

# Programming assignment 10: Matrix Factorization

In [1]:

```
import time
import scipy.sparse as sp
import numpy as np
from scipy.sparse.linalg import svds
from sklearn.linear_model import Ridge

import matplotlib.pyplot as plt
%matplotlib inline
```



## Exporting the results to PDF

Once you complete the assignments, export the entire notebook as PDF and attach it to your homework solutions. The best way of doing that is

1. Run all the cells of the notebook.
2. Download the notebook in HTML (click File > Download as > .html)
3. Convert the HTML to PDF using e.g. <https://www.sejda.com/html-to-pdf> (<https://www.sejda.com/html-to-pdf>) or wkhtmltopdf for Linux ([tutorial \(https://www.cyberciti.biz/open-source/html-to-pdf-freeware-linux-osx-windows-software/\)](https://www.cyberciti.biz/open-source/html-to-pdf-freeware-linux-osx-windows-software/))
4. Concatenate your solutions for other tasks with the output of Step 3. On a Linux machine you can simply use `pdffunite`, there are similar tools for other platforms too. You can only upload a single PDF file to Moodle.

This way is preferred to using `nbconvert`, since `nbconvert` clips lines that exceed page width and makes your code harder to grade.

## Restaurant recommendation

The goal of this task is to recommend restaurants to users based on the rating data in the Yelp dataset. For this, we try to predict the rating a user will give to a restaurant they have not yet rated based on a latent factor model.

Specifically, the objective function (loss) we wanted to optimize is:

$$\mathcal{L} = \min_{P,Q} \sum_{(i,x) \in W} (M_{ix} - \mathbf{q}_i^T \mathbf{p}_x)^2 + \lambda \sum_x \|\mathbf{p}_x\|^2 + \lambda \sum_i \|\mathbf{q}_i\|^2$$

where  $W$  is the set of  $(i, x)$  pairs for which the rating  $M_{ix}$  given by user  $i$  to restaurant  $x$  is known. Here we have also introduced two regularization terms to help us with overfitting where  $\lambda$  is hyper-parameter that control the strength of the regularization.

**Hint 1:** Using the closed form solution for regression might lead to singular values. To avoid this issue perform the regression step with an existing package such as scikit-learn. It is advisable to use ridge regression to account for regularization.

**Hint 2:** If you are using the scikit-learn package remember to set `fit_intercept = False` to only learn the coefficients of the linear regression.

## Load and Preprocess the Data (nothing to do here)

In [2]:

```
ratings = np.load("ratings.npy")
```

In [3]:

```
# We have triplets of (user, restaurant, rating).
ratings
```

Out[3]:

```
array([[101968, 1880, 1],
       [101968, 284, 5],
       [101968, 1378, 2],
       ...,
       [ 72452, 2100, 4],
       [ 72452, 2050, 5],
       [ 74861, 3979, 5]])
```

Now we transform the data into a matrix of dimension  $[N, D]$ , where  $N$  is the number of users and  $D$  is the number of restaurants in the dataset. We store the data as a sparse matrix to avoid out-of-memory issues.

In [4]:

```
n_users = np.max(ratings[:,0] + 1)
n_restaurants = np.max(ratings[:,1] + 1)
M = sp.coo_matrix((ratings[:,2], (ratings[:,0], ratings[:,1])), shape=(n_users, n_restaurants)).tocsc()
M
```

Out[4]:

```
<337867x5899 sparse matrix of type '<class 'numpy.int64'>'
  with 929606 stored elements in Compressed Sparse Row format>
```

To avoid the [cold start problem](https://en.wikipedia.org/wiki/Cold_start_(computing)) ([https://en.wikipedia.org/wiki/Cold\\_start\\_\(computing\)](https://en.wikipedia.org/wiki/Cold_start_(computing))), in the preprocessing step, we recursively remove all users and restaurants with 10 or less ratings.

Then, we randomly select 200 data points for the validation and test sets, respectively.

After this, we subtract the mean rating for each users to account for this global effect.

**Note:** Some entries might become zero in this process -- but these entries are different than the 'unknown' zeros in the matrix. We store the indices for which we the rating data available in a separate variable.

In [5]:



```
def cold_start_preprocessing(matrix, min_entries):
    """
    Recursively removes rows and columns from the input matrix which have less than min_entries nonz

    Parameters
    -----
    matrix      : sp.spmatrix, shape [N, D]
                  The input matrix to be preprocessed.
    min_entries : int
                  Minimum number of nonzero elements per row and column.

    Returns
    -----
    matrix      : sp.spmatrix, shape [N', D']
                  The pre-processed matrix, where N' <= N and D' <= D

    """
    print("Shape before: {}".format(matrix.shape))

    shape = (-1, -1)
    while matrix.shape != shape:
        shape = matrix.shape
        nnz = matrix>0
        row_ixs = nnz.sum(1).A1 > min_entries
        matrix = matrix[row_ixs]
        nnz = matrix>0
        col_ixs = nnz.sum(0).A1 > min_entries
        matrix = matrix[:, col_ixs]
    print("Shape after: {}".format(matrix.shape))
    nnz = matrix>0
    assert (nnz.sum(0).A1 > min_entries).all()
    assert (nnz.sum(1).A1 > min_entries).all()
    return matrix
```

## Task 1: Implement a function that subtracts the mean user rating from the sparse rating matrix

In [6]:



```
def shift_user_mean(matrix):  
    """  
    Subtract the mean rating per user from the non-zero elements in the input matrix.  
  
    Parameters  
    -----  
    matrix : sp.spmatrix, shape [N, D]  
        Input sparse matrix.  
    Returns  
    -----  
    matrix : sp.spmatrix, shape [N, D]  
        The modified input matrix.  
  
    user_means : np.array, shape [N, 1]  
        The mean rating per user that can be used to recover the absolute ratings from the  
  
    """  
  
    # YOUR CODE HERE  
    user_means = np.mean(matrix, axis=1)  
    user_means = user_means.reshape(-1, 1) # reshape to [N, 1]  
  
    matrix = matrix - user_means  
    assert np.all(np.isclose(matrix.mean(1), 0))  
    return matrix, user_means
```

**Split the data into a train, validation and test set (nothing to do here)**

In [7]:



```
def split_data(matrix, n_validation, n_test):
    """
    Extract validation and test entries from the input matrix.

    Parameters
    -----
    matrix          : sp.spmatrix, shape [N, D]
                     The input data matrix.
    n_validation     : int
                     The number of validation entries to extract.
    n_test          : int
                     The number of test entries to extract.

    Returns
    -----
    matrix_split    : sp.spmatrix, shape [N, D]
                     A copy of the input matrix in which the validation and test entries have been

    val_idx         : tuple, shape [2, n_validation]
                     The indices of the validation entries.

    test_idx        : tuple, shape [2, n_test]
                     The indices of the test entries.

    val_values      : np.array, shape [n_validation, ]
                     The values of the input matrix at the validation indices.

    test_values     : np.array, shape [n_test, ]
                     The values of the input matrix at the test indices.

    """

    matrix_cp = matrix.copy()
    non_zero_idx = np.argwhere(matrix_cp)
    ix = np.random.permutation(non_zero_idx)
    val_idx = tuple(ixs[:n_validation].T)
    test_idx = tuple(ixs[n_validation:n_validation + n_test].T)

    val_values = matrix_cp[val_idx].A1
    test_values = matrix_cp[test_idx].A1

    matrix_cp[val_idx] = matrix_cp[test_idx] = 0
    matrix_cp.eliminate_zeros()

    return matrix_cp, val_idx, test_idx, val_values, test_values
```

In [8]:



```
M = cold_start_preprocessing(M, 20)
```

Shape before: (337867, 5899)

Shape after: (3529, 2072)

In [9]:



```
n_validation = 200
n_test = 200
# Split data
M_train, val_idx, test_idx, val_values, test_values = split_data(M, n_validation, n_test)
```

In [10]:



```
# Remove user means.
nonzero_indices = np.argwhere(M_train)
M_shifted, user_means = shift_user_mean(M_train)
# Apply the same shift to the validation and test data.
val_values_shifted = val_values - user_means[np.array(val_idx).T[:, 0]].A1
test_values_shifted = test_values - user_means[np.array(test_idx).T[:, 0]].A1
```

In [19]:



```
print(M_train.shape)
M_shifted.shape
```

(3529, 2072)

Out[19]:

(3529, 2072)

**Compute the loss function (nothing to do here)**

In [11]:



```
def loss(values, ix, Q, P, reg_lambda):
    """
    Compute the loss of the latent factor model (at indices ix).
    Parameters
    -----
    values : np.array, shape [n_ixs,]
        The array with the ground-truth values.
    ix : tuple, shape [2, n_ixs]
        The indices at which we want to evaluate the loss (usually the nonzero indices of the unshifted matrix)
    Q : np.array, shape [N, k]
        The matrix Q of a latent factor model.
    P : np.array, shape [k, D]
        The matrix P of a latent factor model.
    reg_lambda : float
        The regularization strength

    Returns
    -----
    loss : float
        The loss of the latent factor model.

    """
    mean_sse_loss = np.sum((values - Q.dot(P)[ix])**2)
    regularization_loss = reg_lambda * (np.sum(np.linalg.norm(Q, axis=0)**2) + np.sum(np.linalg.norm(P, axis=0)**2))
    return mean_sse_loss + regularization_loss
```

## Alternating optimization

In the first step, we will approach the problem via alternating optimization, as learned in the lecture. That is, during each iteration you first update  $Q$  while having  $P$  fixed and then vice versa.

### Task 2: Implement a function that initializes the latent factors $Q$ and $P$

In [12]:



```
def initialize_Q_P(matrix, k, init='random'):
    """
    Initialize the matrices Q and P for a latent factor model.

    Parameters
    -----
    matrix : sp.spmatrix, shape [N, D]
        The matrix to be factorized.
    k      : int
        The number of latent dimensions.
    init   : str in ['svd', 'random'], default: 'random'
        The initialization strategy. 'svd' means that we use SVD to initialize P and Q, 'random'
        the entries in P and Q randomly in the interval [0, 1).

    Returns
    -----
    Q : np.array, shape [N, k]
        The initialized matrix Q of a latent factor model.

    P : np.array, shape [k, D]
        The initialized matrix P of a latent factor model.
    """
    np.random.seed(0)
    # YOUR CODE HERE

    N = matrix.shape[0]
    D = matrix.shape[1]
    if init == 'random':
        Q = np.random.rand(N, k)/k
        P = np.random.rand(k, D)/k
    elif init == 'svd':
        U, S, V = np.linalg.svd(matrix, full_matrices=False)
        Q = U*S
        P = V

    assert Q.shape == (matrix.shape[0], k)
    assert P.shape == (k, matrix.shape[1])
    return Q, P
```

### Task 3: Implement the alternating optimization approach



In [15]:



```
def latent_factor_alternating_optimization(M, non_zero_idx, k, val_idx, val_values,
                                           reg_lambda, max_steps=100, init='random',
                                           log_every=1, patience=5, eval_every=1):
    """
    Perform matrix factorization using alternating optimization. Training is done via patience,
    i.e. we stop training after we observe no improvement on the validation loss for a certain
    amount of training steps. We then return the best values for Q and P observed during training.

    Parameters
    -----
    M
        : sp.spmatrix, shape [N, D]
        The input matrix to be factorized.

    non_zero_idx
        : np.array, shape [nnz, 2]
        The indices of the non-zero entries of the un-shifted matrix to be factorized.
        nnz refers to the number of non-zero entries. Note that this may be different
        from the number of non-zero entries in the input matrix M, e.g. in the case
        that all ratings by a user have the same value.

    k
        : int
        The latent factor dimension.

    val_idx
        : tuple, shape [2, n_validation]
        Tuple of the validation set indices.
        n_validation refers to the size of the validation set.

    val_values
        : np.array, shape [n_validation, ]
        The values in the validation set.

    reg_lambda
        : float
        The regularization strength.

    max_steps
        : int, optional, default: 100
        Maximum number of training steps. Note that we will stop early if we observe
        no improvement on the validation error for a specified number of steps
        (see "patience" for details).

    init
        : str in ['random', 'svd'], default 'random'
        The initialization strategy for P and Q. See function initialize_Q_P for details.

    log_every
        : int, optional, default: 1
        Log the training status every X iterations.

    patience
        : int, optional, default: 5
        Stop training after we observe no improvement of the validation loss for X
        iterations (see eval_every for details). After we stop training, we restore
        observed values for Q and P (based on the validation loss) and return them.

    eval_every
        : int, optional, default: 1
        Evaluate the training and validation loss every X steps. If we observe no improvement
        of the validation error, we decrease our patience by 1, else we reset it to
        the original value.

    Returns
    -----
    best_Q
        : np.array, shape [N, k]
        Best value for Q (based on validation loss) observed during training

    best_P
        : np.array, shape [k, D]
        Best value for P (based on validation loss) observed during training
```

```

validation_losses : list of floats
                    Validation loss for every evaluation iteration, can be used for plotting the
                    loss over time.

train_losses      : list of floats
                    Training loss for every evaluation iteration, can be used for plotting the t
                    loss over time.

converged_after   : int
                    it - patience*eval_every, where it is the iteration in which patience hits 0
                    or -1 if we hit max_steps before converging.

"""

# YOUR CODE HERE
# initialize the params
Q, P = initialize_Q_P(M, k, init)
best_Q = Q
best_P = P
lost_patience = 0
it = 0
validation_losses=[]
train_losses=[]
clf_p = Ridge(reg_lambda, fit_intercept=False, solver='auto')
clf_q = Ridge(reg_lambda, fit_intercept=False, solver='auto')

idx_x = non_zero_idx[:,0]
idx_y = non_zero_idx[:,1]
value = M[idx_x, idx_y]
value = np.asarray(value).reshape(-1)

idx_tuple = non_zero_idx.transpose()
idx_tuple = tuple(map(tuple, idx_tuple))
# print(value.shape)

# start training
while lost_patience<patience:

    it += 1
    current_train_loss = loss(value, idx_tuple, Q, P, reg_lambda)
    train_losses.append(current_train_loss)

    if it%eval_every==0:
        current_val_loss = loss(val_values, val_idx, Q, P, reg_lambda)
        validation_losses.append(current_val_loss)

    if len(train_losses)>=2 and train_losses[-1]>=train_losses[-2]:
        lost_patience += 1
    else:
        lost_patience = 0

    if len(train_losses)>=2 and train_losses[-1]<train_losses[-2]:
        best_Q = Q
        best_P = P

    # update Q
    clf_q.fit(P.transpose(), M.transpose())
    Q = clf_q.coef_

    # update P

```

```
clf_p.fit(Q, M)
P = clf_p.coef_.transpose()

converged_after = it - patience*eval_every
print("Iteration %s,   trainning loss %.2f,   validation loss %.2f" % (it,
                                                                    train_losses[-1],
                                                                    validation_losses[-1]))

return best_Q, best_P, validation_losses, train_losses, converged_after
```

## Train the latent factor (nothing to do here)

In [16]:



```
Q, P, val_loss, train_loss, converged = latent_factor_alternating_optimization(M_shifted, nonzero_in
                                     k=100, val_idx=val_idx,
                                     val_values=val_values,
                                     reg_lambda=1e-4, init
                                     max_steps=100, patien
```

Iteration 1,	trainning loss 2218582.81,	validation loss 3019.77
Iteration 2,	trainning loss 1577870.02,	validation loss 3532.65
Iteration 3,	trainning loss 1306264.02,	validation loss 3733.20
Iteration 4,	trainning loss 1266097.77,	validation loss 3839.26
Iteration 5,	trainning loss 1255187.13,	validation loss 3880.03
Iteration 6,	trainning loss 1251088.68,	validation loss 3893.93
Iteration 7,	trainning loss 1249164.49,	validation loss 3897.90
Iteration 8,	trainning loss 1248111.33,	validation loss 3898.21
Iteration 9,	trainning loss 1247467.32,	validation loss 3897.06
Iteration 10,	trainning loss 1247035.46,	validation loss 3895.21
Iteration 11,	trainning loss 1246722.09,	validation loss 3893.01
Iteration 12,	trainning loss 1246480.16,	validation loss 3890.59
Iteration 13,	trainning loss 1246284.90,	validation loss 3888.04
Iteration 14,	trainning loss 1246122.53,	validation loss 3885.41
Iteration 15,	trainning loss 1245984.76,	validation loss 3882.73
Iteration 16,	trainning loss 1245866.22,	validation loss 3880.01
Iteration 17,	trainning loss 1245763.15,	validation loss 3877.28
Iteration 18,	trainning loss 1245672.77,	validation loss 3874.53
Iteration 19,	trainning loss 1245592.94,	validation loss 3871.77
Iteration 20,	trainning loss 1245521.97,	validation loss 3869.01
Iteration 21,	trainning loss 1245458.48,	validation loss 3866.24
Iteration 22,	trainning loss 1245401.41,	validation loss 3863.48
Iteration 23,	trainning loss 1245349.86,	validation loss 3860.71
Iteration 24,	trainning loss 1245303.14,	validation loss 3857.94
Iteration 25,	trainning loss 1245260.66,	validation loss 3855.16
Iteration 26,	trainning loss 1245221.96,	validation loss 3852.38
Iteration 27,	trainning loss 1245186.64,	validation loss 3849.60
Iteration 28,	trainning loss 1245154.39,	validation loss 3846.81
Iteration 29,	trainning loss 1245124.92,	validation loss 3844.01
Iteration 30,	trainning loss 1245098.00,	validation loss 3841.21
Iteration 31,	trainning loss 1245073.40,	validation loss 3838.40
Iteration 32,	trainning loss 1245050.95,	validation loss 3835.59
Iteration 33,	trainning loss 1245030.47,	validation loss 3832.78
Iteration 34,	trainning loss 1245011.80,	validation loss 3829.97
Iteration 35,	trainning loss 1244994.81,	validation loss 3827.16
Iteration 36,	trainning loss 1244979.36,	validation loss 3824.34
Iteration 37,	trainning loss 1244965.33,	validation loss 3821.53
Iteration 38,	trainning loss 1244952.61,	validation loss 3818.73
Iteration 39,	trainning loss 1244941.10,	validation loss 3815.93
Iteration 40,	trainning loss 1244930.70,	validation loss 3813.13
Iteration 41,	trainning loss 1244921.33,	validation loss 3810.34
Iteration 42,	trainning loss 1244912.90,	validation loss 3807.56
Iteration 43,	trainning loss 1244905.34,	validation loss 3804.79
Iteration 44,	trainning loss 1244898.57,	validation loss 3802.02
Iteration 45,	trainning loss 1244892.55,	validation loss 3799.27
Iteration 46,	trainning loss 1244887.20,	validation loss 3796.52
Iteration 47,	trainning loss 1244882.47,	validation loss 3793.79
Iteration 48,	trainning loss 1244878.31,	validation loss 3791.06
Iteration 49,	trainning loss 1244874.68,	validation loss 3788.35
Iteration 50,	trainning loss 1244871.54,	validation loss 3785.65
Iteration 51,	trainning loss 1244868.83,	validation loss 3782.96

Iteration 52,	trainning loss	1244866.54,	validation loss	3780.28
Iteration 53,	trainning loss	1244864.61,	validation loss	3777.62
Iteration 54,	trainning loss	1244863.03,	validation loss	3774.97
Iteration 55,	trainning loss	1244861.77,	validation loss	3772.33
Iteration 56,	trainning loss	1244860.79,	validation loss	3769.70
Iteration 57,	trainning loss	1244860.08,	validation loss	3767.09
Iteration 58,	trainning loss	1244859.60,	validation loss	3764.49
Iteration 59,	trainning loss	1244859.35,	validation loss	3761.91
Iteration 60,	trainning loss	1244859.30,	validation loss	3759.33
Iteration 61,	trainning loss	1244859.43,	validation loss	3756.78
Iteration 62,	trainning loss	1244859.73,	validation loss	3754.23
Iteration 63,	trainning loss	1244860.18,	validation loss	3751.70
Iteration 64,	trainning loss	1244860.76,	validation loss	3749.19
Iteration 65,	trainning loss	1244861.48,	validation loss	3746.68
Iteration 66,	trainning loss	1244862.30,	validation loss	3744.20
Iteration 67,	trainning loss	1244863.23,	validation loss	3741.72
Iteration 68,	trainning loss	1244864.25,	validation loss	3739.26
Iteration 69,	trainning loss	1244865.35,	validation loss	3736.82
Iteration 70,	trainning loss	1244866.53,	validation loss	3734.38

In [1]:



```
# I am sure there is sth wrong in the model, but I can't find it.
```

**Plot the validation and training losses over for each iteration (nothing to do here)**

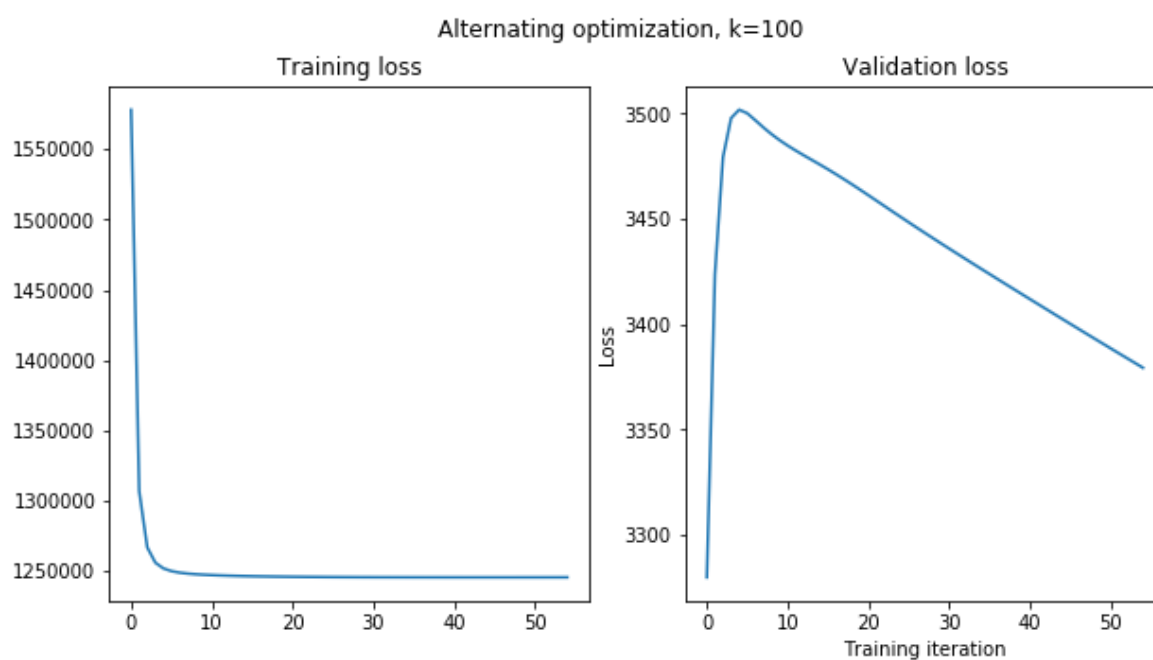
In [199]:

```
fig, ax = plt.subplots(1, 2, figsize=[10, 5])  
fig.suptitle("Alternating optimization, k=100")
```

```
ax[0].plot(train_loss[1::])  
ax[0].set_title('Training loss')  
plt.xlabel("Training iteration")  
plt.ylabel("Loss")
```

```
ax[1].plot(val_loss[1::])  
ax[1].set_title('Validation loss')  
plt.xlabel("Training iteration")  
plt.ylabel("Loss")
```

```
plt.show()
```



In [ ]: