



Other sites

- [Jobs for R-users](#)
- [SAS blogs](#)

Matrix factorization

March 10, 2015

By [saurabhat](#)

Like 65

Share

Tweet

Share

6

(This article was first published on [sane.a.lytics » R](#), and kindly contributed to [R-bloggers](#))



Or fancy words that mean very simple things.

At the heart of most data mining, we are trying to represent complex things in a simple way. The simpler you can explain the phenomenon, the better you understand. It's a little zen – compression is the same as understanding.

Warning: Some math ahead.. but stick with it, it's worth it.

When faced with a matrix of very large number of users and items, we look to some classical ways to explain it. One favorite technique is Singular Value Decomposition, affectionately nicknamed SVD. This says that we can explain *any* matrix A in the form

$$A = USV^T$$

Here

A is of size num_users and num_products

U is of size num_users and num_users

S is a rectangular diagonal matrix of size num_users and num_products

V is of size num_products and num_products

This is cool because the numbers in the diagonal S are decreasing values of variance (roughly speaking). So the first number captures the most variance in A , the second, less so, and so on, till all the numbers put together capture all the variance of A .

You see where this is going. If we are comfortable with explaining only 10% of the variance, we can do so by only taking some of these values. We can then compress A .

The other nice thing about all this is eigenvectors (which those roughly represent) are orthogonal. In other words, they are perpendicular to each other and there aren't any mixed effects. If you don't understand this paragraph, don't worry. Keep on.

Consider that above can be rewritten as

$$A = U\sqrt{S}\sqrt{S}V^T$$

or –

$$Y = X\Theta^T$$



that approximates A into $U \cdot V^T$ is now called a low rank approximation of A , or truncated SVD.

It is important to understand what these matrices represent. X is the representation of users in some low dimension space (say romance, action). And Θ is the representation of products in the very same low dimension space. However, we don't know exactly what these factors mean. They could be romance/action or they could NYC/Dallas. This is also why this method is sometimes called **Latent** Factor Matrix Factorization.. wow, quite a mouthful.

In my chapter in the book [Data Mining Applications with R](#), I go over different themes of matrix factorization models (and other animals as well). But for now, I am only going to cover the basic one that works very well in practice. And yes, won the Netflix prize.

There is one problem with our formulation – SVD is only defined for dense matrices. And our matrix is usually sparse.. very sparse. Netflix's challenge matrix was 1% dense, or 99% sparse. In my job at Rent the Runway, it is only 2% dense. So what will happen to our beautiful formula?

Machine learning is sort of a bastard science. We steal from beautiful formulations, complex biology, probability, Hoeffding's inequality, and derive rules of thumb from it. Or as an artist would say – get inspiration.

So here is what we are going to do. We are going to ignore all the null values when we solve this model. Our cost function now becomes

$$J = R(Y - X\Theta^T)^2 + \lambda(||X||^2 + ||\Theta||^2)$$

Here $Y - X\Theta^T$ is the difference between observed data and our prediction. R is simply a matrix with 1 where we have a value and 0 where we do not. Multiplying these two we are only considering the cost when we observe a value. We are using L2 regularization of magnitude λ . We are going to divide by 2 to make all other math easier. The cost is relative so it doesn't matter.

$$J = \frac{1}{2}R(Y - X\Theta^T)^2 + \frac{1}{2}\lambda(||X||^2 + ||\Theta||^2)$$

Using this our gradients become –

$$\frac{\partial}{\partial X} = R(Y - X\Theta^T)\Theta + \lambda||X||$$

$$\frac{\partial}{\partial \Theta} = R(Y - X\Theta^T)^T X + \lambda||\Theta||$$

If we were wanted to minimize the cost function, we can move in the direction opposite to the gradient at each step, getting new estimates for X and Θ each time.

This looks easy enough. One last thing. We now have what we want to minimize, but how do we do it? R provides many optimization tools. There is a whole CRAN page on it. For our purpose we will use out of the box `optim()` function. This allows us to access a fast optimization method L-BFGS-B. It's not only fast, but also doesn't take too memory intensive which is desirable in our case.

We need to give it the cost function and the gradient that we have above. It also expects one gradient function, so we need to unroll it out to do both our gradients.

```

1  unroll_Vecs <- function (params, Y, R, num_users, num_movies,
2    # Unrolls vector into X and Theta
3    # Also calculates difference between predction and actual
4
5    endIdx <- num_movies * num_features
6
7    X      <- matrix(params[1:endIdx], nrow = num_movies, ncol =
8    Theta <- matrix(params[(endIdx + 1):(endIdx + (num_users *
9      nrow = num_users, ncol = num_features)
10
11  Y_dash  <-  (((X %*% t(Theta)) - Y) * R) # Prediction error
12
13  return(list(X = X, Theta = Theta, Y_dash = Y_dash))
14 }
15
16 J_cost <- function(params, Y, R, num_users, num_movies, num_f
17 # Calculates the cost
18
```

```

21 Theta <- unrolled$Theta
22 Y_dash <- unrolled$Y_dash
23
24 J <- .5 * sum( Y_dash ^2) + lambda/2 * sum(Theta^2) + la
25
26 return (J)
27 }
28
29 grr <- function(params, Y, R, num_users, num_movies, num_featu
30 # Calculates the gradient step
31 # Here lambda is the regularization parameter
32 # Alpha is the step size
33
34 unrolled <- unroll_Vecs(params, Y, R, num_users, num_movies,
35 X <- unrolled$X
36 Theta <- unrolled$Theta
37 Y_dash <- unrolled$Y_dash
38
39 X_grad <- (( Y_dash %%% Theta) + lambda * X )
40 Theta_grad <- (( t(Y_dash) %%% X) + lambda * Theta )
41
42 grad = c(X_grad, Theta_grad)
43 return(grad)
44 }
45
46 # Now that everything is set up, call optim
47 print(
48 res <- optim(par = c(runif(num_users * num_features), runif(
49 fn = J_cost, gr = grr,
50 Y=Y, R=R,
51 num_users=num_users, num_movies=num_movies,num_
52 lambda=lambda, alpha = alpha,
53 method = "L-BFGS-B", control=list(maxit=maxit,
54 )

```

gistfile1.r hosted with by GitHub

[view raw](#)

This is great, we can iterate a few times to approximate users and items into some small number of categories, then predict using that.

I have coded this into another package [recommenderlabrats](#).

Let's see how this does in practice against what we already have. I am going to use the same scheme as last post to evaluate these predictions against some general ones. I am not using Item Based Collaborative Filtering because it is very slow

```

1 require(recommenderlab) # Install this if you don't have it al
2 require(devtools) # Install this if you don't have this alread
3 # Get additional recommendation algorithms
4 install_github("sanealytics", "recommenderlabrats")
5
6 data(MovieLense) # Get data
7
8 # Divvy it up
9 scheme <- evaluationScheme(MovieLense, method = "split", train
10 k = 1, given = 10, goodRating = 4)
11
12 scheme
13
14 # register recommender
15 recommenderRegistry$set_entry(
16 method="RSVD", dataType = "realRatingMatrix", fun=REAL_RSVD,
17 description="Recommender based on Low Rank Matrix Factorizat
18
19 # Some algorithms to test against
20 algorithms <- list(
21 "random items" = list(name="RANDOM", param=list(normalize =
22 "popular items" = list(name="POPULAR", param=list(normalize
23 "user-based CF" = list(name="UBCF", param=list(normalize = "
24 method="Cosin
25 nn=50, minRat
26 "Matrix Factorization" = list(name="RSVD", param=list(catego
27 lambda

```



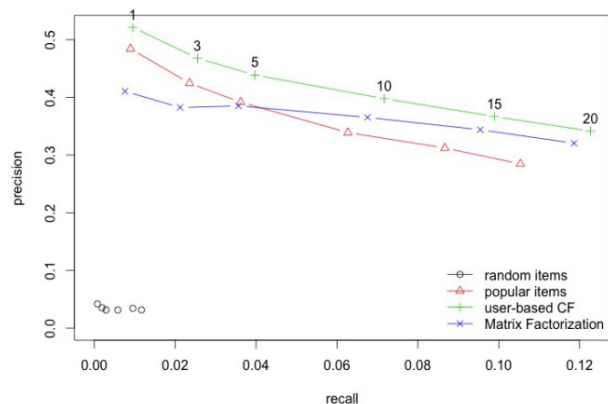
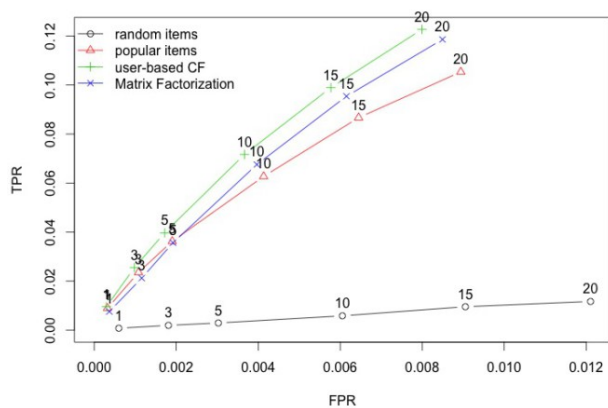
```

30
31 # run algorithms, predict next n movies
32 results <- evaluate(scheme, algorithms, n=c(1, 3, 5, 10, 15, 20))
33
34 # Draw ROC curve
35 plot(results, annotate = 1:4, legend="topleft")
36
37 # See precision / recall
38 plot(results, "prec/rec", annotate=3)

```

recommenderlab-test RSVD hosted with by GitHub

[view raw](#)



It does pretty well. It does better than POPULAR and is equivalent to UBCF. So why use this over UBCF or the other way round?

This is where things get interesting. This algorithm can be sped up quite a lot and more importantly, parallelised. It uses way less memory than UBCF and is more scalable.

Also, if you have already calculated Θ , i.e. your items in some lower dimensional space, you might get away with just putting the users in that space. Now things become really fast because all you have to do is keep Θ fixed and figure out X .

I have gone ahead and implemented a version where we just calculate Θ , I leave it to you as an exercise to modify the code above to test this out as well. The algorithm is called RSVD_SPLIT. Also feel free to try other values of categories, lambda and maxit and see how things change.

On the other hand, the latent categories are very hard to explain. For UBCF you can say this user is similar to these other users. For IBCF, one can say this item that the user picked is similar to these other items. But that not the case for this particular flavor of matrix factorization. We will re-evaluate these limitations later.

The hardest part for a data scientist is to disassociate themselves from their dear models and watch them objectively in the wild. Our simulations and metrics are always imperfect but necessary to optimize. You might see your favorite model crash and burn. And a lot of times simple linear regression will be king. The job is to objectively

Good luck.

Comments: 11

Like 65

Share

Tweet

Share 6

Related



Matrix factorization
In "R bloggers"



The guts of a statistical
factor model
In "R bloggers"

Factor models of
variance in finance
In "R bloggers"

65

Like

Share

Tweet

Share 6

To leave a comment for the author, please follow the link and comment on their blog:
[sane.a.lytics » R.](#)

R-bloggers.com offers [daily e-mail updates](#) about [R](#) news and [tutorials](#) on topics such as: [Data science](#), [Big Data](#), [R jobs](#), visualization ([ggplot2](#), [Boxplots](#), [maps](#), [animation](#)), programming ([RStudio](#), [Sweave](#), [LaTeX](#), [SQL](#), [Eclipse](#), [git](#), [hadoop](#), [Web Scraping](#)) statistics ([regression](#), [PCA](#), [time series](#), [trading](#)) and more...

If you got this far, why not **subscribe for updates** from the site?
Choose your flavor: [e-mail](#), [twitter](#), [RSS](#), or [facebook](#)...

Like 65

Share

Tweet

Share 6

Comments are closed.

Search & Hit Enter

Recent popular posts

- [Latest on the Julia Language \(vs. R\)](#)
- [Bad Coder, Bad Coder!](#)
- [dplyr do: Some Tips for Using and Programming](#)

Most visited articles of the week

- [How to write the first for loop in R](#)
- [Installing R packages](#)
- [R tutorials](#)
- [In-depth introduction to machine learning in 15 hours of expert videos](#)
- [Using apply, sapply, lapply in R](#)
- [Scatterplots](#)
- [How to perform a Logistic Regression in R](#)
- [Latest on the Julia Language \(vs. R\)](#)
- [Computing and visualizing PCA in R](#)

Sponsors