

[Ordinal regression](#) is a classification method for categories on an ordinal scale – e.g. [1, 2, 3, 4, 5] or [G, PG, PG-13, R]. This notebook implements ordinal regression using the method of [Frank and Hal 2001](#), which transforms a k-multiclass classifier into k-1 binary classifiers (each of which predicts whether a data point is above a threshold in the ordinal scale – e.g., whether a movie is "higher" than PG). This method can be used with any binary classification method that outputs probabilities; here L2-regularized binary logistic regression is used.

This notebook trains a model (on `train.txt`), optimizes L2 regularization strength on `dev.txt`, and evaluates performance on `test.txt`. Reports test accuracy with 95% confidence intervals.

▼ 새 섹션

```
from scipy import sparse
from sklearn import linear_model
from collections import Counter
import numpy as np
import operator
import nltk
import math
from scipy.stats import norm
```

```
!python -m nltk.downloader punkt
```

```
/usr/lib/python3.10/runpy.py:126: RuntimeWarning: 'nltk.downloader' found in sys.modules after import of package 'nltk', but prior to execution of 'n
warn(RuntimeWarning(msg))
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
```

```
def load_ordinal_data(filename, ordering):
    X = []
    Y = []
    orig_Y=[]
    for ordinal in ordering:
        Y.append([])

    with open(filename, encoding="utf-8") as file:
        i = 0
        for line in file:
            if i == 0:
                i += 1
                continue

            cols = line.split("Wt")
            idd = cols[0]
            label = cols[2].lstrip().rstrip()
            text = cols[3]

            X.append(text)

            index=ordering.index(label)
            for i in range(len(ordering)):
                if index > i:
                    Y[i].append(1)
                else:
                    Y[i].append(0)
            orig_Y.append(label)

    return X, Y, orig_Y
```

```
class OrdinalClassifier:
```

```
    def __init__(self, ordinal_values, feature_method, trainX, trainY, devX, devY, testX, testY, orig_trainY, orig_devY, orig_testY):
        self.ordinal_values=ordinal_values
        self.feature_vocab = {}
        self.feature_method = feature_method
        self.min_feature_count=2
        self.log_regs = [None]*(len(self.ordinal_values)-1)

        self.trainY=trainY
```

```

self.devY=devY
self.testY=testY

self.orig_trainY=orig_trainY
self.orig_devY=orig_devY
self.orig_testY=orig_testY

self.trainX = self.process(trainX, training=True)
self.devX = self.process(devX, training=False)
self.testX = self.process(testX, training=False)

# Featurize entire dataset
def featurize(self, data):
    featurized_data = []
    for text in data:
        feats = self.feature_method(text)
        featurized_data.append(feats)
    return featurized_data

# Read dataset and returned featurized representation as sparse matrix + label array
def process(self, X_data, training = False):

    data = self.featurize(X_data)

    if training:
        fid = 0
        feature_doc_count = Counter()
        for feats in data:
            for feat in feats:
                feature_doc_count[feat]+= 1

        for feat in feature_doc_count:
            if feature_doc_count[feat] >= self.min_feature_count:
                self.feature_vocab[feat] = fid
                fid += 1

    F = len(self.feature_vocab)
    D = len(data)
    X = sparse.dok_matrix((D, F))
    for idx, feats in enumerate(data):
        for feat in feats:
            if feat in self.feature_vocab:
                X[idx, self.feature_vocab[feat]] = feats[feat]

    return X

def train(self):
    (D,F) = self.trainX.shape

    for idx, ordinal_value in enumerate(self.ordinal_values[:-1]):
        best_dev_accuracy=0
        best_model=None
        for C in [0.1, 1, 10, 100]:

            log_reg = linear_model.LogisticRegression(C = C, max_iter=1000)
            log_reg.fit(self.trainX, self.trainY[idx])
            development_accuracy = log_reg.score(self.devX, self.devY[idx])
            if development_accuracy > best_dev_accuracy:
                best_dev_accuracy=development_accuracy
                best_model=log_reg

        self.log_regs[idx]=best_model

def test(self):
    cor=tot=0
    counts=Counter()
    preds=[None]*(len(self.ordinal_values)-1)
    for idx, ordinal_value in enumerate(self.ordinal_values[:-1]):
        preds[idx]=self.log_regs[idx].predict_proba(self.testX)[:,-1]

    preds=np.array(preds)

    for data_point in range(len(preds[0])):

```

```

        ordinal_preds=np.zeros(len(self.ordinal_values))
        for ordinal in range(len(self.ordinal_values)-1):
            if ordinal == 0:
                ordinal_preds[ordinal]=1-preds[ordinal][data_point]
            else:
                ordinal_preds[ordinal]=preds[ordinal-1][data_point]-preds[ordinal][data_point]

        ordinal_preds[len(self.ordinal_values)-1]=preds[len(preds)-1][data_point]

        prediction=np.argmax(ordinal_preds)
        counts[prediction]+=1
        if prediction == self.ordinal_values.index(self.orig_testY[data_point]):
            cor+=1
        tot+=1

    return cor/tot

import nltk
nltk.download('opinion_lexicon')

from nltk.corpus import opinion_lexicon

def binary_bow_featurize(text):
    feats = {}
    words = nltk.word_tokenize(text)

    #To assess capitalization of text (whether title contains at least one word that's in ALLCAP)
    ALLCAP = 0
    exclam = 0
    qmark = 0
    sarc_1 = 0
    sarc = set(['...', '???', ':', ':(', 'lmao'])

    # Get the positive and negative words from the opinion lexicon
    positive_words = set(opinion_lexicon.positive())
    negative_words = set(opinion_lexicon.negative())

    # Add some custom negative words
    negative_words.update(set(['not', 'no', 'never', 'hate', 'sad', 'Not worth', 'not good']))
    negative = 0
    # Add some custom positive words
    positive_words.update(set(['awesome', 'fantastic', 'amazing', 'love', 'happy', 'fun', 'the best', 'finally made it']))
    positive = 0

    for word in words:
        if len(word) != 1 and word == word.upper():
            ALLCAP = 1

        word=word.lower()
        if word in negative_words:
            negative += 1
        if word in positive_words:
            positive += 1

        if word in sarc:
            sarc_1 = 1
        if word == '?':
            qmark = 1
        if word == '!':
            exclam = 1

    length = len(words)

    #Clarity
    if length > 10 or qmark == 1:
        feats['Clarity'] = 1

    #Sarcasm
    if exclam == 1 or sarc_1 == 1 or ALLCAP == 1:
        feats['Sarcasm'] = 1

    #Tone
    if positive > negative:
        feats['positive'] = 1

```

```

elif negative > positive:
    feats['negative'] = 1
else:
    feats['neutral'] = 1

#Q vs. Statement
if qmark == 1:
    feats['QS'] = 1

#Capitalization
if ALLCAP == 1:
    feats['allcap'] = 1

#Sentence Length
if length <= 10:
    feats['length:0-10'] = 1
elif length <= 20:
    feats['length:11-20'] = 1
else:
    feats['length:20+'] = 1

return feats

[nltk_data] Downloading package opinion_lexicon to /root/nltk_data...
[nltk_data] Package opinion_lexicon is already up-to-date!

def confidence_intervals(accuracy, n, significance_level):
    critical_value=(1-significance_level)/2
    z_alpha=-1*norm.ppf(critical_value)
    se=math.sqrt((accuracy*(1-accuracy))/n)
    return accuracy-(se*z_alpha), accuracy+(se*z_alpha)

def run(trainingFile, devFile, testFile, ordinal_values):

    trainX, trainY, orig_trainY=load_ordinal_data(trainingFile, ordinal_values)
    devX, devY, orig_devY=load_ordinal_data(devFile, ordinal_values)
    testX, testY, orig_testY=load_ordinal_data(testFile, ordinal_values)

    simple_classifier = OrdinalClassifier(ordinal_values, binary_bow_featurize, trainX, trainY, devX, devY, testX, testY, orig_trainY, orig_devY, orig_test)
    simple_classifier.train()
    accuracy=simple_classifier.test()

    lower, upper=confidence_intervals(accuracy, len(testY[0]), .95)
    print("Test accuracy for best dev model: %.3f, 95% CIs: [%.3f %.3f]\n" % (accuracy, lower, upper))

trainingFile = "splits/train.txt"
devFile = "splits/dev.txt"
testFile = "splits/test.txt"

# ordinal values must be in order *as strings* from smallest to largest, e.g.:
# ordinal_values=["G", "PG", "PG-13", "R"]

ordinal_values=['Not popular', 'Average', 'Popular']

run(trainingFile, devFile, testFile, ordinal_values)

Test accuracy for best dev model: 0.450, 95% CIs: [0.352 0.548]

```

