

# Regression Problem: Forest Fires dataset

In [1]:

```

1 import os, sys
2 import re
3 import pandas as pd
4 from scipy import stats
5 import numpy as np
6
7 import urllib.request
8 import json
9
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 # all the scikit-learn libraries
13
14 from sklearn.preprocessing import scale, LabelEncoder
15 import sklearn.linear_model as skl_lm
16 from sklearn.metrics import mean_squared_error, r2_score
17 import statsmodels.api as sm
18 import statsmodels.formula.api as smf
19 from sklearn.preprocessing import label_binarize, PolynomialFeatures
20 from sklearn.model_selection import train_test_split, GridSearchCV
21 from sklearn.metrics import confusion_matrix, roc_curve, auc, classification_report
22 from sklearn.preprocessing import scale, StandardScaler
23 from sklearn import model_selection
24 from sklearn.linear_model import LinearRegression, Ridge, RidgeCV, Lasso, LassoCV
25 from sklearn.decomposition import PCA
26 from sklearn.cross_decomposition import PLSRegression
27 from sklearn.model_selection import KFold, cross_val_score
28 from patsy import dmatrix
29
30 %matplotlib inline
31 plt.style.use('ggplot')

```

/Users/User/anaconda/lib/python3.5/site-packages/statsmodels/compat/pandas.py:56: FutureWarning: The pandas.core.datetools module is deprecated and will be removed in a future version. Please use the pandas.tseries module instead.

```
from pandas.core import datetools
```

## LOAD THE DATASET :

***We import urllib.request to load the dataset dynamically from the repository, instead of having to save it/ downloading it prior to running the notebook***

In [2]:

```
1 url = 'https://archive.ics.uci.edu/ml/machine-learning-databases/forest-fires/forestfi'
```

In [3]:

```
1 raw_firedata = urllib.request.urlopen(url)
```

In [4]:

```

1 fire_dataset = raw_firedata.read().decode('utf-8')
2 with open('forestfires.csv', 'w+') as f:
3     f.write(fire_dataset)
4 f.close()

```

## Dataset Visualization

In [5]:

```

1 fire_df = pd.read_csv('forestfires.csv', na_values='?').dropna()
2 fire_df.head(7)

```

Out[5]:

	X	Y	month	day	FFMC	DMC	DC	ISI	temp	RH	wind	rain	area
0	7	5	mar	fri	86.2	26.2	94.3	5.1	8.2	51	6.7	0.0	0.0
1	7	4	oct	tue	90.6	35.4	669.1	6.7	18.0	33	0.9	0.0	0.0
2	7	4	oct	sat	90.6	43.7	686.9	6.7	14.6	33	1.3	0.0	0.0
3	8	6	mar	fri	91.7	33.3	77.5	9.0	8.3	97	4.0	0.2	0.0
4	8	6	mar	sun	89.3	51.3	102.2	9.6	11.4	99	1.8	0.0	0.0
5	8	6	aug	sun	92.3	85.3	488.0	14.7	22.2	29	5.4	0.0	0.0
6	8	6	aug	mon	92.3	88.9	495.6	8.5	24.1	27	3.1	0.0	0.0

In [6]:

```
1 fire_df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 517 entries, 0 to 516
Data columns (total 13 columns):
X          517 non-null int64
Y          517 non-null int64
month      517 non-null object
day        517 non-null object
FFMC       517 non-null float64
DMC        517 non-null float64
DC         517 non-null float64
ISI        517 non-null float64
temp       517 non-null float64
RH         517 non-null int64
wind       517 non-null float64
rain       517 non-null float64
area       517 non-null float64
dtypes: float64(8), int64(3), object(2)
memory usage: 56.5+ KB

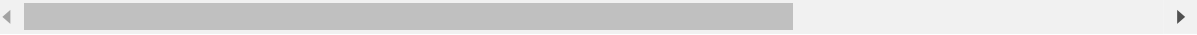
```

In [7]:

```
1 fire_df.describe()
```

Out[7]:

	X	Y	FFMC	DMC	DC	ISI	temp	
count	517.000000	517.000000	517.000000	517.000000	517.000000	517.000000	517.000000	51
mean	4.669246	4.299807	90.644681	110.872340	547.940039	9.021663	18.889168	4
std	2.313778	1.229900	5.520111	64.046482	248.066192	4.559477	5.806625	1
min	1.000000	2.000000	18.700000	1.100000	7.900000	0.000000	2.200000	1
25%	3.000000	4.000000	90.200000	68.600000	437.700000	6.500000	15.500000	3
50%	4.000000	4.000000	91.600000	108.300000	664.200000	8.400000	19.300000	4
75%	7.000000	5.000000	92.900000	142.400000	713.900000	10.800000	22.800000	5
max	9.000000	9.000000	96.200000	291.300000	860.600000	56.100000	33.300000	10

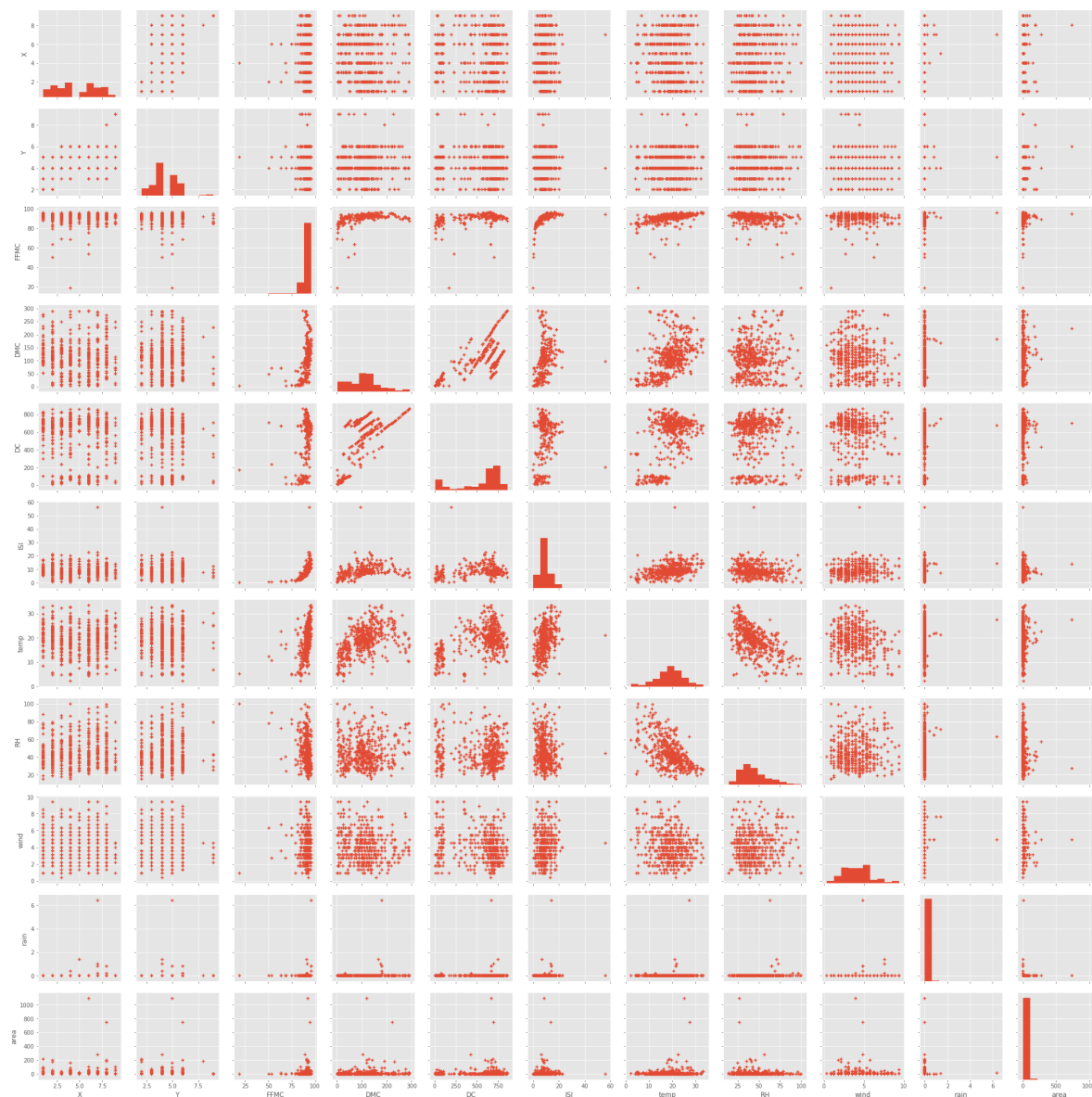


In [8]:

```
1 sns.pairplot(fire_df, markers='+')
```

Out[8]:

&lt;seaborn.axisgrid.PairGrid at 0x110937b00&gt;



## Dataset Cleaning

In [9]:

```
1 # converting categorical variables into numerical values
2 label = LabelEncoder()
3 fire_df.month = label.fit_transform(fire_df.month)
4 fire_df.day = label.fit_transform(fire_df.day)
```

In [10]:

```
1 fire_df.head()
```

Out[10]:

	X	Y	month	day	FFMC	DMC	DC	ISI	temp	RH	wind	rain	area
0	7	5	7	0	86.2	26.2	94.3	5.1	8.2	51	6.7	0.0	0.0
1	7	4	10	5	90.6	35.4	669.1	6.7	18.0	33	0.9	0.0	0.0
2	7	4	10	2	90.6	43.7	686.9	6.7	14.6	33	1.3	0.0	0.0
3	8	6	7	0	91.7	33.3	77.5	9.0	8.3	97	4.0	0.2	0.0
4	8	6	7	3	89.3	51.3	102.2	9.6	11.4	99	1.8	0.0	0.0

In [11]:

```
1 fire_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Int64Index: 517 entries, 0 to 516
```

```
Data columns (total 13 columns):
```

```
X          517 non-null int64
```

```
Y          517 non-null int64
```

```
month      517 non-null int64
```

```
day        517 non-null int64
```

```
FFMC       517 non-null float64
```

```
DMC        517 non-null float64
```

```
DC         517 non-null float64
```

```
ISI        517 non-null float64
```

```
temp       517 non-null float64
```

```
RH         517 non-null int64
```

```
wind       517 non-null float64
```

```
rain       517 non-null float64
```

```
area       517 non-null float64
```

```
dtypes: float64(8), int64(5)
```

```
memory usage: 76.5 KB
```

In [12]:

```
1 fireX = fire_df.drop(["area"], axis=1) # X = all columns excluding the response
2 print(fireX.shape)
3
4 fireY = fire_df.area # Y = response (area column)
5 print(fireY.shape)
6
7 # split the training and test data
8 X_train, X_test, y_train, y_test = train_test_split(fireX, fireY, test_size = 0.3, ran
```

```
(517, 12)
```

```
(517,)
```

## Model I: Multiple Linear Regression

In [13]:

```
1 xcolumns = fireX.iloc[0:0, :]
2 print(xcolumns)
```

Empty DataFrame

Columns: [X, Y, month, day, FFMC, DMC, DC, ISI, temp, RH, wind, rain]

Index: []

In [14]:

```
1 smf_regr = 'area ~ '
2
3 xcols = ['X', 'Y', 'month', 'day', 'FFMC', 'DMC', 'DC', 'ISI', 'temp', 'RH', 'wind', 'rain']
4
5 for i in xcols[:-1]:
6     # script to create multiple linear regression formula for statmodels
7     smf_regr += '{ } + '.format(i)
8 smf_regr += xcols[-1]
9
10 print(smf_regr)
```

area ~ X + Y + month + day + FFMC + DMC + DC + ISI + temp + RH + wind + rain

In [15]:

```
1 # statsmodels -- ordinary least squares without any interaction term
2 # try to fit the model to all the data
3 est = smf.ols(smf_regr, fire_df).fit()
4 est.summary().tables[1]
```

Out[15]:

	coef	std err	t	P> t	[0.025	0.975]
<b>Intercept</b>	-11.5006	63.324	-0.182	0.856	-135.912	112.911
<b>X</b>	1.8812	1.450	1.298	0.195	-0.967	4.729
<b>Y</b>	0.5268	2.740	0.192	0.848	-4.856	5.910
<b>month</b>	0.9733	0.776	1.254	0.210	-0.551	2.498
<b>day</b>	0.4995	1.497	0.334	0.739	-2.442	3.441
<b>FFMC</b>	-0.1074	0.664	-0.162	0.872	-1.411	1.197
<b>DMC</b>	0.1098	0.072	1.524	0.128	-0.032	0.251
<b>DC</b>	-0.0146	0.019	-0.779	0.437	-0.052	0.022
<b>ISI</b>	-0.6108	0.779	-0.784	0.433	-2.141	0.920
<b>temp</b>	0.9801	0.802	1.222	0.222	-0.596	2.556
<b>RH</b>	-0.1849	0.240	-0.770	0.442	-0.657	0.287
<b>wind</b>	1.7823	1.680	1.061	0.289	-1.519	5.084
<b>rain</b>	-3.2517	9.699	-0.335	0.738	-22.306	15.803

In this model where we consider all the input features for regression, some of the features have positive coefficients while others have negative. The p-values are rather high for all the input features.

In [16]:

```
1 est.summary().tables[0]
```

Out[16]:

OLS Regression Results

<b>Dep. Variable:</b>	area	<b>R-squared:</b>	0.025
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.002
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1.065
<b>Date:</b>	Sat, 12 May 2018	<b>Prob (F-statistic):</b>	0.388
<b>Time:</b>	19:18:46	<b>Log-Likelihood:</b>	-2874.0
<b>No. Observations:</b>	517	<b>AIC:</b>	5774.
<b>Df Residuals:</b>	504	<b>BIC:</b>	5829.
<b>Df Model:</b>	12		
<b>Covariance Type:</b>	nonrobust		

A good model will have an F-statistic much greater than 1. The F-statistic for this model(when all the input features are considered) is low, i.e, close to 1. An R-squared statistic near 0 indicates that the regression did not explain much of the variability in the response; this might occur because the linear model is wrong, or the inherent error(variance) is high, or both.

In [17]:

```
1 print('Residual Sum of Squares =', est.mse_resid)
```

Residual Sum of Squares = 4045.97315219

## Feature Selection: Backward Step-wise

**Metric considered: F-statistic**

In [18]:

```

1 # Backward selection for Advertising dataset -- a modification of the Forward Selection
2
3 def backward_selected(data, response):
4     """Linear model designed by backward selection.
5     Parameters:
6     -----
7     data : pandas dataframe (fire_df)
8     response: string, name of response column in fire_df2 ("area")
9     Returns:
10    -----
11    model: an "optimal" fitted statsmodels linear model
12    with an intercept
13    selected by backward selection
14    evaluated by F-statistic
15    """
16
17    # in this algorithm, we start with all the values and begin to eliminate features
18    remaining = list(data.columns) # initializing the list 'remaining' with a set of the
19    remaining.remove(response) # the output/response field is removed from the list of
20
21    selected = [] # list initialized to progressively add best features
22
23    # create the formula to calculate the OLS of all data first (to evaluate Rsquare-value)
24    formula = "{} ~ {} + 1".format(response, ' + '.join(remaining))
25
26    total_score = smf.ols(formula, data).fit().fvalue
27    print('F-statistic with all the features =', total_score, '\n')
28
29    for candidate in remaining:
30        backup = remaining[:] # keep the original list of predictors intact, and use a
31        backup.remove(candidate) # test for Rsqr-value by removing each candidate
32
33        formula = "{} ~ {} + 1".format(response, ' + '.join(backup))
34        #print()
35        #print(formula)
36
37        score = smf.ols(formula, data).fit().fvalue # the rsquared_adj() method is specified
38        print(score)
39        if score > total_score:
40            total_score = score # total score get assigned the new value of the highest
41            remaining.remove(candidate) # only features responsible for this best fit
42
43    formula = "{} ~ {} + 1".format(response, ' + '.join(remaining))
44    model = smf.ols(formula, data).fit() # model.formula
45    return model
46
47 # call the function to perform feature selection on the forest fire dataset
48 b_model = backward_selected(fire_df, 'area') # data = dataset, response = output
49
50 print('\nFinal F-statistic for the model with selected features =', '{:3.10f}'.format(
51 print('\nOriginal feature set: \n', smf_regr)
52 print('\nOptimized feature set: \n', b_model.model.formula)

```

F-statistic with all the features = 1.06472438483

1.00705467957  
1.16036964164  
1.266888283



1.13941624267  
 1.33941934275  
 1.36792604734  
 1.46467104637

Final F-statistic for the model with selected features = 1.4646710464

Original feature set:

area ~ X + Y + month + day + FFMC + DMC + DC + ISI + temp + RH + wind + rain

Optimized feature set:

area ~ X + month + FFMC + DMC + ISI + RH + rain + 1

In [19]:

```
1 est_rsqr = smf.ols(b_model.model.formula, fire_df).fit()
2 est_rsqr.summary().tables[0]
```

Out[19]:

OLS Regression Results

<b>Dep. Variable:</b>	area	<b>R-squared:</b>	0.020
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.006
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1.465
<b>Date:</b>	Sat, 12 May 2018	<b>Prob (F-statistic):</b>	0.177
<b>Time:</b>	19:18:46	<b>Log-Likelihood:</b>	-2875.3
<b>No. Observations:</b>	517	<b>AIC:</b>	5767.
<b>Df Residuals:</b>	509	<b>BIC:</b>	5801.
<b>Df Model:</b>	7		
<b>Covariance Type:</b>	nonrobust		

F-statistic of the final model is higher than the original model.

In [20]:

```
1 print('Residual Sum of Squares =', est_rsqr.mse_resid)
```

Residual Sum of Squares = 4026.68033268

## Model II : kNN

Original feature set: area ~ Y + day + DC + + temp + + wind +

Optimized feature set: area ~ X + month + FFMC + DMC + ISI + RH + rain + 1

In [21]:

```
1 fireX = fire_df.drop(["area"], axis=1) # X = all columns excluding the response
2 print(fireX.shape)
3
4 fireY = fire_df.area # Y = response (area column)
5 print(fireY.shape)
6
7 # split the training and test data
8 X_train, X_test, y_train, y_test = train_test_split(fireX, fireY, test_size = 0.3, rand
```

(517, 12)

(517,)

In [22]:

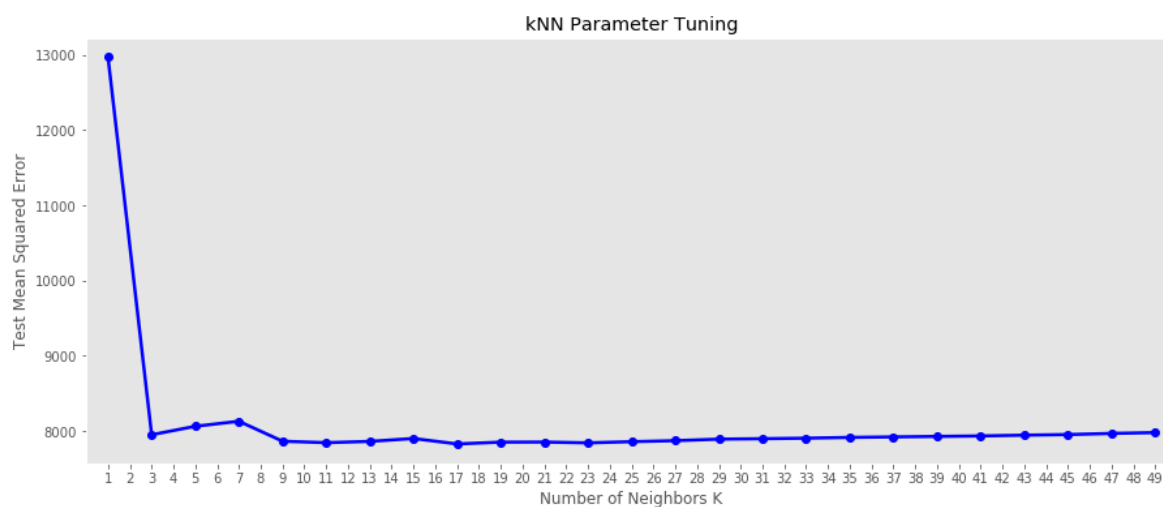
```
1 from sklearn.neighbors import KNeighborsRegressor
2
3 # Creating odd List of K for KNN
4 myList = list(range(1,50))
5
6 # Subsetting just the odd ones
7 neighbors = list(filter(lambda x: x % 2 != 0, myList))
8
9 # Empty List that will hold cv scores
10 MSE = []
11
12 #Find the best k
13 for k in neighbors:
14     knn = KNeighborsRegressor(n_neighbors=k)
15     knn.fit(X_train, y_train)
16     y_test_pred = knn.predict(X_test)
17     scores = mean_squared_error(y_test, y_test_pred)
18     MSE.append(scores)
```

In [23]:

```

1 # Plot misclassification error vs k
2 plt.figure(figsize=(15,6))
3 plt.plot(neighbors, MSE, 'b-o', linewidth =2.5)
4 plt.xlim(0, 50)
5 plt.xticks(range(1,50))
6 plt.xlabel('Number of Neighbors K')
7 plt.ylabel('Test Mean Squared Error')
8 plt.title('kNN Parameter Tuning')
9 plt.grid()
10 plt.show()
11
12 # # Determining best k
13 optimal_k = neighbors[MSE.index(min(MSE))]
14 print ("The optimal number of neighbors is %d" % optimal_k)
15 print("Best Test MSE = ", min(MSE))

```



The optimal number of neighbors is 17  
 Best Test MSE = 7826.06492685

## Model III: Lasso Regression

In [24]:

```

1 fireX = fire_df.drop(["area"], axis=1) # X = all columns excluding the response
2 print(fireX.shape)
3
4 fireY = fire_df.area # Y = response (area column)
5 print(fireY.shape)
6
7 # Feature Scaling
8 scaler = StandardScaler()
9 fireX_scaled = scaler.fit_transform(fireX)
10
11 # split the training and test data
12 X_train, X_test, y_train, y_test = train_test_split(fireX_scaled, fireY, test_size = 0

```

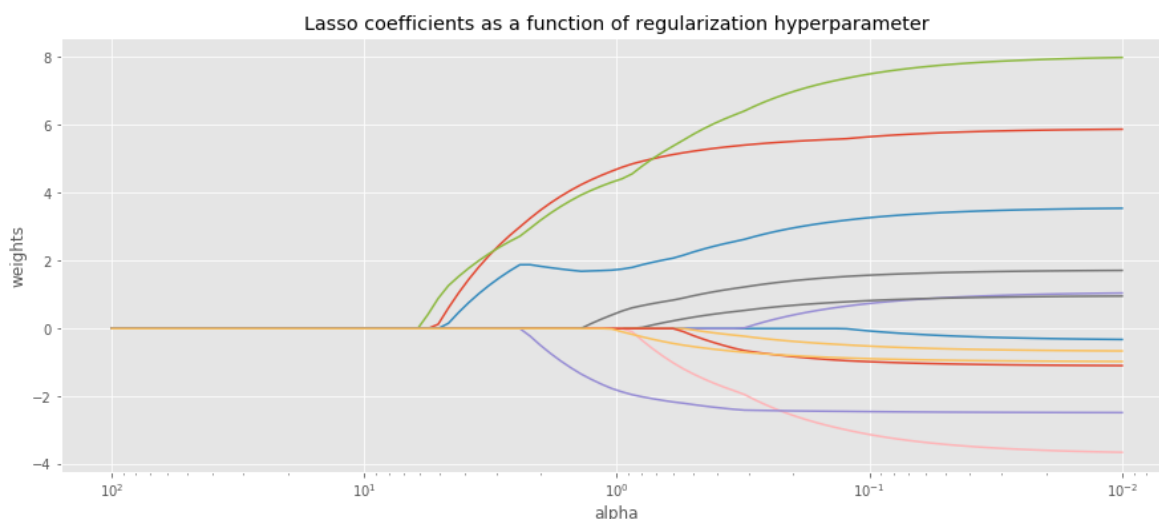
(517, 12)  
 (517,)

In [25]:

```

1 alphas = 10**np.linspace(2,-2,100)*0.5
2
3 lasso = Lasso()
4 coefs = []
5
6 for a in alphas*2:
7     lasso.set_params(alpha=a)
8     lasso.fit(X_train, y_train)
9     coefs.append(lasso.coef_)
10
11 ax = plt.gca()
12 # labels = ['X' , 'Y' , 'month', 'day ', 'FFMC' , 'DMC' , 'DC', 'ISI' , 'temp', 'RH', 'wind', 'rc
13 ax.plot(alphas*2, coefs)
14 ax.set_xscale('log')
15 ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
16 plt.axis('tight')
17 plt.xlabel('alpha')
18 plt.ylabel('weights')
19 # plt.legend()
20 plt.title('Lasso coefficients as a function of regularization hyperparameter')
21 plt.gcf().set_size_inches(15,6)

```



In [26]:

```

1 lassoCV = LassoCV(alphas=None, cv=10)
2 lassoCV.fit(X_train, y_train)

```

Out[26]:

```

LassoCV(alphas=None, copy_X=True, cv=10, eps=0.001, fit_intercept=True,
max_iter=1000, n_alphas=100, n_jobs=1, normalize=False, positive=False,
precompute='auto', random_state=None, selection='cyclic', tol=0.0001,
verbose=False)

```

In [27]:

```

1 print('Best alpha (Lasso) =',lassoCV.alpha_)
2 print('Test MSE =',mean_squared_error(y_test, lassoCV.predict(X_test)))

```

```

Best alpha (Lasso) = 1.71868752562
Test MSE = 7936.52990235

```

In [28]:

```
1 # Best Lasso Coefficients
2 pd.Series(lassocv.coef_, index=fireX.columns)
```

Out[28]:

```
X          3.837107
Y           0.000000
month      -0.000000
day         0.000000
FFMC        0.000000
DMC         3.536034
DC          -0.000000
ISI          0.000000
temp        1.767039
RH          -0.885794
wind         0.000000
rain        -0.000000
dtype: float64
```

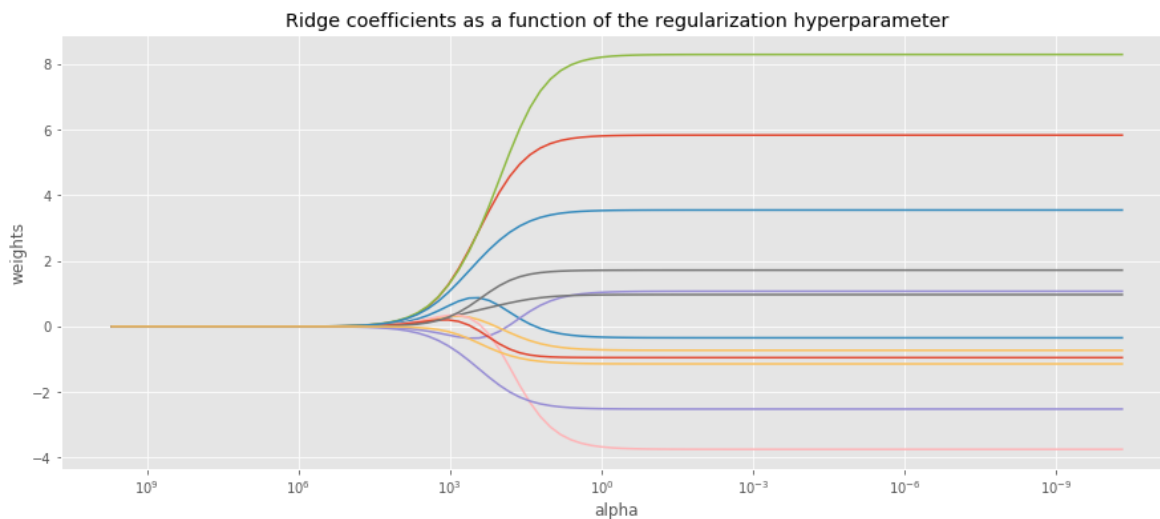
## Model IV: Ridge Regression

In [29]:

```

1 alphas = 10**np.linspace(10,-10,100)*0.5
2
3 ridge = Ridge()
4 coefs = []
5
6 for a in alphas:
7     ridge.set_params(alpha=a)
8     ridge.fit(scale(X_train), y_train)
9     coefs.append(ridge.coef_)
10
11 ax = plt.gca()
12 ax.plot(alphas, coefs)
13 ax.set_xscale('log')
14 ax.set_xlim(ax.get_xlim()[::-1]) # reverse axis
15 plt.axis('tight')
16 plt.xlabel('alpha')
17 plt.ylabel('weights')
18 plt.title('Ridge coefficients as a function of the regularization hyperparameter');
19 plt.gcf().set_size_inches(15,6)

```



In [30]:

```

1 ridgecv = RidgeCV(alphas=alphas, cv=10)
2 ridgecv.fit(scale(X_train), y_train)

```

Out[30]:

```

RidgeCV(alphas=array([ 5.00000e+09,  3.14015e+09, ...,  7.96141e-11,  5.
00000e-11]),
        cv=10, fit_intercept=True, gcv_mode=None, normalize=False,
        scoring=None, store_cv_values=False)

```

In [31]:

```
1 ridgecv.alpha_
```

Out[31]:

```
1715.2346431574597
```

In [32]:

```
1 print('Test MSE =',mean_squared_error(y_test, ridgecv.predict(X_test)))
```

Test MSE = 7945.21888668

In [33]:

```
1 #Best Ridge Coefficients  
2 pd.Series(ridge.coef_, index=fireX.columns)
```

Out[33]:

```
X          5.835932  
Y         -0.344120  
month      1.075099  
day         0.971851  
FFMC       -0.725321  
DMC         8.293224  
DC          -3.744169  
ISI         -0.947774  
temp        3.550765  
RH          -2.517020  
wind        1.718755  
rain        -1.138975  
dtype: float64
```