
xlearn_doc Documentation

Release 0.2.0

Chao Ma

Jan 17, 2018

Contents

1	Link to the Other Helpful Resources	3
2	Quick Install	5
3	Python Demo	7
3.1	Get Started with xLearn !	7
3.2	Installation Guide	9
3.3	xLearn Command Line Guide	11
3.4	xLearn Python Package Guide	17
3.5	xLearn R Package Guide	25
3.6	xLearn API List	25
3.7	Large-Scale Machine Learning	28
3.8	xLearn Demo	31
3.9	xLearn Tutorials	33

xLearn is a high-performance, easy-to-use, and scalable machine learning package, which can be used to solve large-scale machine learning problems, especially for the problems on large-scale sparse data, which is very common in scenes like CTR prediction and recommender system. If you are the user of liblinear, libfm, or libffm, now xLearn is your another better choice. This is because xLearn handles all of these models in a uniform platform and provides better performance and scalability compared to its competitors.

This is a quick start tutorial showing snippets for you to quickly try out xLearn on a small demo dataset (Criteo CTR prediction) for a binary classification task.

Link to the Other Helpful Resources

- See [Installation Guide](#) on how to install xLearn.
- See [Command Line Guide](#) on how to use xLearn command line.
- See [Python API Guide](#) on how to use xLearn Python API.
- See [R API Guide](#) on how to use xLearn R API.
- See [Demo Page Learning to use xLearn by Examples](#).
- See [Tutorial](#) on tutorials on specific tasks.

Here is a simple Python demo on how to use xLearn:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.setValidate("./small_test.txt")
param = {'task':'binary', 'lr':0.2,
         'lambda':0.002, 'metric':'auc'}

ffm_model.fit(param, "./model.out")

# Prediction task
ffm_model.setTest("./small_test.txt")
# Convert output to 0~1
ffm_model.setSigmoid()
ffm_model.predict("./model.out", "./output.txt")
```

This example shows how to use *field-aware factorizations machine (ffm)* to solve a simple binary classification task. You can check out the demo data (`small_train.txt` and `small_test.txt`) from the path `demo/classification/criteo_ctr`.

3.1 Get Started with xLearn !

xLearn is a high-performance, easy-to-use, and scalable machine learning package, which can be used to solve large-scale machine learning problems, especially for the problems on large-scale sparse data, which is very common in scenes like CTR prediction and recommender system. If you are the user of liblinear, libfm, or libffm, now xLearn is your another better choice. This is because xLearn handles all of these models in a uniform platform and provides better performance and scalability compared to its competitors.


```

ffm_model.setValidate("./small_test.txt")
param = {'task':'binary', 'lr':0.2,
        'lambda':0.002, 'metric':'auc'}

ffm_model.fit(param, "./model.out")

# Prediction task
ffm_model.setTest("./small_test.txt")
# Convert output to 0~1
ffm_model.setSigmoid()
ffm_model.predict("./model.out", "./output.txt")

```

This example shows how to use *field-aware factorizations machine (ffm)* to solve a simple binary classification task. You can check out the demo data (small_train.txt and small_test.txt) from the path demo/classification/criteo_ctr.

3.2 Installation Guide

For now, xLearn can support Linux and Mac OS X. We will support it on Windows platform in the near future. This page gives instructions on how to build and install the xLearn package using pip and how to build it from source code. No matter what way you choose, make sure that your OS has already installed GCC (or Clang) and CMake, and your compiler need to support C++11. If you have not installed them, please see [this page](#) on how to install GCC and CMake.

3.2.1 Install xLearn from pip

The easiest way to install xLearn Python package is to use pip. The following command will download the xLearn source code from pip and install Python package locally.

```
sudo pip install xlearn
```

The installation process will take a while to complete. And then you can type the following script in your python shell to check whether the xLearn has been installed successfully:

```
>>> import xlearn as xl
>>> xl.hello()
```

You will see:

```

-----
               |
            _-|_|
           / \ / | /  \ / \ | / | \ | /
          > <|_|_|  / \ / \ / \ | / \
         / \ \ _____ \ \ / \ / \ | / \
               |
          xLearn  -- 0.10 Version  --
-----

```

If you want to build the latest code from [Github](#), or you want to use the xLearn command line instead of the Python API, you can see how to build xLearn from source code as follow.

3.2.2 Install xLearn from Source Code

Building xLearn from source code consists two steps.

First, you need to build the executable files (`xlearn_train` and `xlearn_predict`), as well as the shared library (`libxlearn_api.so` for Linux and `libxlearn_api.dylib` for Mac OSX) from the C++ code.

Then, you can install the Python package through `install-python.sh`.

Fortunately, we write a script `build.sh` to do all the cumbersome work for users.

For users, you just need to clone the code from github

```
git clone https://github.com/aksnzhy/xlearn.git
```

and then build xLearn using the folloing commands:

```
cd xlearn
./build.sh
```

You may be asked to input your password during installation.

3.2.3 Test Your Building

Now you can test your installation by using the following command:

```
cd build
./run_example.sh
```

You can also test the Python package by using the following command:

```
cd python-package/test
python test_python.py
```

3.2.4 Install R Package

The R package installation guide is coming soon.

Install GCC or Clang

If you have already installed your compiler before, you can skip this step.

- On Cygwin, run `setup.exe` and install `gcc` and `binutils`.
- On Debian/Ubuntu Linux, type the command:

```
sudo apt-get install gcc binutils
```

to install GCC (or Clang) by using

```
sudo apt-get install clang
```

- On FreeBSD, type the following command to install Clang

```
sudo pkg_add -r clang
```

- On Mac OS X, install XCode gets you Clang.

Install CMake

If you have already installed CMake before, you can skip this step.

To install CMake from binary packages:

- On Cygwin, run `setup.exe` and install cmake.
- On Debian/Ubuntu Linux, type the command to install cmake:

```
sudo apt-get install cmake
```

- On FreeBSD, type the command:

```
sudo pkg_add -r cmake
```

On Mac OS X, if you have homebrew, you can use the command

```
brew install cmake
```

or if you have MacPorts, run

```
sudo port install cmake
```

You won't want to have both Homebrew and MacPorts installed.

3.3 xLearn Command Line Guide

Once you built xLearn from source code successfully, you will get two executable files `xlearn_train` and `xlearn_predict` in your build directory. Now you can use these two executable files to perform training and prediction tasks.

3.3.1 Quick Start

Make sure that you are in the build directory of xLearn, and you can find the demo data `small_test.txt` and `small_train.txt` in this directory. Now type the following command to train a model:

```
./xlearn_train ./small_train.txt
```

Here, we show a portion of the xLearn's output. Note that the loss value shown in your machine could be different.

Epoch	Train log_loss	Time cost (sec)
1	0.567514	0.00
2	0.516861	0.00
3	0.489884	0.00
4	0.469971	0.00
5	0.452699	0.00

6	0.437590	0.00
7	0.425759	0.00
8	0.415190	0.00
9	0.405954	0.00
10	0.396313	0.00

By default, xLearn will use the logistic regression (LR) to train our model within 10 epoch.

We can see that a new file called `small_train.txt.model` has been generated in the current directory. This file stores the trained model checkpoint, and we can use this model file to make prediction in the future

```
./xlearn_predict ./small_test.txt ./small_train.txt.model
```

After that, we can get a new file called `small_test.txt.out` in the current directory. This is output prediction. Here we show the first five lines of this output by using the following command

```
head -n 5 ./small_test.txt.out  
  
-1.9872  
-0.0707959  
-0.456214  
-0.170811  
-1.28986
```

These lines of data are the prediction score calculated for examples in the test set. The negative data represents the negative example and positive data represents the positive example. In xLearn, you can convert the score to (0-1) by using `--sigmoid` option, or you can convert your result to binary result (0 and 1) by using `--sign` option

```
./xlearn_predict ./small_test.txt ./small_train.txt.model --sigmoid  
head -n 5 ./small_test.txt.out  
  
0.120553  
0.482308  
0.387884  
0.457401  
0.215877  
  
./xlearn_predict ./small_test.txt ./small_train.txt.model --sign  
head -n 5 ./small_test.txt.out  
  
0  
0  
0  
0  
0
```

Users may want to generate different model files, so you can set the name of the model checkpoint file by using `-m` option. By default, the name of the model file equals to `training_data_name + .model`

```
./xlearn_train ./small_train.txt -m new_model
```

Also, users can save the model in txt format by using `-t` option. For example:

```
./xlearn_train ./small_train.txt -t model.txt
```


After that, we get a new file called `model.txt`, which stores the trained model in txt format. For now, xLearn only supports to save the bias and linear term in txt file.

```
head -n 5 ./model.txt

-0.688182
0.458082
0
0
0
```

Users can also set `-o` option to specify the output file. For example:

```
./xlearn_predict ./small_test.txt ./small_train.txt.model -o output.txt
head -n 5 ./output.txt

-2.01192
-0.0657416
-0.456185
-0.170979
-1.28849
```

By default, the name of the output file equals to `test_data_name + .out`.

3.3.2 Choose Machine Learning Algorithm

For now, xLearn can support three different machine learning algorithms, including LR, FM and FFM. Users can choose different machine learning algorithms by using `-s` option:

```
-s <type> : Type of machine learning model (default 0)
  for classification task:
    0 -- linear model (GLM)
    1 -- factorization machines (FM)
    2 -- field-aware factorization machines (FFM)
  for regression task:
    3 -- linear model (GLM)
    4 -- factorization machines (FM)
    5 -- field-aware factorization machines (FFM)
```

For LR and FM, the input data format can be CSV or libsvm. For FFM, the input data should be the libffm format.

```
libsvm format:

  label index_1:value_1 index_2:value_2 ... index_n:value_n

CSV format:

  value_1 value_2 .. value_n label

libffm format:

  label field_1:index_1:value_1 field_2:index_2:value_2 ...
```

Users can also give a libffm file to LR and FM. At that time, xLearn will treat this data as libsvm format. The following command shows how to use different machine learning algorithms to solve the binary classification problem:

```
./xlearn_train ./small_train.txt -s 0 # Linear model
./xlearn_train ./small_train.txt -s 1 # Factorization machine (FM)
./xlearn_train ./small_train.txt -s 2 # Field-aware factorization machine_
↪ (FFM)
```

3.3.3 Set Validation Dataset

A validation dataset is used to tune the hyperparameters of a machine learning model. In xLearn, users can use `-v` option to set the validation dataset. For example:

```
./xlearn_train ./small_train.txt -v ./small_test.txt
```

A portion of xLearn's output:

Epoch	Train log_loss	Test log_loss	Time cost (sec)
1	0.575049	0.530560	0.00
2	0.517496	0.537741	0.00
3	0.488428	0.527205	0.00
4	0.469010	0.538175	0.00
5	0.452817	0.537245	0.00
6	0.438929	0.536588	0.00
7	0.423491	0.532349	0.00
8	0.416492	0.541107	0.00
9	0.404554	0.546218	0.00

Here we can see that the training loss continuously goes down. But the validation loss (test loss) goes down first, and then goes up. This is because our model has already overfitted current training dataset. By default, xLearn will calculate the validation loss in each epoch, while users can also set different evaluation metrics by using `-x` option. For classification problems, the metric can be : `acc` (accuracy), `prec` (precision), `f1` (f1 score), `auc` (AUC score). For example:

```
./xlearn_train ./small_train.txt -v ./small_test.txt -x acc
./xlearn_train ./small_train.txt -v ./small_test.txt -x prec
./xlearn_train ./small_train.txt -v ./small_test.txt -x f1
./xlearn_train ./small_train.txt -v ./small_test.txt -x auc
```

For regression problems, the metric can be `mae`, `mape`, and `rmsd` (rmse). For example:

```
cd demo/house_price/
../../xlearn_train ./house_price_train.txt -s 3 -x rmse --cv
../../xlearn_train ./house_price_train.txt -s 3 -x rmsd --cv
```

3.3.4 Cross-Validation

Cross-validation, sometimes called rotation estimation, is a model validation technique for assessing how the results of a statistical analysis will generalize to an independent dataset. In xLearn, users can use the `--cv` option to use this technique. For example:

```
cd build
./xlearn_train ./small_train.txt --cv
```

On default, xLearn uses 5-folds cross validation, and users can set the number of fold by using `-f` option:

```
./xlearn_train ./small_train.txt -f 3 --cv
```

Here we set the number of folds to 3. The xLearn will calculate the average validation loss at the end of its output message.

```
[-----] Average log_loss: 0.549417
[ ACTION   ] Finish Cross-Validation
[ ACTION   ] Clear the xLearn environment ...
[-----] Total time cost: 0.03 (sec)
```

3.3.5 Choose Optimization Method

In xLearn, users can choose different optimization methods by using `-p` option. For now, users can choose `sgd`, `adagrad`, and `ftrl` method. By default, xLearn uses the `adagrad` method. For example:

```
./xlearn_train ./small_train.txt -p sgd
./xlearn_train ./small_train.txt -p adagrad
./xlearn_train ./small_train.txt -p ftrl
```

Compared to traditional `sgd` method, `adagrad` adapts the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameters. For this reason, it is well-suited for dealing with sparse data. In addition, `sgd` is more sensitive to the learning rate compared with `adagrad`.

FTRL (Follow-the-Regularized-Leader) is also a famous method that has been widely used in the large-scale sparse problem. To use FTRL, users need to tune more hyperparameters compared with `sgd` and `adagrad`.

3.3.6 Hyperparameter Tuning

In machine learning, a *hyperparameter* is a parameter whose value is set before the learning process begins. By contrast, the value of other parameters is derived via training. Hyperparameter tuning is the problem of choosing a set of optimal hyperparameters for a learning algorithm.

First, the `learning_rate` is one of the most important hyperparameters used in machine learning. By default, this value is set to `0.2`, and we can tune this value by using `-r` option:

```
./xlearn_train ./small_train.txt -v ./small_test.txt -r 0.1
./xlearn_train ./small_train.txt -v ./small_test.txt -r 0.5
./xlearn_train ./small_train.txt -v ./small_test.txt -r 0.01
```

We can also use the `-b` option to perform regularization. By default, xLearn uses L2 regularization, and the `regular_lambda` has been set to `0.00002`.

```
./xlearn_train ./small_train.txt -v ./small_test.txt -r 0.1 -b 0.001
./xlearn_train ./small_train.txt -v ./small_test.txt -r 0.1 -b 0.002
./xlearn_train ./small_train.txt -v ./small_test.txt -r 0.1 -b 0.01
```

For the FTRL method, we also need to tune another four hyperparameters, including `-alpha`, `-beta`, `-lambda_1`, and `-lambda_2`. For example:

```
./xlearn_train ./small_train.txt -p ftrl -alpha 0.002 -beta 0.8 -lambda_1 0.
↪001 -lambda_2 1.0
```

For FM and FFM, users also need to set the size of *latent factor* by using `-k` option. By default, xLearn uses 4 for this value.

```
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt -k 2
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt -k 4
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt -k 5
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt -k 8
```

xLearn uses *SSE* instruction to accelerate vector operation, and hence the time cost for $k=2$ and $k=4$ are the same.

For FM and FFM, users can also set the hyperparameter `-u` for model initialization. By default, this value is set to 0.66.

```
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt -u 0.80
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt -u 0.40
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt -u 0.10
```

3.3.7 Set Epoch Number and Early-Stopping

For machine learning, one epoch consists of one full training cycle on the training set. In xLearn, users can set the number of epoch for training by using `-e` option.

```
./xlearn_train ./small_train.txt -e 3
./xlearn_train ./small_train.txt -e 5
./xlearn_train ./small_train.txt -e 10
```

If you set the validation data, xLearn will perform early-stopping by default. For example:

```
./xlearn_train ./small_train.txt -s 2 -v ./small_test.txt -e 10
```

Here, we set epoch number to 10, but xLearn stopped at epoch 7 because we get the best model at that epoch (you may get different a stopping number on your machine)

```
[ ACTION      ] Early-stopping at epoch 7
[ ACTION      ] Start to save model ...
```

Users can disable early-stopping by using `--dis-es` option

```
./xlearn_train ./small_train.txt -s 2 -v ./small_test.txt -e 10 --dis-es
```

At this time, xLearn performed 10 epoch for training.

3.3.8 Lock-Free Training

By default, xLearn performs *Hogwild!* lock-free training, which takes advantages of multiple cores to accelerate training task. But lock-free training is *non-deterministic*. For example, if we run the following command multiple times, we may get different loss value at each epoch.

```
./xlearn_train ./small_train.txt

The 1st time: 0.396352
The 2nd time: 0.396119
The 3rd time: 0.396187
...
```

Users can set the number of thread for xLearn by using `-nthread` option:

```
./xlearn_train ./small_train.txt -nthread 2
```

If you don't set this option, xLearn uses all of the CPU cores by default.

Users can disable lock-free training by using `--dis-lock-free`

```
./xlearn_train ./small_train.txt --dis-lock-free
```

In this time, our result are *deterministic*.

```
The 1st time: 0.396372
The 2nd time: 0.396372
The 3rd time: 0.396372
```

The disadvantage of `--dis-lock-free` is that it is much slower than lock-free training.

3.3.9 Instance-wise Normalization

For FM and FFM, xLearn uses *instance-wise normalization* by default. In some scenes like CTR prediction, this technique is very useful. But sometimes it hurts model performance. Users can disable instance-wise normalization by using `--no-norm` option

```
./xlearn_train ./small_train.txt -s 1 -v ./small_test.txt --no-norm
```

Note that we usually use `--no-norm` in regression tasks.

3.3.10 Quiet Training

When using `--quiet` option, xLearn will not calculate any evaluation information during the training, and it just train the model quietly

```
./xlearn_train ./small_train.txt --quiet
```

In this way, xLearn can accelerate its training speed.

3.4 xLearn Python Package Guide

xLearn supports very easy-to-use Python API for users. Once you install the xLearn Python package successfully, you can try it. Type `python` in your shell and type the following Python code to check your installation:

```
import xlearn as xl
xl.hello()
```

If you install xLearn Python package successfully, you will see

```
-----
      _
     / \
    /   \
   /     \
  /       \
 /         \
/           \
 \         /
  \       /
   \     /
    \   /
     \ /
      _
```

```
> <| |__| __/ (| | | | | | |
/_/\_\____/\____|\_,_|_| |_| |_|

xLearn -- 0.20 Version --
-----
```

3.4.1 Quick Start

Here is a simple Python demo to demonstrate how to use xLearn. You can check out the demo data (small_train.txt and small_test.txt) from the path demo/classification/criteo_ctr.

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm() # Use field-aware factorization machine
ffm_model.setTrain("./small_train.txt") # Training data
ffm_model.setValidate("./small_test.txt") # Validation data

# param:
# 0. binary classification
# 1. learning rate : 0.2
# 2. regular lambda : 0.002
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}

# Train model
ffm_model.fit(param, "./model.out")
```

A portion of the xLearn's output

```
Start to train ...
Epoch      Train log_loss      Test log_loss      Time cost (sec)
  1          0.593750          0.535847           0.00
  2          0.539226          0.543829           0.00
  3          0.520034          0.531732           0.00
  4          0.505186          0.537418           0.00
  5          0.494089          0.533448           0.00
  6          0.483678          0.534629           0.00
  7          0.470848          0.528086           0.00
  8          0.466330          0.533253           0.00
  9          0.456660          0.535635           0.00
Early-stopping at epoch 7
Start to save model ...
```

In this example, xLearn uses *field-aware factorization machines* (ffm) to train our model for solving a binary classification task. If you want train a model for regression task. You can reset the `task` parameter to `reg`.

```
param = {'task':'reg', 'lr':0.2, 'lambda':0.002}
```

We can see that a new file called `model.out` has been generated in the current directory. This file stores the trained model checkpoint, and we can use this model file to make prediction in the future:

```
ffm_model.setTest("./small_test.txt")
ffm_model.predict("./model.out", "./output.txt")
```

After we run this Python code, we can get a new file called `output.txt` in current directory. This is output prediction. Here we show the first five lines of this output by using the following command

```
head -n 5 ./output.txt

-1.66107
-0.616408
-0.815918
-0.608931
-1.30794
```

These lines of data are the prediction score calculated for examples in the test set. The negative data represents the negative example and positive data represents the positive example. In xLearn, you can convert the score to (0-1) by using `setSigmoid()` option:

```
ffm_model.setTest("./small_test.txt")
ffm_model.setSigmoid()
ffm_model.predict("./model.out", "./output.txt")
```

and then we can get the result

```
head -n 5 ./output.txt

0.158613
0.354297
0.310193
0.357449
0.220061
```

We can also convert the score to binary result (0 and 1) by using `setSign()` API

```
# Prediction task
ffm_model.setTest("./small_test.txt")
ffm_model.setSign()
ffm_model.predict("./model.out", "./output.txt")
```

and then we can get the result

```
head -n 5 ./output.txt

0
0
0
0
0
```

3.4.2 Choose Machine Learning Algorithm

For now, xLearn can support three different machine learning algorithms, including LR, FM and FFM. Users can choose different machine learning algorithms by using `create_xxx()` API:

```
import xlearn as xl

ffm_model = xl.create_ffm()
fm_model = xl.create_fm()
lr_model = xl.create_lr()
```

For LR and FM, the input data format can be CSV or libsvm. For FFM, the input data should be the libffm format.

```
libsvm format:
    label index_1:value_1 index_2:value_2 ... index_n:value_n

CSV format:
    value_1 value_2 .. value_n label

libffm format:
    label field_1:index_1:value_1 field_2:index_2:value_2 ...
```

Users can also give a libffm file to LR and FM. At that time, xLearn will treat this data as libsvm format.

3.4.3 Set Validation Dataset

A validation dataset is used to tune the hyperparameters of a machine learning model. In xLearn, users can use `setValidate()` API to set the validation dataset. For example:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.setValidate("./small_test.txt")
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}

ffm_model.fit(param, "./model.out")
```

A portion of xLearn's output:

Epoch	Train log_loss	Test log_loss	Time cost (sec)
1	0.598814	0.536327	0.00
2	0.539872	0.542924	0.00
3	0.521035	0.531595	0.00
4	0.505414	0.536246	0.00
5	0.492150	0.532070	0.00
6	0.482229	0.536482	0.00
7	0.470457	0.528871	0.00
8	0.464445	0.534550	0.00
9	0.456061	0.537320	0.00

Here we can see that the training loss continuously goes down. But the validation loss (test loss) goes down first, and then goes up. This is because our model has already overfitted current training dataset. By default, xLearn will calculate the validation loss in each epoch, while users can also set different evaluation metrics by using `metric` parameter. For classification problems, the metric can be : `acc` (accuracy), `prec` (precision), `f1` (f1 score), and `auc` (AUC score). For example:

```
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'metric': 'acc'}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'metric': 'prec'}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'metric': 'f1'}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'metric': 'auc'}
```

For regression problems, the metric can be `mae`, `mape`, and `rmsd` (rmse). For example:


```
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'metric': 'rmse'}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'metric': 'mae'}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'metric': 'mape'}
```

3.4.4 Cross-Validation

Cross-validation, sometimes called rotation estimation, is a model validation technique for assessing how the results of a statistical analysis will generalize to an independent dataset. In xLearn, users can use the `cv()` API to use this technique. For example:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}

ffm_model.cv(param)
```

On default, xLearn uses 5-folds cross validation, and users can set the number of fold by using the `fold` parameter:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'fold':3}

ffm_model.cv(param)
```

Here we set the number of folds to 3. The xLearn will calculate the average validation loss at the end of its output message.

```
[-----] Average log_loss: 0.547632
[ ACTION   ] Finish Cross-Validation
[ ACTION   ] Clear the xLearn environment ...
[-----] Total time cost: 0.05 (sec)
```

3.4.5 Choose Optimization Method

In xLearn, users can choose different optimization methods by using `opt` parameter. For now, users can choose `sgd`, `adagrad`, and `ftrl` method. By default, xLearn uses the `adagrad` method. For example:

```
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'opt':'sgd'}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'opt':'adagrad'}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'opt':'ftrl'}
```

Compared to traditional `sgd` method, `adagrad` adapts the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameters. For this reason, it is well-suited for dealing with sparse data. In addition, `sgd` is more sensitive to the learning rate compared with `adagrad`.

FTRL (Follow-the-Regularized-Leader) is also a famous method that has been widely used in the large-scale sparse problem. To use FTRL, users need to tune more hyperparameters compared with `sgd` and

adagard.

3.4.6 Hyperparameter Tuning

In machine learning, a *hyperparameter* is a parameter whose value is set before the learning process begins. By contrast, the value of other parameters is derived via training. Hyperparameter tuning is the problem of choosing a set of optimal hyperparameters for a learning algorithm.

First, the `learning_rate` is one of the most important hyperparameters used in machine learning. By default, this value is set to 0.2, and we can tune this value by using `lr` parameter:

```
param = {'task':'binary', 'lr':0.2}
param = {'task':'binary', 'lr':0.5}
param = {'task':'binary', 'lr':0.01}
```

We can also use the `lambda` parameter to perform regularization. By default, xLearn uses L2 regularization, and the `regular_lambda` has been set to 0.00002.

```
param = {'task':'binary', 'lr':0.2, 'lambda':0.01}
param = {'task':'binary', 'lr':0.2, 'lambda':0.02}
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}
```

For the FTRL method, we also need to tune another four hyperparameters, including `alpha`, `beta`, `lambda_1`, and `lambda_2`. For example:

```
param = {'alpha':0.002, 'beta':0.8, 'lambda_1':0.001, 'lambda_2': 1.0}
```

For FM and FFM, users also need to set the size of latent factor by using `k` parameter. By default, xLearn uses 4 for this value.

```
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'k':2}
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'k':4}
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'k':5}
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'k':8}
```

xLearn uses *SSE* instruction to accelerate vector operation, and hence the time cost for `k=2` and `k=4` are the same.

For FM and FFM, users can also set the parameter `init` for model initialization. By default, this value is set to 0.66.

```
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'init':0.5} param = {'task':'binary', 'lr':0.2,
'lambda':0.01, 'init':0.8}
```

3.4.7 Set Epoch Number and Early-Stopping

For machine learning, one epoch consists of one full training cycle on the training set. In xLearn, users can set the number of epoch for training by using `epoch` option.

```
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'epoch':3}
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'epoch':5}
param = {'task':'binary', 'lr':0.2, 'lambda':0.01, 'epoch':10}
```

If you set the validation data, xLearn will perform early-stopping by default. For example:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.setValidate("./small_test.txt")
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'epoch':10}

ffm_model.fit(param, "./model.out")
```

Here, we set epoch number to 10, but xLearn stopped at epoch 7 because we get the best model at that epoch (you may get different a stopping number on your machine)

```
Early-stopping at epoch 7
Start to save model ...
```

Users can disable early-stopping by using `disableEarlyStop()` API:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.setValidate("./small_test.txt")
ffm_model.disableEarlyStop();
param = {'task':'binary', 'lr':0.2, 'lambda':0.002, 'epoch':10}

ffm_model.fit(param, "./model.out")
```

At this time, xLearn performed 10 epoch for training.

3.4.8 Lock-Free Training

By default, xLearn performs *Hogwild!* lock-free training, which takes advantages of multiple cores to accelerate training task. But lock-free training is *non-deterministic*. For example, if we run the following Python code multiple times, we may get different loss value at each epoch.

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}

ffm_model.fit(param, "./model.out")

The 1st time: 0.449056
The 2nd time: 0.449302
The 3rd time: 0.449185
```

Users can disable lock-free training by using `disableLockFree()` API.

```
import xlearn as xl

# Training task
```

```
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.disableLockFree()
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}

ffm_model.fit(param, "./model.out")
```

In this time, our result are *deterministic*.

```
The 1st time: 0.449172
The 2nd time: 0.449172
The 3rd time: 0.449172
```

The disadvantage of `disableLockFree()` is that it is much slower than lock-free training.

3.4.9 Instance-wise Normalization

For FM and FFM, xLearn uses instance-wise normalization by default. In some scenes like CTR prediction, this technique is very useful. But sometimes it hurts model performance. Users can disable *instance-wise normalization* by using `disableNorm()` API.

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.disableNorm()
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}

ffm_model.fit(param, "./model.out")
```

Note that we usually use `disableNorm` in regression tasks.

3.4.10 Quiet Training

When using `setQuiet()` API, xLearn will not calculate any evaluation information during the training, and it just train the model quietly

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.setQuiet()
param = {'task':'binary', 'lr':0.2, 'lambda':0.002}

ffm_model.fit(param, "./model.out")
```

In this way, xLearn can accelerate its training speed.

3.5 xLearn R Package Guide

xLearn R package guide is coming soon.

3.6 xLearn API List

This page gives the xLearn API List for the command line, Python package, and R package.

3.6.1 xLearn Command Line Usage

For Training:

```
xlearn_train <train_file_path> [OPTIONS]
```

Options:

```
-s <type> : Type of machine learning model (default 0)
  for classification task:
    0 -- linear model (GLM)
    1 -- factorization machines (FM)
    2 -- field-aware factorization machines (FFM)
  for regression task:
    3 -- linear model (GLM)
    4 -- factorization machines (FM)
    5 -- field-aware factorization machines (FFM)

-x <metric> : The metric can be 'acc', 'prec', 'recall', 'f1', 'auc
↳ ' for classification, and
↳ 'mae', 'mape', 'rmsd (rmse)' for regression. On
↳ default, xLearn will not print
↳ any evaluation metric information (only loss value).

-p <opt_method> : Choose the optimization method, including 'sgd',
↳ 'adagrad', and 'ftrl'. On default,
↳ we use the 'adagrad' optimization.

-v <validate_file> : Path of the validation data file. This option will
↳ be empty by default,
↳ and in this way, the xLearn will not perform
↳ validation process.

-m <model_file> : Path of the model dump file. On default, the model
↳ file name is 'train_file' + '.model'.
↳ If we set this value to 'none', the xLearn will not
↳ dump the model checkpoint after training.

-t <txt_model_file> : Path of the txt model checkpoint file. On default,
↳ we do not set this option
↳ and xLearn will not dump the txt model. For now,
↳ only the bias and linear term
↳ will be dump to the txt model file.

-l <log_file> : Path of the log file. Using '/tmp/xlearn_log.*' by
↳ default.
```

```

-k <number_of_K>      : Number of the latent factor used by fm and ffm tasks.
  ↳ Using 4 by default.
      Note that, we will get the same model size when
  ↳ setting k to 1 and 4.
      This is because we use SSE instruction and the
  ↳ memory need to be aligned.
      So even you assign k = 1, we still fill some dummy
  ↳ zeros from k = 2 to 4.

-r <learning_rate>    : Learning rate for optimization method. Using 0.2 by
  ↳ default.
      xLearn can use adaptive gradient descent (AdaGrad)
  ↳ for optimization problem,
      if you choose AdaGrad method, the learning rate will
  ↳ be changed adaptively.

-b <lambda_for_regu>  : Lambda for L2 regular. Using 0.00002 by default. We
  ↳ can disable the
      regular term by setting this value to 0.0

-u <model_scale>      : Hyper parameter used for initialize model parameters.
  ↳ Using 0.66 by default.

-e <epoch_number>     : Number of epoch for training. Using 10 by default.
  ↳ Note that, xLearn will
      perform early-stopping by default, so this value is
  ↳ just a upper bound.

-f <fold_number>      : Number of folds for cross-validation. Using 5 by
  ↳ default.

-nthread <thread number> : Number of thread for multi-thread training.

--disk                : Open on-disk training for large-scale machine
  ↳ learning problems.

--cv                  : Open cross-validation in training tasks. If we use
  ↳ this option, xLearn
      will ignore the validation file (-t).

--dis-lock-free       : Disable lock-free training. Lock-free training can
  ↳ accelerate training but
      the result is non-deterministic. Our suggestion is
  ↳ that you can open this flag
      if the training data is big and sparse.

--dis-es              : Disable early-stopping in training. By default,
  ↳ xLearn will use early-stopping
      in training tasks, except for training in cross-
  ↳ validation.

--no-norm              : Disable instance-wise normalization. By default,
  ↳ xLearn will use
      instance-wise normalization for both training and
  ↳ prediction.

--quiet                : Don't print any evaluation information during the
  ↳ training and

```

```

                                just train the model quietly.
-alpha                          : Hyper parameters used by ftrl.
-beta                            : Hyper parameters used by ftrl.
-lambda_1                        : Hyper parameters used by ftrl.
-lambda_2                        : Hyper parameters used by ftrl.

```

For Prediction:

```
xlearn_predict <test_file> <model_file> [OPTIONS]
```

Options:

```

-o <output_file>                : Path of the output file. On default, this value will
  ↳be set                        to 'test_file' + '.out'
-l <log_file_path>              : Path of the log file. Using '/tmp/xlearn_log' by
  ↳default.
-nthread <thread number>       : Number of thread for multi-thread training.
--sign                          : Converting output to 0 and 1.
--sigmoid                       : Converting output to 0~1 (problebility).

```

3.6.2 xLearn Python API

API List:

```

import xlearn as xl          # Import xlearn package
xl.hello()                     # Say hello to user
model = create_linear()        # Create linear model.
model = create_fm()           # Create factorization machines.
model = create_ffm()          # Create field-aware factorizarion machines.
model.show()                  # Show model information.
model.fit(param, "model_path") # Train model.
model.cv(param)               # Perform cross-validation.
model.predict("model_path", "output_path") # Perform prediction.
model.setTrain("data_path")   # Set training data for xLearn.
model.setValidate("data_path") # Set validation data for xLearn.
model.setTest("data_path")     # Set test data for xLearn.

```

```
model.setQuiet()      # Set xlearn to train model quietly.
model.setOnDisk()    # Set xlearn to use on-disk training.
model.setSign()      # Convert prediction to 0 and 1.
model.setSigmoid()   # Convert prediction to (0, 1).
model.disableNorm()  # Disable instance-wise normalization.
model.disableLockFree() # Disable lock-free training.
model.disableEarlyStop() # Disable early-stopping.
```

Parameter List:

```
task      : {'binary', 'reg'} # machine learning task
metric    : {'acc', 'prec', 'recall',
             'f1', 'mae', 'mape', 'rmse', 'rmsd'} # Evaluation metric
lr        : float value # learning rate
lambda   : float value # regular lambda
k         : int value # latent factor
init      : float value # model initialize
alpha    : float value # parameter for ftrl
beta     : float value # parameter for ftrl
lambda_1  : float value # parameter for ftrl
lambda_2  : float value # parameter for ftrl
epoch    : int vlaue # number of epoch
fold     : int value # number of fold for cross-validation
opt      : {'sgd', 'agagrad', 'ftrl'} # optimization method
```

3.6.3 xLearn R API

xLearn R API page is coming soon.

3.7 Large-Scale Machine Learning

This page shows how to use xLearn to solve large-scale machine learning problems. In recent years, challenges arise with the fast-growing data. For “big-data”, we focus on datasets with potentially trillions of training examples, which cannot fit into the memory of a single machine. Motivated by this, we design xLearn to solve large-scale machine learning problems. First, xLearn can handle very large data (TB) on a single PC by using *out-of-core* learning. In addition, xLearn can scale beyond billions of example across many machines to support distributed learning by using the *parameter server* framework.

3.7.1 Out-of-Core Learning

Out-of-core learning refers to the machine learning algorithms working with data cannot fit into the memory of a single machine, but that can easily fit into some data storage such as local hard disk or web repository. Your available RAM, the core memory on your single machine, may indeed range from a few gigabytes (sometimes 2 GB, more commonly 4 GB, but we assume that you have 2 GB at maximum) up

to 256 GB on large server machines. Large servers are like the ones you can get on cloud computing services such as Amazon Elastic Compute Cloud (EC2), whereas your storage capabilities can easily exceed terabytes of capacity using just an external drive (most likely about 1 TB but it can reach up to 4 TB).

Actually, the ability to learn incrementally from a mini-batch of instances is key to *out-of-core* learning as it guarantees that at any given time there will be only a small amount of data in the main memory. Choose a good size for the mini-batch that balances relevancy and memory footprint could involve some tuning.



Out-of-Core Learning Using xLearn Command Line

It's very easy to perform *out-of-core* learning in xLearn command line, where users can just use the `--disk` option, and xLearn will help you do all the other things. For example:

```
./xlearn_train ./big_data.txt -s 2 --disk
```

Epoch	Train log_loss	Time cost (sec)
1	0.483997	4.41
2	0.466553	4.56
3	0.458234	4.88
4	0.451463	4.77
5	0.445169	4.79
6	0.438834	4.71
7	0.432173	4.84
8	0.424904	4.91
9	0.416855	5.03
10	0.407846	4.53

In this example, xLearn can finish the training of each epoch in nearly 4.5 second. If you delete the `--disk` option, xLearn can train faster.

```
./xlearn_train ./big_data.txt -s 2
```

Epoch	Train log_loss	Time cost (sec)
1	0.484022	1.65
2	0.466452	1.64
3	0.458112	1.64
4	0.451371	1.76
5	0.445040	1.83

6	0.438680	1.92
7	0.432007	1.99
8	0.424695	1.95
9	0.416579	1.96
10	0.407518	2.11

In this time, the training of each epoch will only spend nearly 1.8 seconds.

Out-of-Core Learning Using xLearn Python API

In Python, users can use `setOnDisk` API to perform *out-of-core* learning. For example:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setOnDisk()
ffm_model.setTrain("./small_train.txt")
ffm_model.setValidate("./small_test.txt")
param = {'task':'binary', 'lr':0.2,
         'lambda':0.002, 'metric':'auc'}

ffm_model.fit(param, "./model.out")

# Prediction task
ffm_model.setTest("./small_test.txt")
# Convert output to 0~1
ffm_model.setSigmoid()
ffm_model.predict("./model.out", "./output.txt")
```

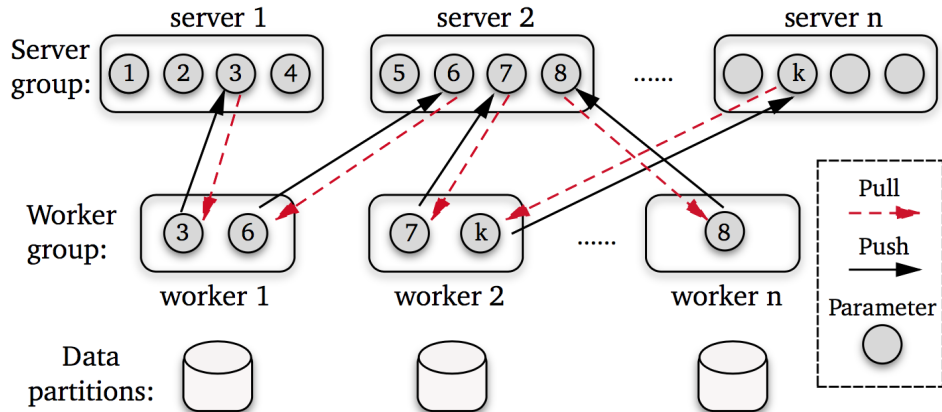
Out-of-Core Learning Using xLearn R API

The R guide is coming soon.

3.7.2 Distributed Learning

As we mentioned before, for some large-scale machine challenges like computational advertising, we focus on the problem with potentially trillions of training examples and billions of model parameters, both of which cannot fit into the memory of a single machine, which brings the *scalability challenge* for users and system designer. For this challenge, parallelizing the training process across machines has become a prerequisite.

The *Parameter Server* (PS) framework has emerged as an efficient approach to solve the “big model” machine learning challenge recently. Under this framework, both the training data and workloads are spread across worker nodes, while the server nodes maintain the globally shared model parameters. The following figure demonstrates the architecture of the PS framework.



As we can see, the *Parameter Server* provides two concise APIs for users.

Push sends a vector of (key, value) pairs to the server nodes. To be more specific – in the distributed gradient descent, the worker nodes might send the locally computed gradients to servers. Due to the data sparsity, only a part of the gradients is non-zero. Often it is desirable to present the gradient as a list of (key, value) pairs, where the feature index is the key and the according gradient item is value.

Pull requests the values associated with a list of keys, which will get the newest parameters from the server nodes. This is particularly useful whenever the main memory of a single worker cannot hold a full model. Instead, workers prefetch the model entries relevant for solving the model only when needed.

The distributed training guide for xLearn is coming soon.

3.8 xLearn Demo

Copyright of the dataset belongs to the original copyright holder.

3.8.1 Criteo CTR Prediction

Predict click-through rates on display ads ([Link](#))

Display advertising is a billion dollar effort and one of the central uses of machine learning on the Internet. However, its data and methods are usually kept under lock and key. In this research competition, CriteoLabs is sharing a week's worth of data for you to develop models predicting ad click-through rate (CTR). Given a user and the page he is visiting, what is the probability that he will click on a given ad?

You can find the data used in this demo in the path `/demo/classification/criteo_ctr/`.

The follow code is the Python demo:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./small_train.txt")
ffm_model.setValidate("./small_test.txt")
param = {'task':'binary', 'lr':0.2,
         'lambda':0.002, 'metric':'auc'}

ffm_model.fit(param, "./model.out")
```

```
# Prediction task
ffm_model.setTest("./small_test.txt")
# Convert output to 0~1
ffm_model.setSigmoid()
ffm_model.predict("./model.out", "./output.txt")
```

3.8.2 Mushroom Classification

This dataset comes from UCI Machine Learning Repository ([Link](#))

This data set includes descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms in the Agaricus and Lepiota Family (pp. 500-525). Each species is identified as definitely edible, definitely poisonous, or of unknown edibility and not recommended. This latter class was combined with the poisonous one. The Guide clearly states that there is no simple rule for determining the edibility of a mushroom; no rule like *leaflets three, let it be* for Poisonous Oak and Ivy.

You can find a small portion of data used in this demo in the path `/demo/classification/mushroom/`.

The follow code is the Python demo:

```
import xlearn as xl

# Training task
linear_model = xl.create_linear()
linear_model.setTrain("./agaricus_train.txt")
linear_model.setValidate("./agaricus_test.txt")
param = {'task':'binary', 'lr':0.2,
         'lambda':0.002, 'metric':'acc',
         'opt':'sgd'}

linear_model.fit(param, './model.out')

# Prediction task
linear_model.setTest("./agaricus_test.txt")
# Convert output to 0-1
linear_model.setSigmoid()
linear_model.predict("./model.out", "./output.txt")
```

3.8.3 Predict Survival in Titanic

This challenge comes from the Kaggle. In this challenge, we ask you to complete the analysis of what sorts of people were likely to survive. In particular, we ask you to apply the tools of machine learning to predict which passengers survived the tragedy. ([Link](#))

You can find the data used in this demo in the path `/demo/classification/titanic/`.

The follow code is the Python demo:

```
import xlearn as xl

# Training task
fm_model = xl.create_fm()
fm_model.setTrain("./titanic_train.txt")
param = {'task':'binary', 'lr':0.2,
         'lambda':0.002, 'metric':'acc'}
```

```
# Cross-validation
ffm_model.cv(param)
```

3.8.4 House Price Prediction

This demo shows how to use xLearn to solve the regression problem, and it comes from the Kaggle. The Ames Housing dataset was compiled by Dean De Cock for use in data science education. It's an incredible alternative for data scientists looking for a modernized and expanded version of the often cited Boston Housing dataset. ([Link](#))

You can find the data used in this demo in the path `/demo/regression/house_price/`.

The follow code is the Python demo:

```
import xlearn as xl

# Training task
ffm_model = xl.create_ffm()
ffm_model.setTrain("./house_price_train.txt")
param = {'task':'reg', 'lr':0.2,
         'lambda':0.002, 'metric':'rmse'}

# Cross-validation
ffm_model.cv(param)
```

More Demo in xLearn is coming soon.

3.9 xLearn Tutorials

1. From linear model to FM and FFM (Chinese version)

Tutorials on specific tasks is coming soon.