

Logistic Regression

NAME

Logistic Regression - A classification algorithm to predict the binary output with logistic loss.

SYNOPSIS

```
import com.nec.frovedis.mllib.classification.LogisticRegressionWithSGD
LogisticRegressionModel
LogisticRegressionWithSGD.train(RDD[LabeledPoint] data,
    Int numIter = 1000,
    Double stepSize = 0.01,
    Double miniBatchFraction = 1.0,
    Double regParam = 0.01)

import com.nec.frovedis.mllib.classification.LogisticRegressionWithLBFGS
LogisticRegressionModel
LogisticRegressionWithLBFGS.train(RDD[LabeledPoint] data,
    Int numIter = 1000,
    Double stepSize = 0.01,
    Int histSize = 10,
    Double regParam = 0.01)
```

DESCRIPTION

Classification aims to divide the items into categories. The most common classification type is binary classification, where there are two categories, usually named positive and negative. Frovedis supports binary classification algorithm only.

Logistic regression is widely used to predict a binary response. It is a linear method with the loss function given by the **logistic loss**:

$$L(\mathbf{w}; \mathbf{x}, y) := \log(1 + \exp(-y\mathbf{w}^T\mathbf{x}))$$

Where the vectors \mathbf{x} are the training data examples and y are their corresponding labels (Frovedis considers negative response as -1 and positive response as 1, but when calling from Spark interface, user should pass 0 for negative response and 1 for positive response according to the Spark requirement) which we want to predict. \mathbf{w} is the linear model (also called as weight) which uses a single weighted sum of features to make a prediction. Frovedis Logistic Regression supports ZERO, L1 and L2 regularization to address the overfit problem. But when calling from Spark interface, it supports the default L2 regularization only.

The gradient of the logistic loss is: $-y(1 - 1 / (1 + \exp(-ywTx))) \cdot x$

The gradient of the L1 regularizer is: $\text{sign}(w)$

And The gradient of the L2 regularizer is: w

For binary classification problems, the algorithm outputs a binary logistic regression model. Given a new data point, denoted by x , the model makes predictions by applying the logistic function:

$$f(z) := 1 / (1 + \exp(-z))$$

Where $z = wTx$. By default (threshold=0.5), if $f(wTx) > 0.5$, the response is positive (1), else the response is negative (0).

Frovedis provides implementation of logistic regression with two different optimizers: (1) stochastic gradient descent with minibatch and (2) LBFGS optimizer.

The simplest method to solve optimization problems of the form $\min f(w)$ is gradient descent. Such first-order optimization methods well-suited for large-scale and distributed computation. Whereas, L-BFGS is an optimization algorithm in the family of quasi-Newton methods to solve the optimization problems of the similar form.

Like the original BFGS, L-BFGS (Limited Memory BFGS) uses an estimation to the inverse Hessian matrix to steer its search through feature space, but where BFGS stores a dense $n \times n$ approximation to the inverse Hessian (n being the number of features in the problem), L-BFGS stores only a few vectors that represent the approximation implicitly. L-BFGS often achieves rapid convergence compared with other first-order optimization.

This module provides a client-server implementation, where the client application is a normal Apache Spark program. Spark has its own mllib providing the Logistic Regression support. But that algorithm is slower when comparing with the equivalent Frovedis algorithm (see frovedis manual for ml/logistic_regression) with big dataset. Thus in this implementation, a spark client can interact with a frovedis server sending the required spark data for training at frovedis side. Spark RDD data is converted into frovedis compatible data internally and the spark ML call is linked with the respective frovedis ML call to get the job done at frovedis server.

Spark side call for Logistic Regression quickly returns, right after submitting the training request to the frovedis server with a dummy LogisticRegressionModel object containing the model information like threshold value etc. with a unique model ID for the submitted training request.

When operations like prediction will be required on the trained model, spark client sends the same request to frovedis server on the same model (containing the unique ID) and the request is served at frovedis server and output is sent back to the spark client.

Detailed Description

LogisticRegressionWithSGD.train()

Parameters

data: A RDD[LabeledPoint] containing spark-side distributed sparse training data

numIter: An integer parameter containing the maximum number of iteration count (Default: 1000)

stepSize: A double parameter containing the learning rate (Default: 0.01)

minibatchFraction: A double parameter containing the minibatch fraction (Default: 1.0)

regParam: A double parameter containing the regularization parameter (Default: 0.01)

Purpose

It trains a logistic regression model with stochastic gradient descent with minibatch optimizer and with default L2 regularizer. It starts with an initial guess of zeros for the model vector and keeps updating the model to minimize the cost function until convergence is achieved (default convergence tolerance value is

0.001) or maximum iteration count is reached. After the training, it returns the trained logistic regression model.

For example,

```
val data = MLUtils.loadLibSVMFile(sc, "./sample")
val splits = data.randomSplit(Array(0.8, 0.2), seed = 11L)
val training = splits(0)
val test = splits(1)
val tvec = test.map(_._features)

// training a logistic regression model with default parameters using SGD
val model = LogisticRegressionWithSGD.train(training)
// cross-validation of the trained model on 20% test data
model.predict(tvec).collect.foreach(println)
```

Note that, inside the `train()` function spark side sparse data is converted into frovedis side sparse data and after the training, frovedis side sparse data is released from the server memory. But if the user needs to store the server side constructed sparse data for some other operations, he may also like to pass the `FrovedisSparseData` object as the value of the “data” parameter. In that case, the user needs to explicitly release the server side sparse data when it will no longer be needed.

For example,

```
val fdata = new FrovedisSparseData(data) // manual creation of frovedis sparse data
val model2 = LogisticRegressionWithSGD.train(fdata) // passing frovedis sparse data
fdata.release() // explicit release of the server side data
```

Return Value

This is a non-blocking call. The control will return quickly, right after submitting the training request at frovedis server side with a `LogisticRegressionModel` object containing a unique model ID for the training request along with some other general information like threshold (default 0.5) etc. But it does not contain any weight values. It simply works like a spark side pointer of the actual model at frovedis server side. It may be possible that the training is not completed at the frovedis server side even though the client spark side `train()` returns with a pseudo model.

LogisticRegressionWithLBFGS.train()

Parameters

data: A `RDD[LabeledPoint]` containing spark-side distributed sparse training data
numIter: An integer parameter containing the maximum number of iteration count (Default: 1000)
stepSize: A double parameter containing the learning rate (Default: 0.01)
histSize: An integer parameter containing the gradient history size (Default: 1.0)
regParam: A double parameter containing the regularization parameter (Default: 0.01)

Purpose

It trains a logistic regression model with LBFGS optimizer and with default L2 regularizer. It starts with an initial guess of zeros for the model vector and keeps updating the model to minimize the cost function until convergence is achieved (default convergence tolerance value is 0.001) or maximum iteration count is reached. After the training, it returns the trained logistic regression model.

For example,

```

val data = MLUtils.loadLibSVMFile(sc, "./sample")
val splits = data.randomSplit(Array(0.8, 0.2), seed = 11L)
val training = splits(0)
val test = splits(1)
val tvec = test.map(_.features)

// training a logistic regression model with default parameters using LBFGS
val model = LogisticRegressionWithLBFGS.train(training)
// cross-validation of the trained model on 20% test data
model.predict(tvec).collect.foreach(println)

```

Note that, inside the `train()` function spark side sparse data is converted into frovedis side sparse data and after the training, frovedis side sparse data is released from the server memory. But if the user needs to store the server side constructed sparse data for some other operations, he may also like to pass the `FrovedisSparseData` object as the value of the “data” parameter. In that case, the user needs to explicitly release the server side sparse data when it will no longer be needed.

For example,

```

val fdata = new FrovedisSparseData(data) // manual creation of frovedis sparse data
val model2 = LogisticRegressionWithLBFGS.train(fdata) // passing frovedis sparse data
fdata.release() // explicit release of the server side data

```

Return Value

This is a non-blocking call. The control will return quickly, right after submitting the training request at frovedis server side with a `LogisticRegressionModel` object containing a unique model ID for the training request along with some other general information like threshold (default 0.5) etc. But it does not contain any weight values. It simply works like a spark side pointer of the actual model at frovedis server side. It may be possible that the training is not completed at the frovedis server side even though the client spark side `train()` returns with a pseudo model.

SEE ALSO

`logistic_regression_model`, `linear_svm`, `frovedis_sparse_data`