G+1 ⟨ 0 ⟩

Create Blog　Sign In

# Kyu-DATA

categories

Sunday, January 31, 2016

## Machine Learning for Data Analysis Week3 Lasso Regression

```python
1    from pandas import Series, DataFrame
2    import matplotlib.pylab as plt
3    from sklearn.cross_validation import train_test_split
4    from sklearn.metrics import classification_report
5    import pandas as pd
6    import numpy as np
7    from sklearn import preprocessing
8
9    data = pd.read_csv('data.csv', low_memory=False)
10
11   #-------------------------------------------------
12   # Data Cleaning
13   #-------------------------------------------------
14   # explanatory variable
15   predictors = data[['gender_respondent',
16   'interest_attention','interest_whovote2008','presapp_track','presapp_job_x','pr
17   'presapp_foreign_x','presapp_health_x','presapp_war_x','finance_finfam','finan
18   'finance_finnext_x','health_insured','health_2010hcr_x','libcpre_choose',
19   'divgov_splitgov','campfin_limcorp','campfin_banads','ineq_incgap_x','effic_und
20   'effic_carestd','econ_ecpast_x',
21   'econ_ecnext_x','econ_unpast_x','ecblame_pres','ecblame_fmpr','ecblame_dem',
22   'tea_supp_x','gun_importance','immig_policy','fedspend_ss','fedspend_schools',
23   'fedspend_scitech','fedspend_crime','fedspend_welfare','fedspend_child',
24   'fedspend_poor','fedspend_enviro','dem_marital','dem_edugroup','dem_eduspgroup
25   'dem_veteran','dem_empstatus_1digitfinal','dem_raceeth','dem_parents','dem2_nur
26   'owngun_owngun', 'orientn_rgay', 'happ_lifesatisf']]
27
28   # target variable
29   targets = data['prevote_regpty']
30
31   # Convert categorical variable to numpy arrays and fill NaN values to zero.
32   # predictors[col] = number.fit_transform(predictors[col].replace(np.nan,'0', re
33   def convert(dta):
34       number = preprocessing.LabelEncoder()
35       for col in dta.columns:
36           dta[col] = number.fit_transform(dta[col].fillna(''))
37       return dta
38
39           # Catagorizing income group function
40   def incgroup_prepost(row):
41       if type(row) == float and np.isnan(row):
42           return float('NaN')
43       elif row == "$15,000-$17,499" or row == "$10,000-$12,499" or row == "$5,00
44           return 1
45       elif row == "$27,500-$29,999" or row == "$25,000-$27,499" or row == "$20,0
46           return 2
47       elif row == "$35,000-$39,999" or row == "$30,000-$34,999":
48           return 3
49       elif row == "$45,000-$49,999" or row == "$40,000-$44,999":
50           return 4
51       elif row == "$50,000-$54,999" or row == "$55,000-$59,999":
52           return 5
53       elif row == "$60,000-$64,999" or row == "$65,000-$69,999":
54           return 6
55       elif row == "$70,000-$74,999" or row == "$75,000-$79,999":
56           return 7
57       elif row == "$80,000-$89,999":
58           return 8
```

```
59          elif row == "$90,000-$99,999":
60              return 9
61          elif row == "$100,000-$109,999":
62              return 10
63          elif row == "$110,000-$124,999" or row == "$125,000-$149,999":
64              return 11
65          elif row == "$150,000-$174,999" or row == "$175,000-$249,999":
66              return 15
67          elif row == "$250,000 Or More":
68              return 25
69
70   # Explantory var Cleaning
71   predictors = convert(predictors)
72   predictors['incgroup_prepost'] = (data['incgroup_prepost'].apply(lambda row: i
73
74   # Response var Cleaning
75   number = preprocessing.LabelEncoder()
76   targets = number.fit_transform(targets.fillna(''))
77
78   # Spliting Data
79   pred_train, pred_test, tar_train, tar_test  = train_test_split(predictors, tar
80
81
82   #-------------------------------------------------------------------------
83   # Building Lasso Model
84   #-------------------------------------------------------------------------
85   # standardize predictors to have mean=0 and sd=1
86   import matplotlib.pylab as plt
87   from sklearn.linear_model import LassoLarsCV
88   from sklearn import preprocessing
89
90   # standardize clustering variables to have mean=0 and sd=1
91   predictors = predictors.copy()
92   def stdNscale(dta):
93       for col in dta.columns:
94           predictors[col] = preprocessing.scale(predictors[col].astype('float64')
95       return dta
96   predictors = stdNscale(predictors)
97   targets = preprocessing.scale(targets.astype('float64'))
98
99   # split data into train and test sets
100  pred_train, pred_test, tar_train, tar_test = train_test_split(predictors, targe
101                                                   test_size=.3, ra
102
103
104  # specify the lasso regression model
105  model = LassoLarsCV(cv=10, precompute=False).fit(pred_train,tar_train)
106
107  # print variable names and regression coefficients
108  dict(zip(predictors.columns, model.coef_))
109
110  # plot coefficient progression
111  m_log_alphas = -np.log10(model.alphas_)
112  ax = plt.gca()
113  plt.plot(m_log_alphas, model.coef_path_.T)
114  plt.axvline(-np.log10(model.alpha_), linestyle='--', color='k',
115              label='alpha CV')
116  plt.ylabel('Regression Coefficients')
117  plt.xlabel('-log(alpha)')
118  plt.title('Regression Coefficients Progression for Lasso Paths')
119
120  # plot mean square error for each fold
121  m_log_alphascv = -np.log10(model.cv_alphas_)
122  plt.figure()
123  plt.plot(m_log_alphascv, model.cv_mse_path_, ':')
124  plt.plot(m_log_alphascv, model.cv_mse_path_.mean(axis=-1), 'k',
125           label='Average across the folds', linewidth=2)
126  plt.axvline(-np.log10(model.alpha_), linestyle='--', color='k',
127              label='alpha CV')
128  plt.legend()
129  plt.xlabel('-log(alpha)')
130  plt.ylabel('Mean squared error')
131  plt.title('Mean squared error on each fold')
132
133
```

```
134    # MSE from training and test data
135    from sklearn.metrics import mean_squared_error
136    train_error = mean_squared_error(tar_train, model.predict(pred_train))
137    test_error = mean_squared_error(tar_test, model.predict(pred_test))
138    print ('training data MSE')
139    print(train_error)
140    print ('test data MSE')
141    print(test_error)
142
143    # R-square from training and test data
144    rsquared_train = model.score(pred_train, tar_train)
145    rsquared_test = model.score(pred_test, tar_test)
146    print ('training data R-square')
147    print(rsquared_train)
148    print ('test data R-square')
149    print(rsquared_test)
```
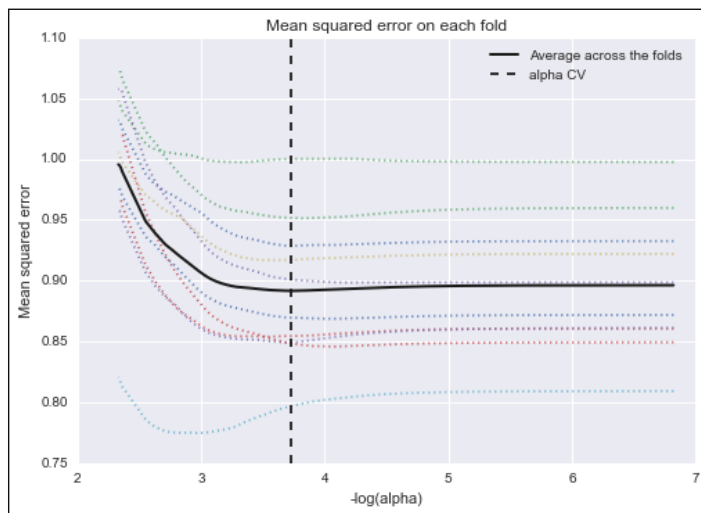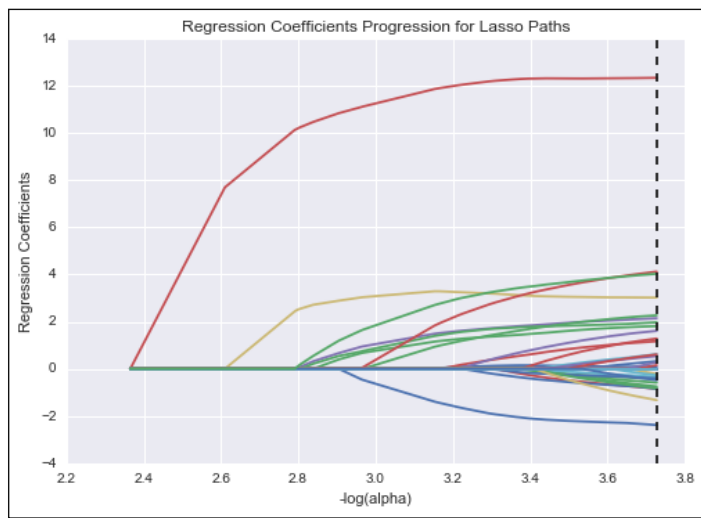
**project12.py** hosted with ♥ by GitHub                    **view raw**

```
Out[69]:
'campfin_banads': 0.0,
'campfin_limcorp': 0.0,
'dem2_numchild': 0.0,
'dem_edugroup': 0.0,
'dem_eduspgroup': 0.0,
'dem_empstatus_1digitfinal': 0.0,
'dem_marital': 0.0093908673093667499,
'dem_parents': 0.064314130130658162,
'dem_raceeth': 0.030557155678705504,
'dem_veteran': 0.0,
'divgov_splitgov': 0.0,
'ecblame_dem': 0.0024174774571679579,
'ecblame_fmpr': 0.035363241515626649,
'ecblame_pres': -0.013176501904743277,
'econ_ecnext_x': 0.0,
'econ_ecpast_x': 0.0,
'econ_unpast_x': -0.0042517012173674644,
'effic_carestd': 0.0,
'effic_undstd': -0.011690301880860978,
'fedspend_child': -0.0053244859938198202,
'fedspend_crime': 0.0042714539023118291,
'fedspend_enviro': -0.012919955395559209,
'fedspend_poor': -0.0063130602487481121,
'fedspend_schools': -0.0068057924372977568,
'fedspend_scitech': 0.019934013244821611,
'fedspend_ss': -0.037135780387225703,
'fedspend_welfare': 0.0,
'finance_finfam': 0.025162446494403065,
'finance_finnext_x': 0.0096358048042361599,
'finance_finpast_x': -0.0034004738990444745,
'gender_respondent': 0.0050017603010246436,
'gun_importance': 0.0,
'happ_lifesatisf': -0.0074631875707886104,
'health_2010hcr_x': 0.028169615424706553,
'health_insured': 0.0,
'immig_policy': 0.0,
'incgroup_prepost': 0.062892261894390789,
'ineq_incgap_x': 0.0,
'interest_attention': -0.0090041515906069888,
'interest_whovote2008': 0.19173136306395783,
'libcpre_choose': -0.013095974298389566,
'orientn_rgay': 0.0,
'owngun_owngun': -0.020832005352398528,
'presapp_econ_x': 0.0,
'presapp_foreign_x': 0.0013025839386105601,
'presapp_health_x': 0.0,
'presapp_job_x': 0.047018777798969492,
'presapp_track': 0.033074628032719842,
'presapp_war_x': 0.01812625523301812,
'tea_supp_x': 0.0083493998108720071
```

Regression Coefficients Progression for Lasso Paths



Mean squared error on each fold

We can see that 19 variables are removed out of 50 variables after i applied the lasso penalty. We can see that income group is positively related with respondent's political party and gun owned is negatively related with it.

training data MSE
0.877166028246
test data MSE
0.890551394001

training data R-square
0.118174140339
test data R-square
0.120284409471

We have similar MSE from training set and testing set which mean the prediction is pretty stable.

If you want to know the detail variables that is been used for this analysis, check the following link
https://d396qusza40orc.cloudfront.net/statistics%2Fproject%2Fanes1.html#incgroup_prepost

Posted by Kyu Cho at 6:51 PM

G+1  Recommend this on Google

Labels: categories

No comments:

Post a Comment

Enter your comment...

**Comment as:** Google Accou ▼

Publish    Preview

Subscribe to: Post Comments (Atom)

Simple template. Powered by Blogger.