
SMURFF

Release v0.17.0-1-g00f12c67

Apr 09, 2021

Contents

1	Contents	3
1.1	What is SMURFF	3
1.2	Getting Started using IPython Notebooks	3
1.3	Compilation of SMURFF	16
1.4	Python API Reference	17
	Python Module Index	27
	Index	29

Welcome to SMURFF's documentation. SMURFF is a highly optimized and parallelized framework for Bayesian Matrix and Tensors Factorization.

The easiest way to install SMURFF is to use conda:

```
conda install -c vanderaa smurff
```

To get started with SMURFF using IPython notebooks, have a look at our [Getting Started using IPython Notebooks](#). To learn more about Matrix Factorization, have a look at [What is SMURFF](#).

1.1 What is SMURFF

SMURFF is a highly optimized and parallelized framework for Bayesian Matrix and Tensors Factorization.

1.1.1 What is Bayesian Matrix Factorization

Matrix factorization is a common machine learning technique for recommender systems, like books for Amazon or movies for Netflix.

The idea of these methods is to approximate the user-movie rating matrix R as a product of two low-rank matrices U and V such that $R \approx U \times V$. In this way U and V are constructed from the known ratings in R , which is usually very sparsely filled. The recommendations can be made from the approximation $U \times V$ which is dense. If $M \times N$ is the dimension of R then U and V will have dimensions $M \times K$ and $N \times K$. Noise is added to the model to avoid overfitting.

1.1.2 What can SMURFF do for you

SMURFF supports multiple matrix factorization methods:

- **BPMF**, the basic version;
- **Macau**, adding support for high-dimensional side information to the factorization;
- **GFA**, doing Group Factor Analysis.

Macau and BPMF can also perform **tensor** factorization.

1.2 Getting Started using IPython Notebooks

This section contains documentation generated from IPython notebooks that discuss different aspects of SMURFF.

Note: This page was generated from [notebooks/a_first_example.ipynb](#).

1.2.1 A first example running SMURFF

In this notebook we will run the BPMF algorithm using SMURFF, on compound-activity data.

Downloading the data files

In these examples we use ChEMBL dataset for compound-proteins activities (IC50). The IC50 values and ECFP fingerprints can be downloaded using this smurff function:

```
[ ]: import logging
      logging.basicConfig(level = logging.INFO)

      import smurff

      ic50_train, ic50_test, ecfp = smurff.load_chembl()
```

The resulting variables are all `scipy.sparse` matrices: `ic50` is a sparse matrix containing interactions between chemical compounds (in the rows) and protein targets (called essays - in the columns). The matrix is already split in as train and test set.

The `ecfp` contains compound features. These features will not be used in this example.

Having a look at the data

The `spy` function in `matplotlib` is a handy function to plot sparsity pattern of a matrix.

```
[ ]: %matplotlib inline

      import matplotlib.pyplot as plt

      fig = plt.figure()
      ax = fig.add_subplot(111)
      ax.spy(ic50_train.tocsr()[0:1000,:].T, markersize = 1)
```

Running SMURFF

Finally we run make a BPMF training `trainSession` and call `run`. The `run` function builds the model and returns the predictions of the test data.

```
[ ]: trainSession = smurff.BPMFSession(
      Ytrain      = ic50_train,
      Ytest       = ic50_test,
      num_latent  = 16,
      burnin      = 40,
      nsamples    = 200,
      verbose     = 0,
      checkpoint_freq = 1,
      save_freq   = 1,)
```

(continues on next page)

(continued from previous page)

```
predictions = trainSession.run()
```

We can use the `calc_rmse` function to calculate the RMSE.

```
[ ]: rmse = smurff.calc_rmse(predictions)
      rmse
```

Plotting predictions versus actual values

Next to RMSE, we can also plot the predicted versus the actual values, to see how well the model performs.

```
[ ]: import numpy
      from matplotlib.pyplot import subplots, show

      y = numpy.array([ p.val for p in predictions ])
      predicted = numpy.array([ p.pred_avg for p in predictions ])

      fig, ax = subplots()
      ax.scatter(y, predicted, edgecolors=(0, 0, 0))
      ax.plot([y.min(), y.max()], [y.min(), y.max()], 'k--', lw=4)
      ax.set_xlabel('Measured')
      ax.set_ylabel('Predicted')
      show()
```

```
[ ]:
```

Note: This page was generated from [notebooks/input_matrices_and_tensors.ipynb](#).

```
[ ]: import numpy as np
      import scipy.sparse as sp

      import logging
      logging.basicConfig(level = logging.INFO)

      import smurff
```

1.2.2 Input to SMURFF

In this notebook we will look at how to provide input to SMURFF with dense and sparse matrices;

SMURFF accepts the following matrix files for train, test and side-info data:

- for dense matrix or tensor input: `numpy.ndarrays`
- for sparse matrices input: `scipy Sparse matrices` in COO, CSR or CSC format
- for sparse tensors: a wrapper around a `pandas.DataFrame`

Let's have a look on how this could work.

Dense Train Input

```
[ ]: # dense input
Ydense = np.random.rand(10, 20)
trainSession = smurff.TrainSession(burnin = 5, nsamples = 5)
trainSession.addTrainAndTest(Ydense)
trainSession.run()
```

Sparse Matrix Input

The so-called *zero* elements in sparse matrices can either represent

1. missing values, also called ‘unknown’ or ‘not-available’ (NA) values.
2. actual zero values, to optimize the space that stores the matrix

Important: * when calling `addTrainAndTest(Ytrain, Ytest, is_scarce)` the `is_scarce` refers to the `Ytrain` matrix. `Ytest` is *always* scarce. * when calling `addSideInfo(mode, sideinfoMatrix)` with a sparse `sideinfoMatrix`, this matrix is always fully known.

```
[ ]: # sparse matrix input with 20% zeros (fully known)
Ysparse = sp.rand(15, 10, 0.2)
trainSession = smurff.TrainSession(burnin = 5, nsamples = 5)
trainSession.addTrainAndTest(Ysparse, is_scarce = False)
trainSession.run()
```

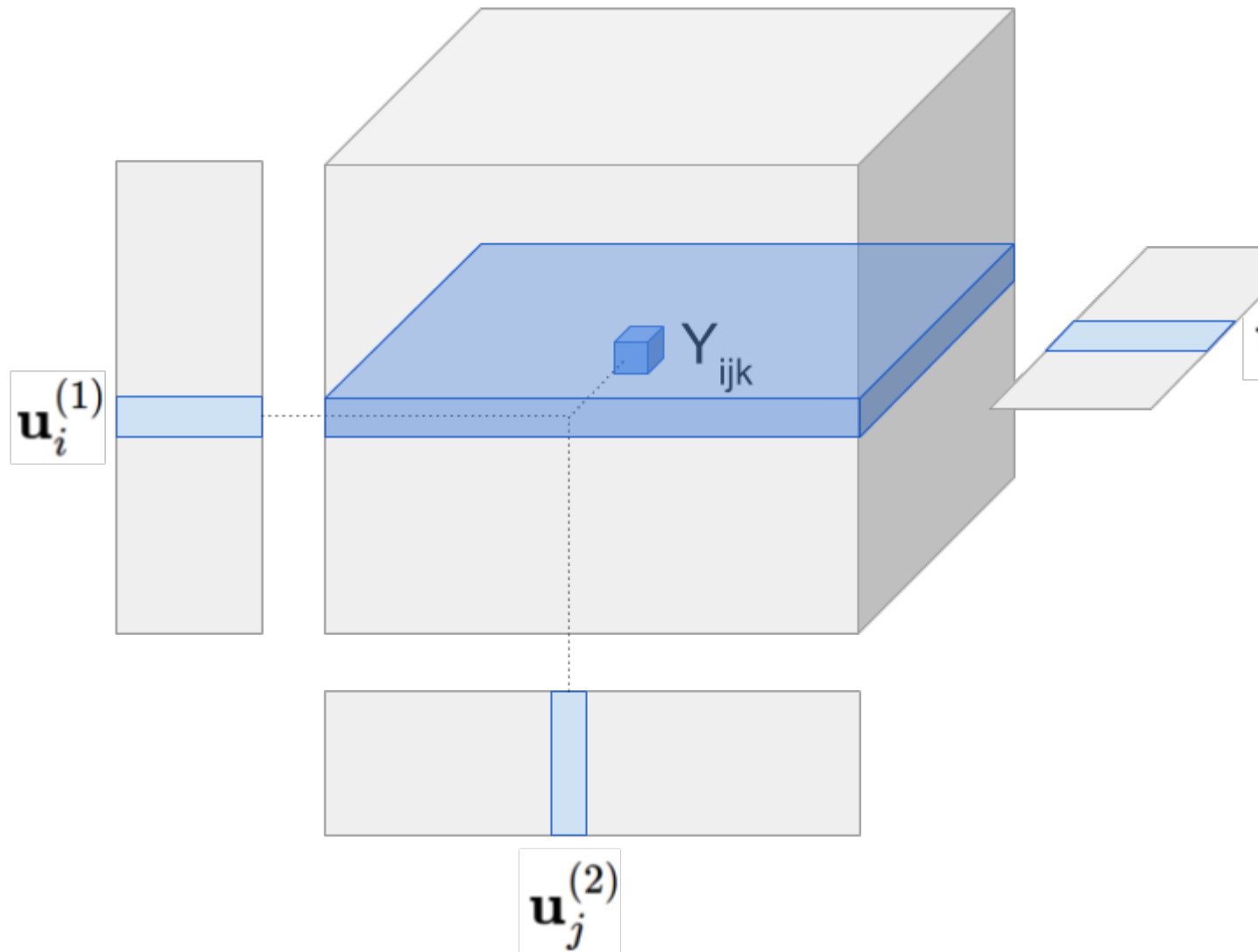
```
[ ]: # sparse matrix input with unknowns (the default)
Yscarce = sp.rand(15, 10, 0.2)
trainSession = smurff.TrainSession(burnin = 5, nsamples = 5)
trainSession.addTrainAndTest(Yscarce, is_scarce = True)
trainSession.run()
```

Tensor input

SMURFF also supports tensor factorization with and without side information on any of the modes. Tensor can be thought as generalization of matrix to relations with more than two items. For example 3-tensor of drug x cell x gene could express the effect of a drug on the given cell and gene. In this case the prediction for the element `Yhat[i, j, k]` is given by

$$\hat{Y}_{ijk} = \sum_{d=1}^D u_{d,i}^{(1)} u_{d,j}^{(2)} u_{d,k}^{(3)} + mean$$

Visually the model can be represented as follows:



Tensor model predicts \hat{Y}_{ijk} by multiplying all latent vectors together element-wise and then taking the sum along the latent dimension (figure omits the global mean).

For tensors SMURFF implements a `SparseTensor` class. `SparseTensor` can be constructed from a pandas `DataFrame` where each row stores the coordinate and the value of a known cell in the tensor. Specifically, the integer columns in the `DataFrame` give the coordinate of the cell and `float` (or `double`) column stores the value in the cell (the order of the columns does not matter). The coordinates are 0-based. The shape of the `SparseTensor` can be provided, otherwise it is inferred from the maximum index in each mode.

Here is a simple toy example with factorizing a 3-tensor with side information on the first mode.

```
[ ]: import numpy as np
import pandas as pd
import scipy.sparse
import smurff
import itertools

## generating toy data
A = np.random.randn(15, 2)
B = np.random.randn(3, 2)
```

(continues on next page)

(continued from previous page)

```

C = np.random.randn(2, 2)

idx = list( itertools.product(np.arange(A.shape[0]),
                             np.arange(B.shape[0]),
                             np.arange(C.shape[0])) )
df = pd.DataFrame( np.asarray(idx), columns=["A", "B", "C"])
df["value"] = np.array([ np.sum(A[i[0], :] * B[i[1], :] * C[i[2], :]) for i in idx ])

## assigning 20% of the cells to test set
Ytrain, Ytest = smurff.make_train_test(df, 0.2)

print("Ytrain = ", Ytrain)

## for artificial dataset using small values for burnin, nsamples and num_latents is_
↪ fine
predictions = smurff.BPMFSession(
    Ytrain=Ytrain,
    Ytest=Ytest,
    num_latent=4,
    burnin=20,
    nsamples=20).run()

print("First prediction of Ytest tensor: ", predictions[0])

```

[]:

Note: This page was generated from [notebooks/different_methods.ipynb](#).

1.2.3 Trying different Matrix Factorization Methods

In this notebook we will try out several MF methods supported by SMURFF.

Downloading the files

As in the previous example we download the ChemBL dataset. The resulting IC50 matrix is a compound x protein matrix, split into train and test. The ECFP matrix has features as side information on the compounds.

```

[ ]: import smurff
import logging

logging.basicConfig(level = logging.INFO)

ic50_train, ic50_test, ecfp = smurff.load_chembl()

```

Matrix Factorization without Side Information (BPMF)

As a first example we can run SMURFF without side information. The method used here is BPMF.

Input matrix for Y is a sparse scipy matrix (either `coo_matrix`, `csr_matrix` or `csc_matrix`). The test matrix Y_{test} also needs to be sparse matrix of the same size as Y . Here we have used burn-in of 20 samples for the Gibbs sampler and then collected 80 samples from the model. We use 16 latent dimensions in the model.

For good results you will need to run more sampling and burnin iterations (≥ 1000) and maybe more latent dimensions.

We create a `trainSession`, and the `run` method returns the predictions of the `Ytest` matrix. `predictions` is a list of of type `Prediction`.

```
[ ]: trainSession = smurff.BPMFSession(
        Ytrain      = ic50_train,
        Ytest       = ic50_test,
        num_latent  = 16,
        burnin      = 20,
        nsamples    = 80,
        verbose     = 0,)

predictions = trainSession.run()
print("First prediction element: ", predictions[0])

rmse = smurff.calc_rmse(predictions)
print("RMSE =", rmse)
```

Matrix Factorization with Side Information (Macau)

If we want to use the compound features we can use the Macau algorithm.

The parameter `side_info = [ecfp, None]` sets the side information for rows and columns, respectively. In this example we only use side information for the compounds (rows of the matrix).

Since the `ecfp` sideinfo is sparse and large, we use the CG solver from Macau to reduce the memory footprint and speedup the computation.

```
[ ]: predictions = smurff.MacauSession(
        Ytrain      = ic50_train,
        Ytest       = ic50_test,
        side_info    = [ecfp, None],
        direct       = False, # use CG solver instead of Cholesky_
        ↪decomposition
        num_latent  = 16,
        burnin      = 40,
        nsamples    = 100).run()

smurff.calc_rmse(predictions)
```

Macau univariate sampler

SMURFF also includes an option to use a *very fast* univariate sampler, i.e., instead of sampling blocks of variables jointly it samples each individually. An example:

```
[ ]: predictions = smurff.MacauSession(
        Ytrain      = ic50_train,
        Ytest       = ic50_test,
        side_info    = [ecfp, None],
        direct       = True,
        univariate   = True,
        num_latent  = 32,
```

(continues on next page)

(continued from previous page)

```
        burnin      = 500,  
        nsamples    = 3500,  
        verbose     = 0,) .run()  
smurff.calc_rmse(predictions)
```

When using it we recommend using larger values for `burnin` and `nsamples`, because the univariate sampler mixes slower than the blocked sampler.

```
[ ]:
```

Note: This page was generated from [notebooks/different_noise_models.ipynb](#).

1.2.4 Different noise models

In this notebook we look at the different noise models.

Prepare train, test and side-info

We first need to download and prepare the data files. This can be accomplished using this a built-in function is `smurff.IC50` is a compound x protein matrix, The ECFP matrix as features as side information on the compounds.

```
[ ]: import smurff  
import logging  
  
logging.basicConfig(level = logging.INFO)  
  
ic50_train, ic50_test, ecfp = smurff.load_chembl()
```

Fixed noise

The noise model of observed data can be annotated by calling `addTrainAndTest` with the optional parameter `noise_model`. The default for this parameter is `FixedNoise` with precision 5.0

```
[ ]: trainSession = smurff.TrainSession(  
        priors = ['normal', 'normal'],  
        num_latent=32,  
        burnin=100,  
        nsamples=500)  
  
# the following line is equivalent to the default, not specifying noise_model  
trainSession.addTrainAndTest(ic50_train, ic50_test, smurff.FixedNoise(5.0))  
predictions = trainSession.run()  
print("RMSE = %.2f" % smurff.calc_rmse(predictions))
```

Adaptive noise

Instead of a fixed precision, we can also allow the model to automatically determine the precision of the noise by using `AdaptiveNoise`, with signal-to-noise ratio parameters `sn_init` and `sn_max`.

- `sn_init` is an initial signal-to-noise ratio.

- `sn_max` is the maximum allowed signal-to-noise ratio. This means that if the updated precision would imply a higher signal-to-noise ratio than `sn_max`, then the precision value is set to $(sn_max + 1.0) / Yvar$ where $Yvar$ is the variance of the training dataset Y .

```
[ ]: trainSession = smurff.TrainSession(
    priors = ['normal', 'normal'],
    num_latent=32,
    burnin=100,
    nsamples=500)

trainSession.addTrainAndTest(ic50_train, ic50_test, smurff.AdaptiveNoise(1.0, 10.))
predictions = trainSession.run()
print("RMSE = %.2f" % smurff.calc_rmse(predictions))
```

Binary matrices

SMURFF can also factorize binary matrices (with or without side information). The input matrices can contain arbitrary values, and are converted to 0's and 1' by means of a threshold. To factorize them we employ probit noise model `ProbitNoise`, taking this threshold as a parameter.

To evaluate binary factorization, we recommend to use *ROC AUC*, which can be enabled by providing a threshold also to the `TrainSession`.

```
[ ]: ic50_threshold = 6.
trainSession = smurff.TrainSession(
    priors = ['normal', 'normal'],
    num_latent=32,
    burnin=100,
    nsamples=100,
    # Using threshold of 6. to calculate AUC on test data
    threshold=ic50_threshold)

## using activity threshold pIC50 > 6. to binarize train data
trainSession.addTrainAndTest(ic50_train, ic50_test, smurff.ProbitNoise(ic50_
    ↪threshold))
predictions = trainSession.run()
print("RMSE = %.2f" % smurff.calc_rmse(predictions))
print("AUC = %.2f" % smurff.calc_auc(predictions, ic50_threshold))
```

The input train and test sets are converted to -1 and +1 values, if the original values are below or above the threshold (respectively). Similarly, the resulting predictions will be negative, if the model predicts the value to be below the threshold, or positive, if the model predicts the value to be above the threshold.

```
[ ]: predictions
```

Binary matrices with Side Info

It is possible to enhance the model for binary matrices by adding side information using the Macau algorithm. Note that the *binary* here refers to the train and test data, not to the side information.

```
[ ]: ic50_threshold = 6.
trainSession = smurff.TrainSession(
    priors = ['macau', 'normal'],
    num_latent=32,
```

(continues on next page)

(continued from previous page)

```

        burnin=100,
        nsamples=100,
        # Using threshold of 6. to calculate AUC on test data
        threshold=ic50_threshold)

## using activity threshold pIC50 > 6. to binarize train data
trainSession.addTrainAndTest(ic50_train, ic50_test, smurff.ProbitNoise(ic50_
    threshold))
trainSession.addSideInfo(0, ecfp, direct = True)
predictions = trainSession.run()
print("RMSE = %.2f" % smurff.calc_rmse(predictions))
print("AUC = %.2f" % smurff.calc_auc(predictions, ic50_threshold))

```

[]:

Note: This page was generated from [notebooks/inference_with_smurff.ipynb](#).

1.2.5 Inference with SMURFF

In this notebook we will continue on the first example. After running a training `trainSession` again in SMURFF, we will look deeper into how to use SMURFF for making predictions. The full Python API for predictions is available in [Python API Reference » Inference](#).

To make predictions we recall that the value of a tensor model is given by a tensor contraction of all latent matrices. Specifically, the prediction for the element \hat{Y}_{ijk} of a rank-3 tensor is given by

$$\hat{Y}_{ijk} = \sum_{d=1}^D u_{d,i}^{(1)} u_{d,j}^{(2)} u_{d,k}^{(3)} + mean$$

Since a matrix is a rank-2 tensor the prediction for a matrix is given by:

$$\hat{Y}_{ij} = \sum_{d=1}^D u_{d,i}^{(1)} u_{d,j}^{(2)} + mean$$

These inner products are computed by SMURFF automagically, as we will see below.

Saving models

We run a `Macau` training `trainSession` using side information (`ecfp`) from the chembl dataset. We make sure we *save every 10th sample*, such that we can load the model afterwards. This run will take some minutes to run.

```

[ ]: import smurff
import os
import logging

ic50_train, ic50_test, ecfp = smurff.load_chembl()

# limit to 100 rows and 100 features to make things go faster
ic50_train = ic50_train.tocsr()[ :100, :]
ic50_test = ic50_test.tocsr()[ :100, :]

```

(continues on next page)

(continued from previous page)

```
ecfp = ecfp.tocsr()[0:100,:].tocsc()[0:,:100]

trainSession = smurff.MacauSession(
    Ytrain      = ic50_train,
    Ytest       = ic50_test,
    side_info    = [ecfp, None],
    num_latent   = 16,
    burnin       = 200,
    nsamples     = 100,
    save_freq    = 10,
    save_name    = "ic50-macau.hdf5",
    verbose      = 0,)

predictions = trainSession.run()
```

Saved Model

The model is saved in an HDF5 file, in this case `ic50-macau.hdf5`. The file contains all saved info from this training run. For example:

```
[ ]: %%bash

h5ls -r ic50-macau.hdf5 | head -n 30
```

The structure of the HDF5 file is:

- Datasets in `/config` contain the input data and configuration provided to the `TrainSession`
- The different `/sample_*` datasets contain for each posterior sample:
 - `predictions/pred_1sample`: Predictions from this sample
 - `predictions/pred_avg`: Predictions average across this and all previous samples
 - `predictions/pred_var`: Predictions variance across this and all previous samples
- `latents_*`: Latent samples for each dimension
- `link_matrices/`: When `sideinfo` is used with the `MacauPrior`, this HDF5 group contains the β link matrix, and the μ HyperPrior sample. This will allow to make predictions from unseen `sideinfo`.

Sparse matrices and tensors are stored using the [h5sparse-tensor Python package](#) which is automatically installed as a dependency of `smurff`.

Making predictions from a TrainSession

The easiest way to make predictions is from an existing `TrainSession`:

```
[ ]: predictor = trainSession.makePredictSession()
print(predictor)
```

Once we have a `PredictSession`, there are several ways to make predictions:

- From a sparse matrix
- For all possible elements in the matrix (the complete $U \times V$)

- For a single point in the matrix
- Using only side-information

Predict all elements

We can make predictions for all rows \times columns in our matrix

```
[ ]: p = predictor.predict_all()
print(len(p)) # p is a list of numpy arrays of size num samples,
print(p[0].shape) # each array's shape = (num rows) x (num columns)
```

Predict element in a sparse matrix

We can make predictions for a sparse matrix, for example our `ic50_test` matrix:

```
[ ]: p = predictor.predict_sparse(ic50_test)
# p is a list of samples
print(len(p), " samples")
# each sample contains a sparse matrix with predictions for this sample
print("predictions 1:", type(p[0]), p[0])
```

Predict just one element

Or just one element. Let's predict the first element of our `ic50_test` matrix:

```
[ ]: from scipy.sparse import find
import numpy as np

(i,j,v) = find(ic50_test)
predictions = predictor.predict((i[0],j[0]))
print(predictions) # list of N samples, each sample is an array each containing a_
↳ single prediction
predictions = [ p.item() for p in predictions]
print("as a list of floats:", predictions)
```

And plot the histogram of predictions for this element.

```
[ ]: %matplotlib inline
import matplotlib.pyplot as plt

# Plot a histogram of the samples.
plt.subplot(111)
plt.hist(predictions, bins=10, density=True, label = "predictions's histogram")
plt.plot(v[0], 1., 'ro', markersize=5, label = 'actual value')
plt.legend()
plt.title('Histogram of ' + str(len(predictions)) + ' predictions')
plt.show()
```

Make predictions using side information

We can make predictions for rows/columns not in our train matrix, using only side info:

```
[ ]: import numpy as np
      from scipy.sparse import find

      (i,j,v) = find(ic50_test)
      row_side_info = ecfp.tocsr().getrow(i[0])
      p = predictor.predict((row_side_info,j[0]))
      print(p)
```

It is also possible to provide sideinfo for the columns, if the MacauPrior was used for the columns. See `smurff.PredictSession.predict` for the full documentation.

Making predictions from saved run

One can also make a `PredictSession` from a saved HDF5 file:

```
[ ]: import smurff

      predictor = smurff.PredictSession("ic50-macau.hdf5")
      print(predictor)
```

```
[ ]:
```

Note: This page was generated from [notebooks/centering.ipynb](#).

1.2.6 Centering

In this notebook we look at why it is important to center your matrices before using them with SMURFF.

SMURFF provides a function similar to `sklearn.preprocessing.scale` with the difference that SMURFF also supports scaling of sparse matrices. This makes sense only when the matrix is scarce, i.e. when the zero-elements represent unknown values.

The Matrix Factorization methods are meant for modeling the variance in the data. Hence it makes sense to subtract the mean first.

```
[ ]: import logging
      logging.basicConfig(level = logging.INFO)

      import smurff

      ic50_train, ic50_test, ecfp = smurff.load_chembl()

      ic50_train_centered, global_mean, _ = smurff.center_and_scale(ic50_train, "global",
      ↪with_mean = True, with_std = False)

      ic50_test_centered = ic50_test
      ic50_test_centered.data -= global_mean # only touch non-zeros
```

When we now run a SMURFF train `trainSession`, we can see from the `TrainSession` information, that the data has been centered:

```
PythonSession {  
  Data: {  
    Type: ScarceMatrixData [with NAs]  
    Component-wise mean: 3.86555e-16  
    ...  
  }
```

```
[ ]: trainSession = smurff.BPMFSession(  
    Ytrain      = ic50_train_centered,  
    Ytest       = ic50_test,  
    num_latent  = 16,  
    burnin     = 40,  
    nsamples    = 200,  
    verbose     = 0,  
)  
  
predictions = trainSession.run()  
rmse = smurff.calc_rmse(predictions)  
print(rmse)
```

1.3 Compilation of SMURFF

Note: the easiest way to install SMURFF is not to build it yourself. Install the binary [Conda](#) package:

```
conda install -c vanderaa smurff
```

1.3.1 Compilation using *conda build*

Conda build works on Linux, macOS and Windows. Execute

```
conda build smurff
```

in the *conda-recipes* directory.

1.3.2 Compilation using CMake

C++ Requirements

- CMake 3.6 or later
- Eigen3 version 3.3.7 or later
- HighFive 2.2. from <https://github.com/BlueBrain/HighFive/>
- Boost 1.5x or newer

Python Requirements

As in setup.py:

```
install_requires = [ 'numpy', 'scipy', 'pandas', 'scikit-learn', 'h5sparse-tensor' ], setup_requires = [ 'se-  
tuptools_scm', 'pybind11' ],
```

Compile using setup.py

Running `setup.py install`

will run CMake to configure, compile and install SMURFF. Extra arguments to CMake can be passed with

`setup.py --extra-cmake-args <...> install`

or by setting the `CMAKE_ARGS` environment variables.

CMake Options

- **Build type switches:**
 - `CMAKE_BUILD_TYPE` - Debug/Release
- **Algebra library switches (select only one):**
 - When no switches are specified, CMake will try to find any LAPACK and BLAS library on your system.
 - `ENABLE_OPENBLAS` - ON/OFF (should include openblas library when linking. openblas also contains implementation of lapack called relapack)
 - `ENABLE_MKL` - ON/OFF: tries to find the [MKL single dynamic library](#).
- **Python:**
 - `ENABLE_PYTHON`

Linux and macOS Specific

Have a look in [ci/](#) for Docker build scripts and for Linux+macOS wheel scripts. These scripts should give you a good idea on how to compile on an Ubuntu and macOS system.

Windows Specific

Work for a vcpkg-based build is in progress.

1.4 Python API Reference

This section of the documentation details modules, classes, and functions of SMURFF.

1.4.1 Preprocessing and Input Files

This section talks about:

Contents

- *Preprocessing and Input Files*
 - *Sparse Tensors*
 - *Split Train and Test*

- *Scaling and Centering*
- *Example ChEMBL dataset*

Sparse Tensors

class `smurff.SparseTensor` (*data, shape=None*)

Wrapper around a pandas DataFrame to represent a sparse tensor

The DataFrame should have N index columns (int type) and 1 value column (float type) N is the dimensionality of the tensor

You can also specify the shape of the tensor. If you don't it is detected automatically.

Split Train and Test

`smurff.make_train_test` (*Y, ntest, shape=None, seed=None*)

Splits a sparse matrix Y into a train and a test matrix.

Parameters **Y** – `numpy.ndarray` or `pandas.DataFrame` or `smurff.SparseTensor`

Matrix/Array/Tensor to split

Returns

- **Ytrain** (*csr_matrix*) – train part
- **Ytest** (*csr_matrix*) – test part

Scaling and Centering

`smurff.center.center_and_scale` (*m, mode, with_mean=True, with_std=True*)

Center and/or scale the matrix m to the mean and/or standard deviation.

Parameters

- **m** (*{array-like, sparse matrix}*) – The data to center and scale.
- **mode** (*{ "rows", "cols", "global" }*) –
 - “rows”: center/scale each row indepently
 - “cols”: center/scale each column idependently
 - “global”: center/scale using global meand and/or standard deviation/
- **with_mean** (*boolean, True by default*) – If True, center the data before scaling.
- **with_std** (*boolean, True by default*) – If True, scale the data to unit variance (or equivalently, unit standard deviation).

Returns

- **m** (*array-like*) – Transformed array.
- **mean** (*array-like or double or None*) – Computed mean depending on mode
- **std** (*array-like or double or None*) – Computed standard deviation depending on mode

Notes

Also supports scaling of sparse matrices. This makes sense only when the matrix is scarce, i.e. when the zero-elements represent unknown values.

Example ChEMBL dataset

```
smurff.load_chembl()
```

Downloads a small subset of the ChEMBL dataset.

Returns

- **ic50_train** (*sparse matrix*) – sparse train matrix
- **ic50_test** (*sparse matrix*) – sparse test matrix
- **feat** (*sparse matrix*) – sparse row features

1.4.2 Training

Contents

- *Training*
 - *TrainSession*
 - *MacauSession*
 - *BPMFSession*
 - *StatusItem*

The most versatile class is `TrainSession`. `MacauSession` and `BPMFSession` provide a simpler interface.

TrainSession

```
class smurff.TrainSession(priors=['normal', 'normal'], num_latent=None, num_threads=None,
                           burnin=None, nsamples=None, seed=None, threshold=None,
                           verbose=None, save_name=None, save_freq=None, check-
                           point_freq=None)
```

Class for doing a training run in smurff

A simple use case could be:

```
>>> trainSession = smurff.TrainSession(burnin = 5, nsamples = 5)
>>> trainSession.setTrain(Ydense)
>>> trainSession.run()
```

priors

The type of prior to use for each dimension

Type list, where element is one of { “normal”, “normalone”, “macau”, “macauone”, “spikeand-slab” }

num_latent

Number of latent dimensions in the model

Type int

burnin
Number of burnin samples to discard

Type int

nsamples
Number of samples to keep

Type int

num_threads
Number of OpenMP threads to use for model building

Type int

verbose
Verbosity level for C++ library

Type {0, 1, 2}

seed
Random seed to use for sampling

Type float

save_name
HDF5 filename to store the samples.

Type path

save_freq

- N>0: save every Nth sample
- N==0: never save a sample
- N== -1: save only the last sample

Type int

checkpoint_freq
Save the state of the trainSession every N seconds.

Type int

addData (*pos*, *Y*, *noise*=<smurff.helper.FixedNoise object>, *is_scarce*=False)
Stacks more matrices/tensors next to the main train matrix.

pos [shape] Block position of the data with respect to train. The train matrix/tensor has implicit block position (0, 0).

Y [:class: *numpy.ndarray*, *scipy.sparse* matrix or :class: *SparseTensor*] Data matrix/tensor to add

is_scarce [bool] When *Y* is sparse, and *is_scarce* is *True* the missing values are considered as *unknown*. When *Y* is sparse, and *is_scarce* is *False* the missing values are considered as *zero*. When *Y* is dense, this parameter is ignored.

noise [:class: *NoiseConfig*] Noise model to use for *Y*

addPropagatedPosterior (*mode*, *mu*, *Lambda*)
Adds *mu* and *Lambda* from propagated posterior

mode [int] dimension to add side info (rows = 0, cols = 1)

mu [:class: *numpy.ndarray* matrix] mean matrix mu should have as many rows as *num_latent* mu should have as many columns as size of dimension *mode* in *train*

Lambda [:class: *numpy.ndarray* matrix] co-variance matrix Lambda should be shaped like $K \times K \times N$ Where $K == num_latent$ and $N == dimension\ mode\ in\ train$

addSideInfo (*mode*, *Y*, *noise*=<*smurff.helper.SampledNoise* object>, *direct*=True)

Adds fully known side info, for use in with the macau or macauone prior

mode [int] dimension to add side info (rows = 0, cols = 1)

Y [:class: *numpy.ndarray*, *scipy.sparse* matrix] Side info matrix/tensor Y should have as many rows in Y as you have elemnts in the dimension selected using *mode*. Columns in Y are features for each element.

noise [:class: *NoiseConfig*] Noise model to use for *Y*

direct [boolean]

- When True, uses a direct inversion method.
- When False, uses a CG solver

The direct method is only feasible for a small (< 100K) number of features.

init ()

Initializes the *TrainSession* after all data has been added.

You need to call this method befor calling *step* (), unless you call *run* ()

Returns

Return type *StatusItem* of the trainSession.

makePredictSession ()

Makes a *PredictSession* based on the model that as built in this *TrainSession*.

run ()

Equivalent to:

```
self.init()
while self.step():
    pass
```

setTrain (*Y*, *noise*=<*smurff.helper.FixedNoise* object>, *is_scarce*=True)

Adds a train and optionally a test matrix as input data to this *TrainSession*

Parameters

- **Y** – Train matrix/tensor
- **noise** – Noise model to use for *Y*
- **is_scarce** (*bool*) – When *Y* is sparse, and *is_scarce* is *True* the missing values are considered as *unknown*. When *Y* is sparse, and *is_scarce* is *False* the missing values are considered as *zero*. When *Y* is dense, this parameter is ignored.

step ()

Does on sampling or burnin iteration.

Returns

- - **When a step was executed** (*StatusItem* of the trainSession.)
- - **After the last iteration, when no step was executed** (*None*.)

MacauSession

```
class smurff.MacauSession (Ytrain, is_scarce=True, Ytest=None, side_info=None, univariate=False, direct=True, *args, **kwargs)
```

A train trainSession specialized for use with the Macau algorithm

Ytrain

Train matrix/tensor

Ytest [scipy.sparse matrix or :class: *SparseTensor*] Test matrix/tensor. Mainly used for calculating RMSE.

side_info [list of :class: *numpy.ndarray*, scipy.sparse matrix or None] Side info matrix/tensor for each dimension If there is no side info for a certain mode, pass *None*. Each side info should have as many rows as you have elemnts in corresponding dimension of *Ytrain*.

direct [bool] Use Cholesky instead of CG solver

univariate [bool] Use univariate or multivariate sampling.

****args:** Extra arguments are passed to the *TrainSession*

Type

class *numpy.ndarray*, scipy.sparse matrix or :class: *SparseTensor*

BPMFSession

```
class smurff.BPMFSession (Ytrain, is_scarce=True, Ytest=None, univariate=False, *args, **kwargs)
```

A train trainSession specialized for use with the BPMF algorithm

Ytrain

Train matrix/tensor

Ytest [scipy.sparse matrix or :class: *SparseTensor*] Test matrix/tensor. Mainly used for calculating RMSE.

univariate [bool] Use univariate or multivariate sampling.

****args:** Extra arguments are passed to the *TrainSession*

Type

class *numpy.ndarray*, scipy.sparse matrix or :class: *SparseTensor*

StatusItem

```
class smurff.StatusItem
```

Short set of parameters indicative for the training progress.

auc_1sample

ROC AUC of the test matrix of the last sampleOnly available if you provided a threshold

auc_avg

Average ROC AUC of the test matrix across all samplesOnly available if you provided a threshold

elapsed_iter

Number of seconds the last sampling iteration took

iter
Current iteration in current phase

nnz_per_sec
Compute performance indicator; number of non-zero elements in train processed per second

phase
{ "Burnin", "Sampling" }

rmse_1sample
RMSE for test matrix of last sample

rmse_avg
Averag RMSE for test matrix across all samples

samples_per_sec
Compute performance indicator; number of rows and columns in U/V processed per second

train_rmse
RMSE for train matrix of last sample

1.4.3 Inference

Contents

- *Inference*
 - *PredictSession*
 - *Prediction*

A `PredictSession` provides access to making prediction from saved models. Predictions for a single point are stored in the `Prediction`.

PredictSession

class `smurff.PredictSession(h5_fname)`

TrainSession for making predictions using a model generated using a *TrainSession*.

A *PredictSession* can be made directly from a *TrainSession*

```
>>> predict_session = train_session.makePredictSession()
```

or from an HDF5 file

```
>>> predict_session = PredictSession("saved_output.hdf5")
```

predict (*operands*, *samples=None*)

Generate predictions on *operands*. Parameters specify coordinates of sideinfo/features for each dimension.

Parameters

- **operands** (*tuple*) – A combination of coordindates in the matrix/tensor and/or features you want to use to make predictions. *len(operands)* should be equal to number of dimensions in the sample. Each element *operands* can be a:
 - `int` : a single element in this dimension is selected. For example, a single row or column in a matrix.

- `slice`: a slice is selected in this dimension. For example, a number of rows or columns in a matrix.
- Ellipsis or `None`: all elements in this dimension are selected. For example, all rows or columns in a matrix.
- `numpy.ndarray`: 2D numpy array used as dense sideinfo. Each row vector is used as side-info.
- `scipy.sparse.spmatrix`: sparse matrix used as sideinfo. Each row vector is used as side-info.
- **samples** (*range or None*) – Range of samples to use for prediction, or None for all samples

Returns list of N `numpy.ndarray`'s of shape `[T1 x T2 x ...]` where N is the number of samples in this *PredictSession* and *T1 x T2 x ...*

Return type list of `numpy.ndarray`

predict_all()

Computes prediction matrix/tensor for full train data shape.

Returns N `numpy.ndarray`'s of shape `[T1 x T2 x ...]` where N is the number of samples in this *PredictSession* and *T1 x T2 x ...* is the shape of the train data

Return type list of `numpy.ndarray`

predict_sparse(test_matrix)

Computes prediction for all elements in a sparse test matrix

Parameters **test_matrix** (*scipy sparse matrix*) –

Returns list of N `scipy.sparse.sp_matrix` objects, where N is the number of samples

Return type list

Prediction

class `smurff.Prediction` (*coords, val, pred_1sample=nan, pred_avg=nan, var=nan, nsamples=-1*)

Stores predictions for a single point in the matrix/tensor

coords

Position of this prediction in the train matrix/tensor

Type shape

val

True value or “nan” if no true value is known

Type float

nsamples

Number of samples collected to make this prediction

Type int

pred_1sample

Predicted value using only the last sample

Type float

pred_avg

Predicted value using the average prediction across all samples

Type float

var

Variance amongst predictions across all samples

Type float

pred_all

List of predictions, one for each sample

Type list

static fromTestMatrix (*test_matrix_or_tensor*)

Creates a list of predictions from a scipy sparse matrix”

Parameters **test_matrix** (*scipy sparse matrix*)-

Returns List of *Prediction*. Only the coordinate and true value is filled.

Return type list

TODO: FixedNoise, AdaptiveNoise, ProbitNoise

S

`smurff.center`, [18](#)

A

addData() (*smurff.TrainSession* method), 20
 addPropagatedPosterior() (*smurff.TrainSession* method), 20
 addSideInfo() (*smurff.TrainSession* method), 21
 auc_1sample (*smurff.StatusItem* attribute), 22
 auc_avg (*smurff.StatusItem* attribute), 22

B

BPMFSession (class in *smurff*), 22
 burnin (*smurff.TrainSession* attribute), 20

C

center_and_scale() (in module *smurff.center*), 18
 checkpoint_freq (*smurff.TrainSession* attribute), 20
 coords (*smurff.Prediction* attribute), 24

E

elapsed_iter (*smurff.StatusItem* attribute), 22

F

fromTestMatrix() (*smurff.Prediction* static method), 25

I

init() (*smurff.TrainSession* method), 21
 iter (*smurff.StatusItem* attribute), 22

L

load_chembl() (in module *smurff*), 19

M

MacauSession (class in *smurff*), 22
 make_train_test() (in module *smurff*), 18
 makePredictSession() (*smurff.TrainSession* method), 21

N

nnz_per_sec (*smurff.StatusItem* attribute), 23

nsamples (*smurff.Prediction* attribute), 24
 nsamples (*smurff.TrainSession* attribute), 20
 num_latent (*smurff.TrainSession* attribute), 19
 num_threads (*smurff.TrainSession* attribute), 20

P

phase (*smurff.StatusItem* attribute), 23
 pred_1sample (*smurff.Prediction* attribute), 24
 pred_all (*smurff.Prediction* attribute), 25
 pred_avg (*smurff.Prediction* attribute), 24
 predict() (*smurff.PredictSession* method), 23
 predict_all() (*smurff.PredictSession* method), 24
 predict_sparse() (*smurff.PredictSession* method), 24

Prediction (class in *smurff*), 24
 PredictSession (class in *smurff*), 23
 priors (*smurff.TrainSession* attribute), 19

R

rmse_1sample (*smurff.StatusItem* attribute), 23
 rmse_avg (*smurff.StatusItem* attribute), 23
 run() (*smurff.TrainSession* method), 21

S

samples_per_sec (*smurff.StatusItem* attribute), 23
 save_freq (*smurff.TrainSession* attribute), 20
 save_name (*smurff.TrainSession* attribute), 20
 seed (*smurff.TrainSession* attribute), 20
 setTrain() (*smurff.TrainSession* method), 21
 smurff.center (module), 18
 SparseTensor (class in *smurff*), 18
 StatusItem (class in *smurff*), 22
 step() (*smurff.TrainSession* method), 21

T

train_rmse (*smurff.StatusItem* attribute), 23
 TrainSession (class in *smurff*), 19

V

val (*smurff.Prediction* attribute), 24

`var` (*smurff.Prediction attribute*), [25](#)

`verbose` (*smurff.TrainSession attribute*), [20](#)

Y

`Ytrain` (*smurff.BPMFSession attribute*), [22](#)

`Ytrain` (*smurff.MacauSession attribute*), [22](#)