

Image Compression & Recommendation Systems

Tool : Singular Value Decomposition

The Analytics Edge

Images are an important type of unstructured data on the Internet. Companies that share user photo such as Instagram, Facebook, Flickr, Twitter need efficient ways to deal with photos of users. Data Compression is an important tool that companies use to tradeoff quality of photos and file size (memory). Tools of matrix factorization such as singular value decomposition is a useful tool. Similarly such techniques are useful in recommendation systems where singular value decomposition on matrices with missing entries are used to make user recommendations.

Photos

Instagram is an online photo sharing site that enables users to take and share photos. It was created in 2010 as a free mobile app and acquired by Facebook in 2012. The company valuation is today between 30 to 37 billion. The total number of photos uploaded has been over 1 billion.

Just recently (May 2015), Google has launched Google Photos, a photo and video sharing and storage service in an attempt to improve on the popularity of Google+.

"A picture is worth a thousand words"

So just how many pictures (photos) are uploaded to the Internet?

It is estimated that 2400 Instagram photos are uploaded in 1 second (internetlivestat.com/one-second)

Each day, a total of 1.8 billion photos are estimated to be uploaded to the Internet

Assuming that an image (photo) on average is 100 Kb. This means that for just one day, the amount of bytes needed = 1.8×10^{14} bytes which is 180 Terabytes.

Image representation

Image \longrightarrow Matrix
Pixel representation

Gray scale image

A gray scale image is an image where the value of each pixel is a single number that carries the intensity information.

For example the images have intensity going from white (0) to black (1) with the different shades of gray taking values in (0, 1).

Color image

Color images can be built as a set of stacked color channels, each of them representing value levels of the given channel.

RGB images \longrightarrow Three channels
Red, Green, Blue
Three primary colors

Singular Value Decomposition (SVD)

Given a rectangular matrix X of dimension $m \times n$, a singular value decomposition of X is of the form: (say $m \geq n$)

$$\underbrace{X}_{\substack{m \times n \\ n}} = \underbrace{U}_{\substack{m \times m \\ m}} \underbrace{S}_{\substack{m \times n \\ n}} \underbrace{V^T}_{\substack{n \times n \\ n}}$$
$$m \begin{pmatrix} n \\ \end{pmatrix} = m \begin{pmatrix} m \\ \end{pmatrix} m \begin{pmatrix} \sigma_1 & \dots & \sigma_n \\ \hline 0 \end{pmatrix} n \begin{pmatrix} n \\ \end{pmatrix}$$

where:

- 1) Both matrices U and V are orthogonal matrices, namely $U^T U = U U^T = I$ and $V^T V = V V^T = I$ where I is an identity matrix
- 2) Matrix S has diagonal entries $\sigma_1, \dots, \sigma_r \geq 0$ where $r = \min(m, n)$ at the top and 0's filling the rest of the matrix. ($\sigma_1, \dots, \sigma_r$ are the singular values)

Equivalently matrix X can be expressed as:

$$X = \sum_{j=1}^r \underbrace{\sigma_j U_j V_j^T}_{\text{Sum of rank one matrices}}$$

where $\sigma_j \geq 0$, U_j is the j th column of U and V_j is the j th column of V

Example:

$$X = \begin{pmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{pmatrix}$$

↓ svd

$$X = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & \sqrt{5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ \sqrt{0.2} & 0 & 0 & 0 & \sqrt{0.8} \\ 0 & 0 & 0 & 1 & 0 \\ -\sqrt{0.8} & 0 & 0 & 0 & \sqrt{0.2} \end{pmatrix}$$

Low rank approximations

Without loss of generality say $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$

The main idea is that we develop a simpler version (low rank approximation) of X by dropping the smaller singular values.

A rank k approximation to X is given by:

$$\hat{X} = \sum_{j=1}^k \sigma_j u_j v_j^T \quad (\text{where } k \leq n)$$

$$m \begin{pmatrix} \vdots \\ \hat{X} \\ \vdots \end{pmatrix} = m \underbrace{\begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}}_{\substack{\text{Use first} \\ k \text{ columns} \\ \text{of } U}} \begin{pmatrix} \sigma_1 & \dots & \sigma_k \end{pmatrix} \underbrace{\begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix}}_{\substack{\text{Use first } k \\ \text{columns of } V \\ \text{to create } V^T}}$$

Example : Rank 2 approximation to X on the previous page is :

$$\hat{X} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \\ 0 & 0 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 4 & 0 \\ 0 & 3 \end{pmatrix} \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

Svd
approximation

$$= \begin{pmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \end{pmatrix}$$

One of the big advantages of performing a low rank approximation of X is that to represent a $m \times n$ matrix, the number of parameters needed = $m k + k + n k$
 $= k(m+n+1)$

For example :

Say $m = 200$, $n = 320$. This matrix has 64000 entries.

If we use $k = 20$, we would need $20(200+320+1) = 10420$ entries to approximate the matrix.

Relation to eigenvalues

The SVD applies to any $m \times n$ (rectangular) matrix where eigenvalue decomposition applies to square matrices

$$X = U S V^T$$

$$X X^T = U S V^T V S^T U^T$$

$$= U S I S^T U^T \quad (\text{since } V^T V = I)$$

$$= U (S S) U^T \quad (\text{Since } S^T = S \text{ is diagonal})$$

$$= U S^2 U^T$$

Thus the column vectors in U are the eigenvectors of $X X^T$ and the nonzero elements σ_k are the square root of the nonzero eigenvalues of $X X^T$

$$\text{Similarly } X^T X = V S U^T U S V^T$$

$$= V (S S) V^T$$

$$= V S^2 V^T$$

$$\text{Here } S^2 = \begin{pmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \\ & & & 0 \end{pmatrix}$$

The Frobenius norm of a matrix is defined as,

$$\|X\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n x_{ij}^2}$$

Equivalently $\|X\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_n^2}$ where $\sigma_1, \dots, \sigma_n$ are the singular values of X .

In approximating a matrix by a low rank matrix, the accuracy of the approximation (in terms of the explained variance) is

given as:

$$\frac{\|\hat{X}\|_F^2}{\|X\|_F^2} = \frac{\sigma_1^2 + \dots + \sigma_r^2}{\sigma_1^2 + \dots + \sigma_n^2}$$

Analytics on image data

install packages ("jpeg")
library(jpeg)

Package to read, write
and display images in
JPEG format

Read in data from image

`lly <- readJPEG("lly.jpg")`

This reads the image into a three dimensional array. The figure is of size 100 KB. This is read into an array of size $410 \times 640 \times 3$ where the dimensions are height \times width \times channels. (You can use `file.info("lly.jpg")`)

`min(lly[, , 1])`
`max(lly[, , 1])`

Values in the data representation
of image range from 0.0313
to 1

`max(abs(lly[, , 1] - lly[, , 2]))`
`max(abs(lly[, , 2], lly[, , 3]))`

It can be verified that
all three channels have
same values

Perform SVD

`S <- svd(lly[, , 1])`

Perform a singular value
decomposition on the
matrix

`str(S)`

`S$d` consists of the singular values (410)

`S$u` consists of left singular vectors (410×410)

`S$v` consists of right singular vectors (640×410)

Develop a low rank approximation

```
lky10 ← S$U[,1:10] %*% diag(S$d[1:10]) %*% t(S$v[,1:10])
```

Computes a rank 10 approximation to the matrix by using the top 10 singular values
Note `%*%` is used here for matrix multiplication

```
writeJPEG(lky10, "lky10.jpg")
```

Note that this writes a grayscale image of size about 16 KB. However this file is lossy and loses a lot of clarity.

```
lky50 ← S$U[,1:50] %*% diag(S$d[1:50])  
%*% t(S$v[,1:50])
```

Computes a rank 50 approximation to the matrix

```
writeJPEG(lky50, "lky50.jpg")
```

Note that this writes a grayscale image of size about 24 KB. which is already very clear

```
var ← cumsum(S$d^2)
```

```
plot(1:410, var)
```

```
plot(1:410, var/max(var))
```

In about 18 singular values out of 410, the approximation contains about 99% of the total variation in the picture

Read in from image

`parsy ← readJPEG("parsy.jpg")`

This reads the image of the flower into an array of size $600 \times 465 \times 3$.

Each of the numbers in the third dimension corresponds to an intensity in the R, G or B Channel.

Note that in this case `parsy[, , 1] ≠ parsy[, , 2] ≠ parsy[, , 3]`.

Perform SVD

`S1 ← svd(parsy[, , 1])`

`S2 ← svd(parsy[, , 2])`

`S3 ← svd(parsy[, , 3])`

`str(S1)`

List with `S1$d` consisting of 465 singular values, `S1$u` consisting of left singular vectors (600×465) and `S1$v` consisting of right singular vectors (465×465). Similarly for `S2` and `S3`.

Develop low rank approximation

`parsy50 ← array(dim = c(600, 465, 3))`

`parsy50[, , 1] ← S1$u[, 1:50] %*% diag(S1$d[1:50])
%*% t(S1$v[, 1:50])`

:

`parsy50[, , 3] ← S3$u[, 1:50] %*% diag(S3$d[1:50]) %*% t(S3$v[, 1:50])`

write JPEG(pnsy50, "pnsy50.jpg")

This creates an image of size about 40 KB which has some blurring along edges of the flower and leaves

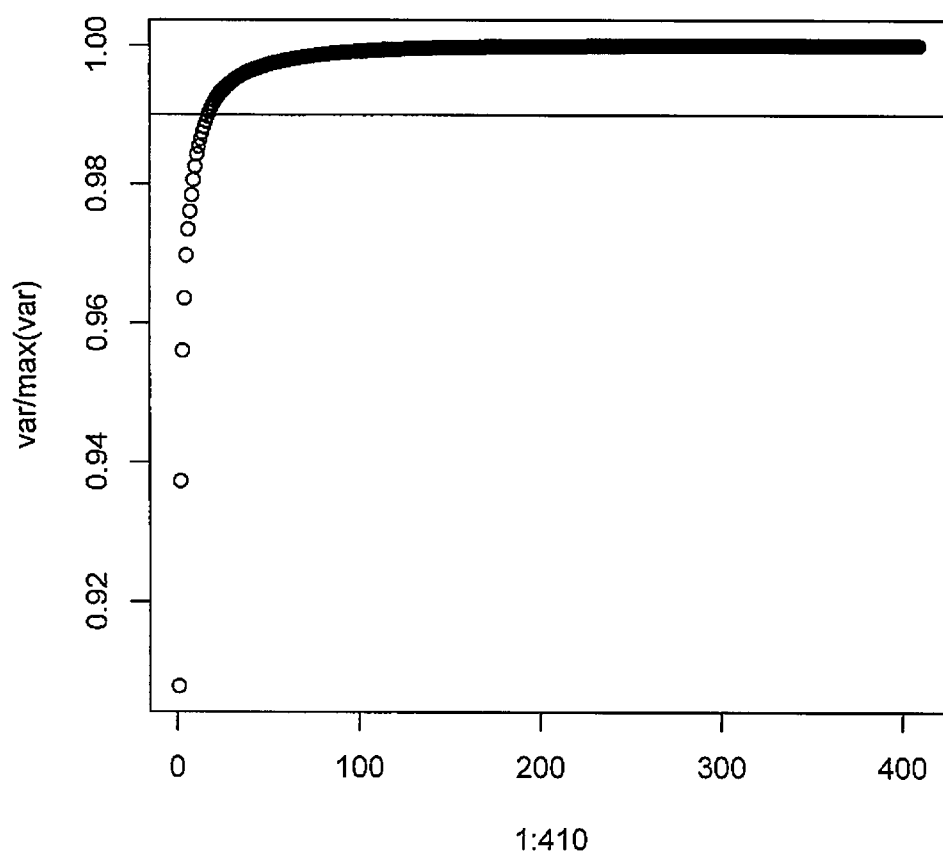
pnsy350 ← array(dim = c(600, 465, 3))

pnsy350[, , 1] ← SI\$V[, , 1:350] %*% diag(SI\$d[1:350])
%*% t(SI\$V[, , 1:350])

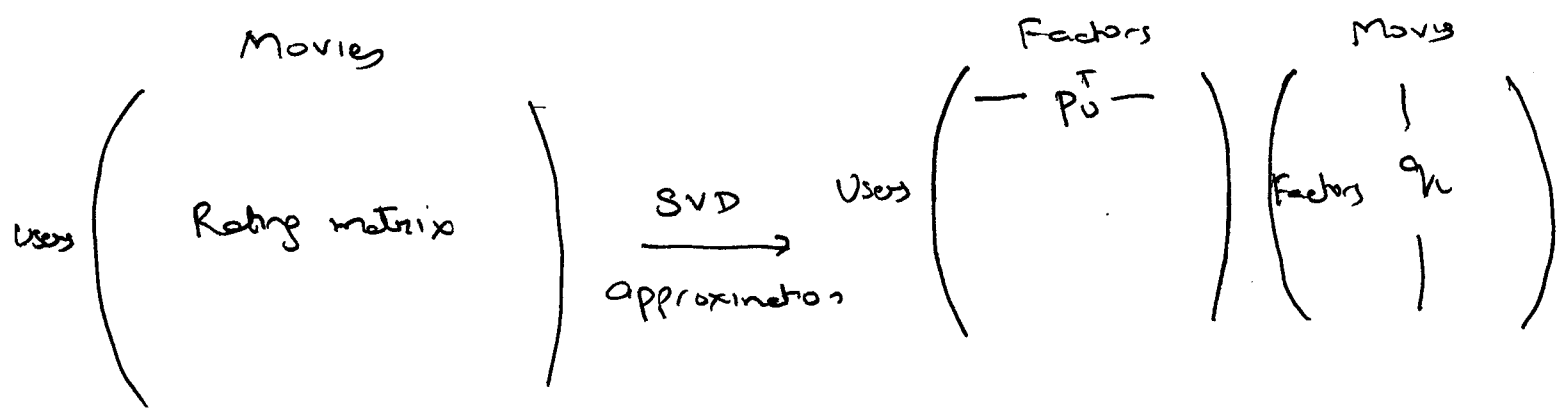
pnsy350[, , 3] ← . . .

write JPEG(pnsy350, "pnsy350.jpg")

This creates an image of size about 45 KB which is quite a good representation of the original image. The original file size here is around 80 KB.



SVD in recommendation systems



Users are associated with a vector p_u and movies (items) with a vector q_i , both of dimension " k ".

Rating is approximated as:

$$\hat{r}_{ui} = p_u^T q_i \quad (\text{or } q_i^T p_u)$$

Such a factorization is computed as a SVD where the diagonal matrix can be absorbed into either the U or V matrix. However the major challenge is that in recommendation systems such as Netflix, many entries are missing and hence conventional SVD fails.

There are a few ways to tackle this:

1) Impute (fill in) the entries using a reasonable method such as taking the average.

Apply SVD to the full matrix.

In the Netflix application, however this implies the matrix to be factorized is very large. Also inaccurate imputation can affect the data.

2) An alternate approach is to work only with the observed ratings and using regularization as in the LASSO model for overfitting. This was one of the approaches (models) used in the final Netflix winning contribution.

$$\min_{\substack{q_i, p_o \\ \forall (i,o) \\ \text{User-rating pairs}}} \sum_{(i,i) \in \text{Observed}} \left[(r_{oi} - q_i^T p_o)^2 + \lambda (\|q_i\|^2 + \|p_o\|^2) \right]$$

One can use simple stochastic gradient methods to solve this problem. (also known as incremental gradient method)

For example:

Start with q_i & p_u chosen arbitrarily.

For each training set user-item pair,

$$\text{Compute error } e_{ui} = r_{ui} - q_i^T p_u$$

Use a gradient method by

$$q_i \leftarrow q_i + \gamma \left((r_{ui} - q_i^T p_u) p_u - \lambda q_i \right)$$

$$p_u \leftarrow p_u + \gamma \left((r_{ui} - q_i^T p_u) q_i - \lambda p_u \right)$$

Gradients negative

This can be implemented using the `funkSVD(.)` function in the `recommenderlab` package.

`Data1 <- Data`

`Data1[spl1c, spl2c] <- NA`

`P <- funkSVD (Data1, k=2)` \nearrow k sets the rank of the SVD approximation

`Predict <- tcrossprod (P %*% U, P %*% V)`

This is equivalent to `P %*% U %*% t(P %*% V)`

$$\text{Err} \leftarrow \sqrt{\text{mean} \left(\left(\text{Predict}[\text{spl1c}, \text{spl2c}] - \text{Data}[\text{spl1c}, \text{spl2c}] \right)^2 \right)}, \text{na.rm=TRUE}$$

You can vary k to get better approximation.

As $k \uparrow$, running time increases

In the training set, a smaller probe set was excluded. The models got trained using this data so as to minimize the errors in the probe set.

Afterwards the probe set was included & then using the same parameters, the second training was done on the entire training set. This was used to make predictions on the qualifying set.

Initially each predictive model was chosen to minimize the RMSE on the probe set.

Afterwards the team realized that the best blend of models is the one where the models were less correlated with the rest of the ensemble (blend) & achieved a low RMSE individually.

So instead of building each model individually, they built it iteratively each time adding new predictive models to blend with the earlier set of predictive models to minimize blend RMSE. This helped in the final winning entry. Details can be found on netflixprize.com.

Final Comments on Netflix Prize

The final solution that won the Netflix prize involved a blend of various models developed by three teams. This involved models such as:

- Nearest neighbor models (K nearest neighbors)
- SVD models (matrix factorization)
- Models that incorporated temporal variations
- Regression models - Classification models

Over 500 models were used in the final solution. This was important in winning the final competition but in the winners own words, a few good models could already land them on the leaderboard.

The approach in blending the various predictive models was an important step in the Netflix prize winning entry.