# featexp

Feature exploration for supervised learning. Helps with feature understanding, identifying noisy features, feature debugging, leakage detection and model monitoring.
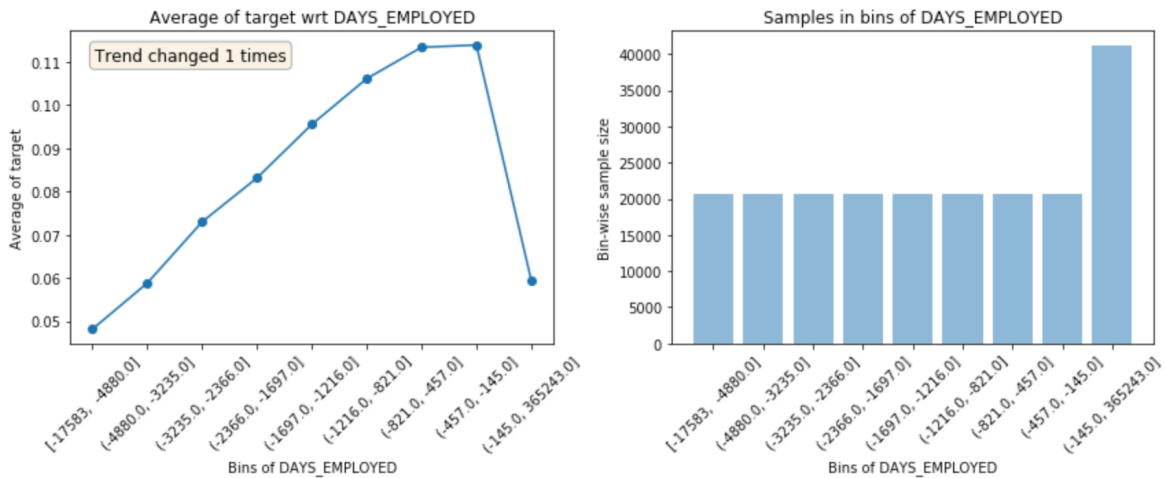
## Installation

```
pip install featexp
```

## Using featexp

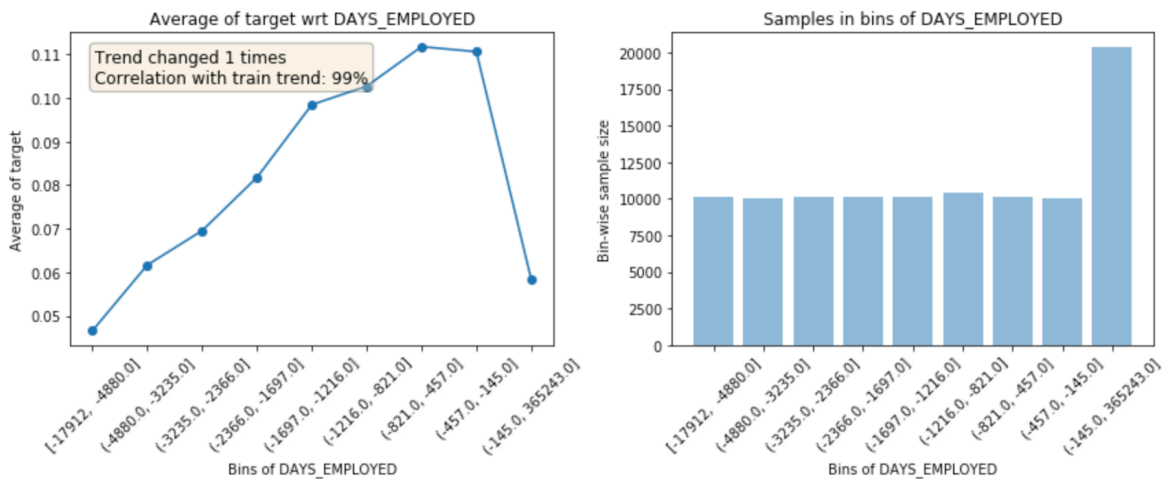Detailed Medium post on using featexp. Translations from web: Chinese, Russian

featexp draws plots similar to partial dependence plots, but directly from data instead of using a trained model like current implementations of pdp do. Since it draws plots from data directly, it helps with understanding the features well and building better ML models.

```
1  from featexp import get_univariate_plots
2  get_univariate_plots(data=data_train, target_col='target', data_test=
      data_test, features_list=['DAYS_EMPLOYED'])
3
4  # data_test and features_list are optional.
5  # Draws plots for all columns if features_list not passed
6  # Draws only train data plots if no test_data passed
```
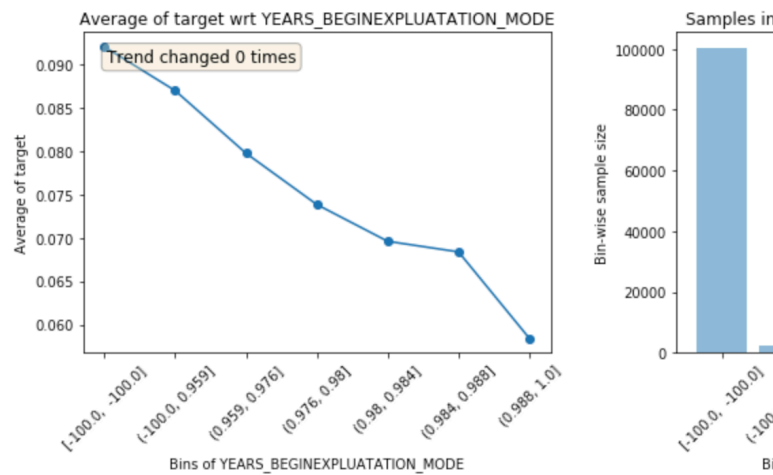
Plots for DAYS_EMPLOYED
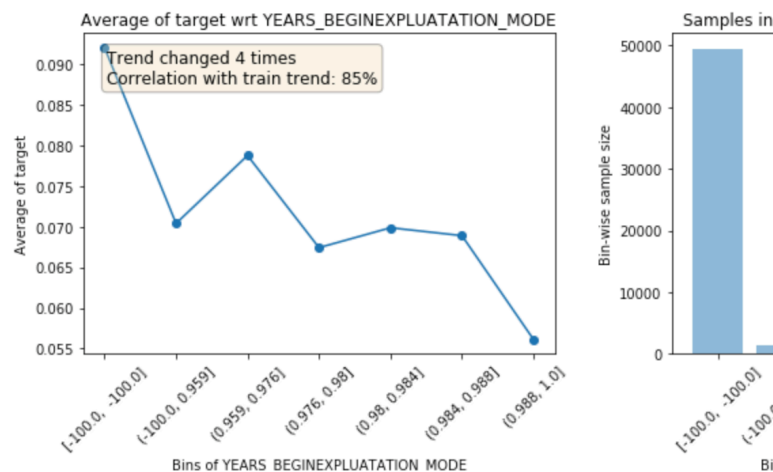Train data plots

Test data plots

featexp bins a feature into equal population bins and shows the mean value of the dependent variable (target) in each bin. Here's how to read these plots:

1. The trend plot on the left helps you understand the relationship between target and feature.
2. Population distribution helps you make sure the feature is correct.
3. Also, shows the number of trend direction changes and the correlation between train and test trend which can be used to identify noisy features. A high number of trend changes or low trend correlation implies high noise.

Example of a noisy feature: Has low trend correlation

## Getting binned feature stats

Returns mean target and population in each bin of a feature

```
1  from featexp import univariate_plotter
2  binned_data_train, binned_data_test = univariate_plotter(data=
       data_train, target_col='target', feature='DAYS_EMPLOYED', data_test=
       data_test)
3  # For only train data
4  binned_data_train = univariate_plotter(data=data_train, target_col='
       target', feature='DAYS_EMPLOYED')
```

## Getting stats for all features

Returns trend changes and trend correlation for all features in a dataframe

```python
1  from featexp import get_trend_stats_feature
2  stats = get_trend_stats(data=data_train, target_col='target', data_test
       =data_test)
3
4  # data_test is optional. If not passed, trend correlations aren't
       calculated.
```

Returns a dataframe with trend changes and trend correlation which can be used for dropping the
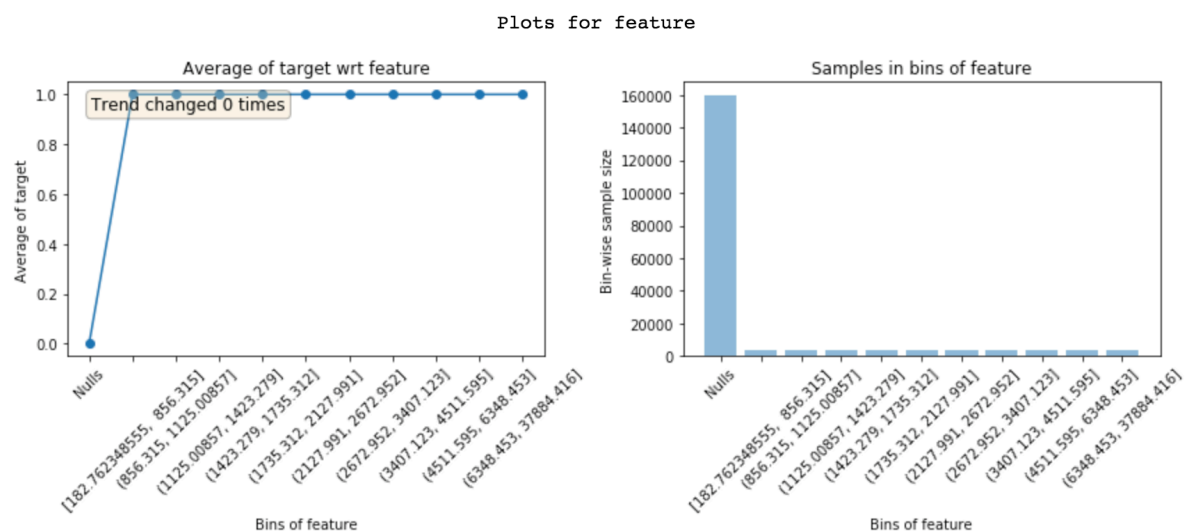
In [79]:   `stats`

Out[79]:

| | Feature | Trend_changes | Trend_changes_test | Trend_correlation |
|---|---|---|---|---|
| 0 | CNT_CHILDREN | 2 | 2 | 0.975688 |
| 1 | AMT_INCOME_TOTAL | 4 | 3 | 0.921382 |
| 2 | AMT_CREDIT | 3 | 3 | 0.988779 |
| 3 | AMT_ANNUITY | 4 | 4 | 0.972325 |
| 4 | AMT_GOODS_PRICE | 7 | 7 | 0.994683 |
| 5 | REGION_POPULATION_RELATIVE | 3 | 3 | 0.987832 |
| 6 | DAYS_BIRTH | 0 | 0 | 0.992783 |
| 7 | DAYS_EMPLOYED | 1 | 1 | 0.995426 |
| 8 | DAYS_REGISTRATION | 2 | 2 | 0.976891 |
| 9 | DAYS_ID_PUBLISH | 0 | 2 | 0.985101 |
| 10 | OWN_CAR_AGE | 1 | 1 | 0.994656 |
| 11 | FLAG_MOBIL | 0 | 0 | 0.000000 |

noisy features, etc.

## Leakage detection

It helps with identifying why a feature is leaky which helps with debugging.

Plots for feature



Nulls have 0% mean target and 100% mean target in other bins. Implies this feature is populated only for target = 1.

## Citing featexp