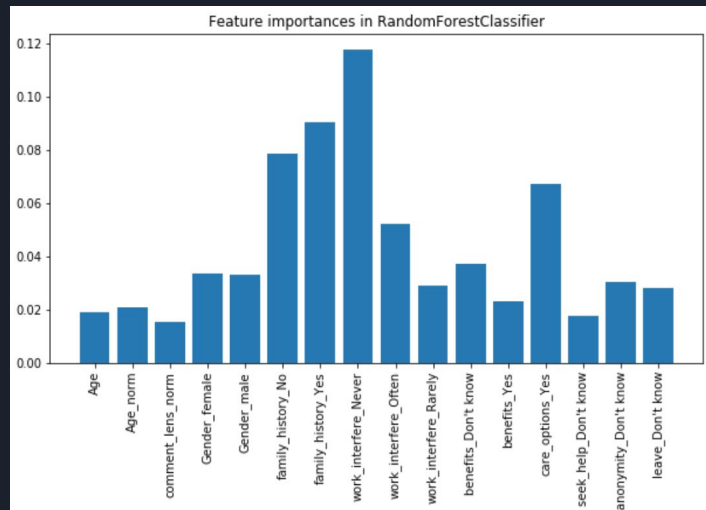- Support Vector Classifier

# 8. MODEL EVALUATION

| | LASSO Logistic Regression | Random Forest Classifier | Naive Bayes Classifier | Gradient Boosting Classifier | Support Vector Classifier |
|---|---|---|---|---|---|
| $R^2$ (train, test, whole) | 74.8% 80.8% 78.2% | 72.3% 79.0% 78.0% | 71.4% 75.3% 72.7% | 70.6% 90.3% 82.3% | 74.8% 91.9% 80.0% |
| Precision & Recall (train, test, whole) | 85.0%, 80.6% 86.0%, 77.7% 82.1%, 76.5% | 85.0%, 80.6% 82.4%, 78.2% 80.9%, 78.0% | 76.9%, 76.9% 77.3%, 77.7% 75.7%, 73.4% | 100.0%, 100.0% 92.2%, 89.9% 84.7%, 82.4% | 100.0%, 100.0% 95.2%, 89.9% 84.8%, 76.7% |
| Type I & II Errors (train, test, whole) | 6.9%, 9.5% 6.9%, 12.2% 9.0%, 12.8% | 6.9%, 9.5% 9.2%, 11.9% 10.0%, 12.0% | 11.3%. 11.3% 12.5%, 12.2% 12.8%, 14.5% | 0.0%, 0.0% 4.2%, 5.6% 8.1%, 9.6% | 0.0%, 0.0% 2.5%, 5.6% 7.5%, 12.7% |

# STRATEGY #1:
# USE PREVIOUS IMPORTANT FEATURES:

- Get most important features from previously fitted Random Forest Classifier
- Feed these features to most successful models (RFC, GBM, & SVC)
- Results
  - **RFC** – slightly lower $R^2$ & precision / recall rates; slightly higher errors (higher feature importances)
  - **GBM & SVC** – lower & more erratic $R^2$ & precision / recall rates



Feature importances in RandomForestClassifier



Feature importances in RandomForestClassifier

## COMPARISON CONCLUSION:

- Will choose original model over new model with previously selected important features
    - Slightly lower $R^2$ (insignificant difference)
    - Slightly lower precision & recall rates (insignificant difference)
    - Errors all get slightly higher

**REGULAR MODEL**

```
***TRAIN***
R² for train: 0.8103448275862069
predicted      0     1   All
actual
0            109    24   133
1             24   113   137
All          133   137   270

Type I errors: 8.89%
Type II errors: 8.89%

Precision: 82.48%
Recall: 82.48%

***TEST***
R² for test: 0.8164383561643835
predicted      0     1   All
actual
0            135    37   172
1             30   163   193
All          165   200   365

Type I errors: 10.14%
Type II errors: 8.22%

Precision: 81.5%
Recall: 84.46%

***WHOLE***
R² for whole: 0.7780821917808219
predicted      0     1   All
actual
0            252    82   334
1             80   316   396
All          332   398   730

Type I errors: 11.23%
Type II errors: 10.96%

Precision: 79.4%
Recall: 79.8%
```

**MODEL WITH PREVIOUS IMPORTANT FEATURES**

```
***TRAIN***
R² for train: 0.8017241379310345
predicted      0     1   All
actual
0            107    26   133
1             29   108   137
All          136   134   270

Type I errors: 9.63%
Type II errors: 10.74%

Precision: 80.6%
Recall: 78.83%

***TEST***
R² for test: 0.7917808219178082
predicted      0     1   All
actual
0            132    40   172
1             36   157   193
All          168   197   365

Type I errors: 10.96%
Type II errors: 9.86%

Precision: 79.7%
Recall: 81.35%

***WHOLE***
R² for whole: 0.7589041095890411
predicted      0     1   All
actual
0            247    87   334
1             89   307   396
All          336   394   730

Type I errors: 11.92%
Type II errors: 12.19%

Precision: 77.92%
Recall: 77.53%
```
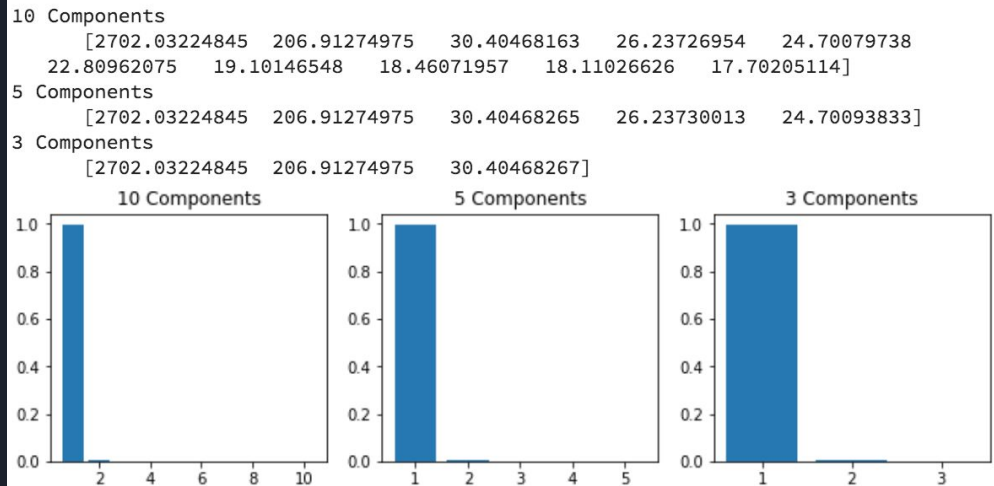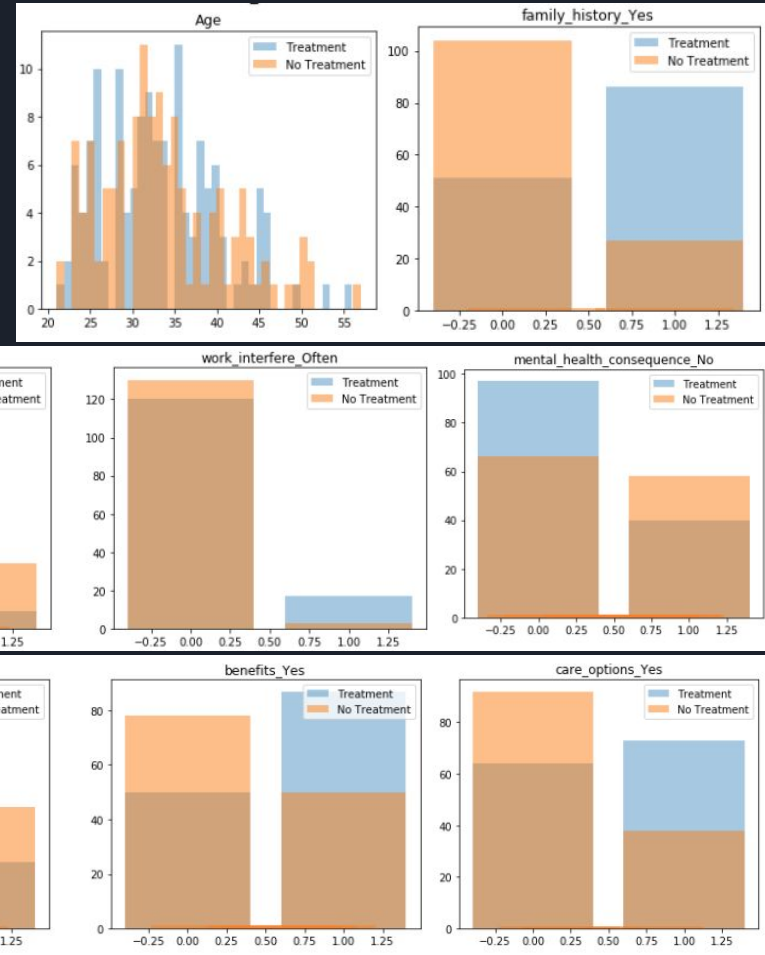
## STRATEGY #2:
## USE PCA TO REDUCE FEATURES:

- Reduced to 3 components
  - Only 1 important component
- Feed PCA features to most successful models (RFC, GBM, & SVC)
- **All models** – much lower & more erratic $R^2$, precision & recall rates (sign of overfitting)
- PCA won't work because most variables are categorical

# 9. FINDINGS



**Employees most likely to seek mental health treatment:**

- Women in late 20s to mid-30s
- Family history of mental health
- Mental health interfering with work
- Have & are aware of mental health benefits & care options at workplace
- Want to avoid mental health discussions with employer
  a. Most interesting discovery
  b. Avoid cause-effect conclusions

# 10. NEXT STEPS

**<u>Takeaways for Employers & Employees:</u>**

- More conversations ("openness") with coworkers & bosses may not be as helpful or appealing to employees seeking treatment as simple availability & awareness of mental health resources
- Understanding how untreated mental health can  interfere with work and  impact the company may incentivize employers to invest in their teams' mental health

**<u>Further Study:</u>**

- Incomplete data picture
  - Specific mental illness employees suffer from
  - Frequency, severity, & kinds of treatment being administered (talk therapy, group therapy, medical procedures, medication, etc.)

# THANK YOU

ESTHER WEON
esther.s.weon@gmail.com