# kaggle

**Competitions**      **Datasets**      **Kernels**      **Forums**      **Jobs**      ▼

# Bias Correction + XGBoost

by Tilii · last run 2 days ago · Python script · 2019 views
using data from Allstate Claims Severity · 👁 Public

▲
**37**
**voters**

---

Code          This script has been released under the Apache 2.0 open source license.      **Download Code**

```python
1    import numpy as np
2    import pandas as pd
3    from datetime import datetime
4    from sklearn.preprocessing import StandardScaler
5    from sklearn.cross_validation import KFold
6    from sklearn.metrics import mean_absolute_error
7    from scipy.stats import skew, boxcox
8    from math import exp, log
9    import xgboost as xgb
10
11
12   def timer(start_time=None):
13       if not start_time:
14           start_time = datetime.now()
15           return start_time
16       elif start_time:
17           tmin, tsec = divmod((datetime.now() - start_time).total_seconds(),
18   60)
19           print(' Time taken: %i minutes and %s seconds.' %
20                   (tmin, round(tsec, 2)))
21
22
23   def scale_data(X, scaler=None):
24       if not scaler:
25           scaler = StandardScaler()
26           scaler.fit(X)
27       X = scaler.transform(X)
28       return X, scaler
29
30
```

```python
31    DATA_TRAIN_PATH = '../input/train.csv'
32    DATA_TEST_PATH = '../input/test.csv'
33
34
35    def load_data(path_train=DATA_TRAIN_PATH, path_test=DATA_TEST_PATH):
36        train_loader = pd.read_csv(path_train, dtype={'id': np.int32})
37        train = train_loader.drop(['id', 'loss'], axis=1)
38        test_loader = pd.read_csv(path_test, dtype={'id': np.int32})
39        test = test_loader.drop(['id'], axis=1)
40        ntrain = train.shape[0]
41        ntest = test.shape[0]
42        train_test = pd.concat((train, test)).reset_index(drop=True)
43        numeric_feats = train_test.dtypes[train_test.dtypes != "object"].index
44
45        # compute skew and do Box-Cox transformation
46        skewed_feats = train[numeric_feats].apply(lambda x: skew(x.dropna()))
47        print("\nSkew in numeric features:")
48        print(skewed_feats)
49        # transform features with skew > 0.25 (this can be varied to find
50    optimal value)
51        skewed_feats = skewed_feats[skewed_feats > 0.25]
52        skewed_feats = skewed_feats.index
53        for feats in skewed_feats:
54            train_test[feats] = train_test[feats] + 1
55            train_test[feats], lam = boxcox(train_test[feats])
56        features = train.columns
57        cats = [feat for feat in features if 'cat' in feat]
58        # factorize categorical features
59        for feat in cats:
60            train_test[feat] = pd.factorize(train_test[feat], sort=True)[0]
61        x_train = train_test.iloc[:ntrain, :]
```

**Code**    Output (3)    Comments (21)    Log    Versions (1)    Forks (13)                **Fork Script**

```python
63        train_test_scaled, scaler = scale_data(train_test)
64        train, _ = scale_data(x_train, scaler)
65        test, _ = scale_data(x_test, scaler)
66
67        train_labels = np.log(np.array(train_loader['loss']))
68        train_ids = train_loader['id'].values.astype(np.int32)
69        test_ids = test_loader['id'].values.astype(np.int32)
70
71        return train, train_labels, test, train_ids, test_ids
72
73    ################################# Actual Run Code
74    #################################
75
76    # enter the number of folds from xgb.cv
77    folds = 5
```

```
 78    cv_sum = 0
 79    early_stopping = 25
 80    fpred = []
 81    xgb_rounds = []
 82
 83    start_time = timer(None)
 84
 85    # Load data set and target values
 86    train, target, test, _, ids = load_data()
 87    d_train_full = xgb.DMatrix(train, label=target)
 88    d_test = xgb.DMatrix(test)
 89
 90    # set up KFold that matches xgb.cv number of folds
 91    kf = KFold(train.shape[0], n_folds=folds)
 92    for i, (train_index, test_index) in enumerate(kf):
 93        print('\n Fold %d\n' % (i + 1))
 94        X_train, X_val = train[train_index], train[test_index]
 95        y_train, y_val = target[train_index], target[test_index]
 96
 97    #########################################
 98    #
 99    # Define cross-validation variables
100    #
101    #########################################
102
103        params = {}
104        params['booster'] = 'gbtree'
105        params['objective'] = "reg:linear"
106        params['eval_metric'] = 'mae'
107        params['eta'] = 0.1
108        params['gamma'] = 0.5290
```

Code        Output (3)      Comments (21)      Log      Versions (1)      Forks (13)                    Fork Script

```
110        params['colsample_bytree'] = 0.5665
111        params['subsample'] = 0.9930
112        params['max_depth'] = 7
113        params['max_delta_step'] = 0
114        params['silent'] = 1
115        params['random_state'] = 1001
116
117        d_train = xgb.DMatrix(X_train, label=y_train)
118        d_valid = xgb.DMatrix(X_val, label=y_val)
119        watchlist = [(d_train, 'train'), (d_valid, 'eval')]
120
121    ###################################
122    #  Build Model
123    ###################################
124
125        clf = xgb.train(params,
```

```
125        clf = xgb.train(params,
126                         d_train,
127                         100000,
128                         watchlist,
129                         early_stopping_rounds=early_stopping)
130
131        ######################################
132        #   Evaluate Model and Predict
133        ######################################
134
135        xgb_rounds.append(clf.best_iteration)
136        scores_val = clf.predict(d_valid, ntree_limit=clf.best_ntree_limit)
137        cv_score = mean_absolute_error(np.exp(y_val), np.exp(scores_val))
138        print(' eval-MAE: %.6f' % cv_score)
139        y_pred = np.exp(clf.predict(d_test, ntree_limit=clf.best_ntree_limit))
140
141        ######################################
142        #   Add Predictions and Average Them
143        ######################################
144
145        if i > 0:
146            fpred = pred + y_pred
147        else:
148            fpred = y_pred
149        pred = fpred
150        cv_sum = cv_sum + cv_score
151
152    mpred = pred / folds
153    score = cv_sum / folds
154    print('\n Average eval-MAE: %.6f' % score)
155    n_rounds = int(np.mean(xgb_rounds))
```

Code    Output (3)    Comments (21)    Log    Versions (1)    Forks (13)          Fork Script

```
158    #   Make Full Dataset Predictions
159    ######################################
160
161    print('\n Training full dataset for %d rounds ...' % n_rounds)
162    watchlist = [(d_train_full, 'train')]
163    clf_full = xgb.train(
164        params, d_train_full,
165        n_rounds,
166        watchlist,
167        verbose_eval=False,)
168    y_pred_full = np.exp(clf_full.predict(d_test))
169
170    # enter the number of iterations from xgb.cv with early_stopping turned on
171    n_fixed = 376
172
```

```
172
173    nfixed = int(n_fixed * (1 + (1. / folds)))
174    print('\n Training full dataset for %d rounds ...\n' % nfixed)
175    clf_fixed = xgb.train(
176        params, d_train_full,
177        nfixed,
178        watchlist,
179        verbose_eval=False,)
180    y_pred_fixed = np.exp(clf_fixed.predict(d_test))
181    timer(start_time)
182
183    print("#\n Writing results")
184    result = pd.DataFrame(mpred, columns=['loss'])
185    result["id"] = ids
186    result = result.set_index("id")
187    print("\n %d-fold average prediction:\n" % folds)
188    print(result.head())
189    result_full = pd.DataFrame(y_pred_full, columns=['loss'])
190    result_full["id"] = ids
191    result_full = result_full.set_index("id")
192    print("\n Full dataset prediction:\n")
193    print(result_full.head())
194    result_fixed = pd.DataFrame(y_pred_fixed, columns=['loss'])
195    result_fixed["id"] = ids
196    result_fixed = result_fixed.set_index("id")
197    print("\n Full datset (at CV #iterations) prediction:\n")
198    print(result_fixed.head())
199
200    now = datetime.now()
201    score = str(round((cv_sum / folds), 6))
202    sub_file = 'submission_5fold-average-xgb_' + str(score) + '_' + str(
```

**Code**   Output (3)   Comments (21)   Log   Versions (1)   Forks (13)            Fork Script

```
205    result.to_csv(sub_file, index=True, index_label='id')
206    sub_file = 'submission_full-average-xgb_' + str(now.strftime(
207        "%Y-%m-%d-%H-%M")) + '.csv'
208    print("\n Writing submission: %s" % sub_file)
209    result_full.to_csv(sub_file, index=True, index_label='id')
210    sub_file = 'submission_full-CV-xgb_' + str(now.strftime(
       "%Y-%m-%d-%H-%M")) + '.csv'
       print("\n Writing submission: %s" % sub_file)
       result_fixed.to_csv(sub_file, index=True, index_label='id')
```

show less

Output                                                              **Download All**

## 3 files

🔲 submission_5fold-ave...

🔲 submission_full-aver...

🔲 submission_full-CV-x...

# submission_5fold-average-xgb_1146.10852_2016-10-13-02-40.csv

**Download File**

| id | loss |
|----|------|
| 4 | 1474.300048828125 |
| 6 | 2071.08837890625 |
| 9 | 7915.3203125 |
| 12 | 5771.5576171875 |
| 15 | 811.8914794921875 |
| 17 | 2246.956298828125 |
| 21 | 2125.173583984375 |
| 28 | 881.2107543945312 |

## Comments

▲
6
▼

**Tilii** · last edited 2 days ago by **Tilii**
2 days ago

This script features **Box-Cox** transformation to correct the skew in continuous variables. XGBoost is used to mop up the mess.

There are 3 different ways of making submission files. One is a simple numeric average of 5 models, each of which is trained on 80% of the data until early stopping criterion is reached. Next, full dataset

| **Code** | Output (3) | Comments (21) | Log | Versions (1) | Forks (13) | *Fork Script* |
|----------|------------|---------------|-----|--------------|------------|---------------|

Leaderboard:

5fold-average **1118.10987**

full-average 1122.27500

full-CV 1122.42771

At the time of this writing, 5-fold average score is #14 on the leaderboard.

▲
0
▼

**modkzs**
a day ago

I am little confused why Box-Cox transformation can make such improvement. Can normal distributed feature cause better xgboost performance?

**▲**
**2** **LzjtuEr**
**▼** a day ago

@modkzs, since XGB uses linear regression as the objective, the assumption of linear regression is that the features follow normal distribution.

**▲**
**5** **Tilii** · last edited a day ago by Tilii
**▼** a day ago

> I am little confused why Box-Cox transformation can make such improvement. Can normal distributed feature cause better xgboost performance?

Violations of normal distribution can increase the error in regression analysis. The skew in data is frequently "fixed" by log(N) (or log(N+1)) transformation. However, log(N) transformation is the best only when Box-Cox lambda factor is close to zero. As it is, most of the skewed columns in our dataset have lambda values that are far from zero (see below). There are better **transformations** than log(N) for different values of lambda and Box-Cox procedure finds the best one for a given dataset.

Edit: Below are lambda values of skewed features. log(N) would not be a best transformation for most of them.

---

-0.585084895197

-0.995197981798

-3.00787936887

| **Code** | Output (3) | Comments (21) | Log | Versions (1) | Forks (13) | Fork Script |
|---|---|---|---|---|---|---|

-2.2837778627

-2.27901724839

-2.56311666326

-0.613891833259

-0.437372340485

-0.47055182116

-1.44082551851

**▲**
**0** **Andrey Shkabko**
**▼**

a day ago

Training full dataset for 366 rounds ...

---

ZeroDivisionError Traceback (most recent call last) in () 162 n_rounds, 163 watchlist, --> 164 verbose_eval=False,) 165 y_pred_full = np.exp(clf_full.predict(d_test)) 166

ZeroDivisionError: integer division or modulo by zero

^^^ python 2.7

**Tilii**
a day ago

> Training full dataset for 366 rounds ...

I have no idea what causes this error. If you look at the **log** file (scroll all the way down), it runs just fine on Kaggle. The same script runs on my computer as well, and I have python 2.7.

Maybe you have one of the packages (or more) that are not recent enough. My setup is:

numpy==1.11.0rc2 pandas==0.18.0 scikit-learn==0.17.1 scipy==0.17.0 xgboost==0.4

I don't have the latest xgboost, but that shouldn't be a problem as Kaggle does, and this script works on Kaggle.

**Andrey Shkabko**
a day ago

numpy 1.11.1 pandas 0.18.1 sklearn 0.17.1 scipy 0.18.1 xgboost 0.4

**Code**     Output (3)     Comments (21)     Log     Versions (1)     Forks (13)        Fork Script

17 hours ago

Nice script, I learned a lot. One question, when I train GBMs I usually use subsample of something around .5 because it allows for stable oob error measures and because of this example: http://scikit-learn.org/stable/auto_examples/ensemble/plot_gradient_boosting_regularization.html. Admittedly, I haven't had a lot of experience with them and the example above is just anecdotal. Would you be so kind as to explain the choice of params['subsample'] = 0.9930? Thanks again

**Tony S**
15 hours ago

I've been thinking about the normality assumption that was brought up here behind the reasoning for the transformation and I think there is something more mysterious going on here.

Isn't the normality assumption in linear regression only assuming that the residual errors are normally distributed? If so, that doesn't necessarily imply that the variables themselves need to be normally

distributed. For example, a simple linear correlation such as y=2x+epsilon (a small normally distributed error) could have an almost perfect one-to-one correlation and work great in linear regression; x in this case could be any distribution, not just a normal distribution.

I'd love for someone to give me more information or correct me if I am making an incorrect assumption, but it seems to me more likely that you stumbled onto a transformation for one or more of the features that has a better linear correlation with the "loss" variable.

▲
**11**
▼

**Alexandru Papiu**
13 hours ago

Nice script and very good score :) I am however a little confused by the comments here. From what I can tell this model is not in any way related to linear regression. It is a boosted model of trees of depth 7. The `['objective'] = "reg:linear"` line is simply telling xgboost that the target is continuous.

Tree-based models do not make any assumptions on the normality of the features. In fact tree based model are **invariant** to monotonic transformations such as taking logs or box-cox. No matter how crazy the distribution of your feature is, the tree will just pick a point to split on.

Out of curiosity I actually ran the script without any log or box-cox transforms and got very similar results: 1122.44270 with full cv, 1118.89063 with 5-fold average.

I am afraid the good score here comes mostly from ensembling a bunch of well-tuned xgboosts.

▲
**5**
▼

**Faron**
12 hours ago

@Alexandru I would change it to: "In fact tree based model are **almost invariant** to monotonic transformations".

Why monotonic transforms do affect decision trees

▲
**PythonCracker**

| Code | Output (3) | Comments (21) | Log | Versions (1) | Forks (13) | Fork Script |
|------|-----------|---------------|-----|--------------|------------|-------------|

I have the same question as @AdamWolfeLevin, how did you get the final parameters like gamma=0.5290, min_child_weight=4.2922?

▲
**1**
▼

**Tilii**
7 hours ago

> explain the choice of params['subsample'] = 0.9930? I have the same question as @AdamWolfeLevin, how did you get the final parameters like gamma=0.5290, min_child_weight=4.2922?

The parameters were found using global Bayesian optimization as implemented here. See here for an example. Note that most of the code in that example is unnecessary if you run the script locally and have the bayes_opt module installed.

## Tilii
7 hours ago

> I'd love for someone to give me more information or correct me if I am making an incorrect assumption, but it seems to me more likely that you stumbled onto a transformation for one or more of the features that has a better linear correlation with the "loss" variable.

This is quite possible and I think that Alexandru added similar point of view to your argument.

## Tilii
7 hours ago

> Tree-based models do not make any assumptions on the normality of the features. In fact tree based model are invariant to monotonic transformations such as taking logs or box-cox. No matter how crazy the distribution of your feature is, the tree will just pick a point to split on.

Funny you should mention this, because I was doing Box-Cox feature manipulations as a preparation for linear models you put to such a good use in this **script**. I tried XGBoost just because it seems to be fashionable these days, and it worked. I am still exploring linear models.

> Out of curiosity I actually ran the script without any log or box-cox transforms and got very similar results: 1122.44270 with full cv, 1118.89063 with 5-fold average.

Yet another thing to try would be a gblinear booster in XGBoost.

> I am afraid the good score here comes mostly from ensembling a bunch of well-tuned xgboosts.

Though I will have a better feel after completing my linear models, this seems like a valid argument.

**Code**      Output (3)      Comments (21)      Log      Versions (1)      Forks (13)                    Fork Script

5 hours ago

> @Alexandru I would change it to: "In fact tree based model are almost invariant to monotonic transformations".

But do they always get better? Or they vary for the same reason they vary when the seed is changed?

## Tony S
5 hours ago

Thanks @Alexandru, @Faron, and @Tilii for clarification. Makes sense now. The transformation likely has little to no no impact on the gbtree algorithm (or other tree-based algorithms), but it potentially could have some impact on linear models if it provides better linear correlation to loss. I guess I have some food for thought on next steps.

**0**

**algo.ai**
3 hours ago

@Tilii

Thanks for another nice contribution!

I would like to know roughly min_child_weight = 4 in XGB is equalivant to what value in sklearn's GradientBoostingRegressor's min_samples_leaf?

**0**

**algo.ai**
3 hours ago

**1**

**Tilii**
2 hours ago

> I would like to know roughly min_child_weight = 4 in XGB is equalivant to what value in sklearn's GradientBoostingRegressor's min_samples_leaf?

I don't know of any simple conversion procedure between min_child_weight and min_samples_leaf.

It is not that difficult to tune a single parameter, if min_samples_leaf is what interests you. See **here** and **here**.

**1**

**kiddo**
2 hours ago

I can confirm that linear regression does not assume predictors are normally distributed. The Box-Cox transformations are nonlinear and can help with the linearity of the relationship between response and predictor in linear regression, but in a regression tree or RF, there should be no difference except for the seed for RF.

**Code**     Output (3)     Comments (21)     Log     Versions (1)     Forks (13)        Fork Script

**M↓** Styling with Markdown supported

Enter your comments.

Post Reply

Our Team    Careers    Terms    Privacy    Contact/Support

**Code**     Output (3)     Comments (21)     Log     Versions (1)     Forks (13)          Fork Script