

SVD, Toplu Tavsiye (Collaborative Filtering)

Diyelim ki Star Trek (ST) dizisini ne kadar begendigini 4 tane kullanıcı sezonlara göre isaretlemiş. Bu örnek veriyi alttaki gibi gösterelim.

```
from pandas import *

d = np.array(
    [[5, 5, 0, 5],
     [5, 0, 3, 4],
     [3, 4, 0, 3],
     [0, 0, 5, 3],
     [5, 4, 4, 5],
     [5, 4, 5, 5]])

data = DataFrame (d.T,
    columns=['S1', 'S2', 'S3', 'S4', 'S5', 'S6'],
    index=['Ben', 'Tom', 'John', 'Fred'])
print data
```

	S1	S2	S3	S4	S5	S6
Ben	5	5	3	0	5	5
Tom	5	0	4	0	4	4
John	0	3	0	5	4	5
Fred	5	4	3	3	5	5

Veriye göre Tom, ST dizisinin 3. sezonunu 4 seviyesinde sevmiş. 0 değeri o sezonun seyredilmediğini gösteriyor.

Toplu Tavsiye algoritmaları verideki diğer kişilerin bir ürünü, diziyi, vs. ne kadar begendığının verisinin diğer "benzer" kişilere tavsiye olarak sunulabilir, ya da ondan önce, bir kişinin daha almadığı ürünü, seyretmediği sezonu, dinlemediği müziği ne kadar "begeneceğini" tahmin eder. Kaggle sitesi üzerinden yapılan unlu Netflix yarışmasının amacı buydu - ayrıca tahmin edilen ve gerçek beğeni notunun hata payının hesabi için RMSE hesabi kullanılmıştı.

Peki benzerliğin kriteri nedir, ve benzerlik nelerin arasında ölçülür?

Benzerlik, ürün seviyesinde, ya da kişi seviyesinde yapılabilir. Eğer ürün seviyesinde ise, tek bir ürün için tüm kullanıcıların verdiği nota bakılır. Eğer kullanıcı seviyesinde ise, tek kullanıcının tüm ürünlere verdiği beğeni notları vektörü kullanılır.

Mesela 1. sezondan hareketle, o sezonu begenen kişilere o sezona benzer diğer sezonlar tavsiye edilebilir. Kisiden hareketle, mesela John'a benzeyen diğer kişiler bulunarak onların begendigi ürünler John'a tavsiye edilebilir.

Ürün ya da kişi bazında olsun, benzerliği hesaplamanın da farklı yolları var. Genel olarak benzerlik ölçütünün 0 ile 1 arasında değişen bir sayı olmasını tercih ediyoruz ve tüm ayarları ona göre yapıyoruz. Diyelim ki elimizde beğeni notlarını taşıyan A, B vektörleri var (baska veri turu de taşıyor olabilir tabii), ve bu vektörlerin içinde beğeni notları var. Benzerlik çeşitleri şöyle:

Oklit Benzerligi (Euclidian Similarity)

Bu benzerlik $1/(1 + \text{mesafe})$ olarak hesaplanır. Mesafe karelerin toplamının karekoku (yani Oklitsel mesafe, ki isim buradan geliyor). Bu yuzden mesafe 0 ise (yani iki "sey" arasinda hic mesafe yok, birbirlerine cok yakinlar), o zaman hesap 1 dondurur (mukemmel benzerlik). Mesafe arttikca bolen buyudugu icin benzerlik sifira yaklasir.

Pearson Benzerligi

Bu benzerligin Oklit'ten farklilik, sayi buyuklugune hassas olmamasidir. Diye-lim ki birisi her sezonu 1 ile begenmis, digeri 5 ile begenmis, bu iki vektorun Pearson benzerligine gore birbirine esit cika-r. Pearson -1 ile +1 arasinda bir deger dondurur, alttaki hesap onu normalize ederek 0 ile 1 arasina ceker.

Kosinus Benzerligi (Cosine Similarity)

iki vektoru geometrik vektor olarak gorur ve bu vektorlerin arasinda olusan aci-yi (daha dogrusu onun kosinusunu) farklilik olcutu olarak kullanir.

$$\cos \theta = \frac{A \cdot B}{\|A\| \|B\|}$$

```
from numpy import linalg as la
def euclid(inA,inB):
    return 1.0/(1.0 + la.norm(inA - inB))

def pearson(inA,inB):
    if len(inA) < 3 : return 1.0
    return 0.5+0.5*np.corrcoef(inA, inB, rowvar = 0)[0][1]

def cos_sim(inA,inB):
    num = float(np.dot(inA.T,inB))
    denom = la.norm(inA)*la.norm(inB)
    return 0.5+0.5*(num/denom)

print np.array(data.ix['Fred'])
print np.array(data.ix['John'])
print np.array(data.ix['Ben'])
print pearson(data.ix['Fred'],data.ix['John'])
print pearson(data.ix['Fred'],data.ix['Ben'])

[5 4 3 3 5 5]
[0 3 0 5 4 5]
[5 5 3 0 5 5]
0.551221949943
0.906922851283

print cos_sim(data.ix['Fred'],data.ix['John'])
print cos_sim(data.ix['Fred'],data.ix['Ben'])

0.898160909799
0.977064220183
```

Simdi tavsiye mekanigine gelelim. En basit tavsiye yontemi, mesela kisi bazli olarak, bir kisiye en yakin diger kisileri bulmak (matrisin tamamina bakarak) ve onların begendikleri urunu istenilen kisiye tavsiye etmek. Benzerlik icin ustteki olcutlerden birini kullanmak.

Fakat belki de elimizde cok fazla urun, ya da kullanıcı var. Bir boyut azaltma islemi yapamaz miyiz?

Evet. SVD yontemi burada da isimize yarar.

$$A = USV$$

elde edecegimiz icin, ve S icindeki en buyuk degerlere tekabul eden U, V degerleri siralanmis olarak geldigi icin U, V'nin en bastaki degerlerini almak bize "en onemli" bloklari verir. Bu en onemli kolon ya da satirlari alarak azaltilmis bir boyut icinde benzerlik hesabi yapmak islemlerimizi hizlandirir. Bu azaltilmis boyutta kumeleme algoritmalarini devreye sokabiliriz; U'nun mesela en onemli iki kolonu bize iki boyuttaki sezon kumelerini verebilir, V'nin en onemli iki (en ust) satiri bize iki boyutta bir kisi kumesi verebilir.

O zaman begeni matrisi uzerinde SVD uygulayalim,

```
from numpy.linalg import linalg as la
U,Sigma,V=la.svd(data, full_matrices=False)
print data.shape
print U.shape, Sigma.shape, V.shape
u = U[:, :2]
vt=V[:, :2].T
print 'u', u
print 'vt', vt
print u.shape, vt.shape

(4, 6)
(4, 4) (4,) (4, 6)
u [[-0.57098887 -0.22279713]
   [-0.4274751  -0.51723555]
   [-0.38459931  0.82462029]
   [-0.58593526  0.05319973]]
vt [[-0.44721867 -0.53728743]
     [-0.35861531  0.24605053]
     [-0.29246336 -0.40329582]
     [-0.20779151  0.67004393]
     [-0.50993331  0.05969518]
     [-0.53164501  0.18870999]]
(4, 2) (6, 2)
```

degerleri elimize gecer. U ve VT matrisleri

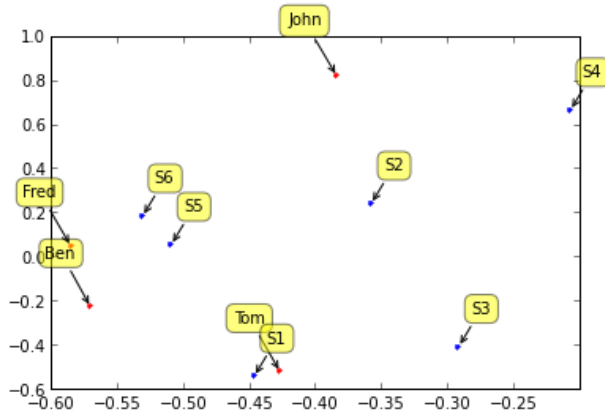
```
def label_points(d,xx,yy,style):
    for label, x, y in zip(d, xx, yy):
        plt.annotate(
            label,
```

```

xy = (x, y), xytext = style,
textcoords = 'offset points', ha = 'right', va = 'bottom',
bbox = dict(boxstyle = 'round,pad=0.5', fc = 'yellow', alpha = 0.5),
arrowprops = dict(arrowstyle = '->', connectionstyle = 'arc3,rad=0'))

plt.plot(u[:,0],u[:,1], 'r.')
label_points(data.index, u[:, 0], u[:, 1], style=(-10, 30))
plt.plot(vt[:,0],vt[:,1], 'b.')
label_points(data.columns, vt[:, 0], vt[:, 1], style=(20, 20))
plt.savefig('svdrecom_1.png')

```



Cok guzel! SVD bize urun bazinda sezon 5 ve 6'nin bir kume olusturdugunu, Ben ve Fred'in de kisi bazinda ayri bir kume oldugunu gosterdi.

Azaltilmis boyutlari nasil kullaniriz? Yeni bir kisiyi (mesela Bob) ele alinca, bu kisinin verisini oncelikle aynen diger verilerin indirgendigi gibi azaltilmis boyuta "indirgememiz" gerekiyor. Cunku artik islem yaptigimiz boyut orasi. Peki bu indirgemeyi nasil yapariz? SVD genel formulu hatirlarsak,

$$A = USV$$

Azaltilmis ortamda

$$A = U_k S_k V_k$$

Diyelim ki gitmek istedigimiz nokta azaltilmis U , o zaman U_k 'yi tek basina birakalim,

$$AV_k^{-1} = U_k S_k V_k^{-1}$$

U_k, V_k matrisleri ortonormal, o zaman $V_k^{-1}V_k = I$ olacak, yani yokolacak

$$AV_k^{-1} = U_k S$$

Benzer sekilde

$$AV_k^{-1}S^{-1} = U_k$$

Cok fazla ters alma islemi var, her iki tarafın devrigini alalım

$$(S^{-1})^T(V_k^{-1})^T A^T = U_k^T$$

$V_k^{-1} = V_k^T$ o zaman üstteki formül devrigin devrigini almak demektir, yani tekrar basa donmus oluyoruz, demek ki V_k degismeden kalıyor

$$(S^{-1})^T V_k A^T = U_k^T$$

S ise kosegen matris, onun tersi yine kosegen, kosegen matrisin devrigi yine kendisi

$$S^{-1} V_k A^T = U_k^T$$

Bazi kod ispatlari, u 'nun ortonormal olmasi:

```
print np.dot(u.T,u)

[[ 1.00000000e+00  4.83147593e-18]
 [ 4.83147593e-18  1.00000000e+00]]
```

Dogal olarak $..e-17$ gibi bir sayi sifira cok yakin, yani sifir kabul edilebilir. Devrik ve tersin ayni oldugunu gosterelim: Iki matrisi birbirinden cikartip, cok kucuk bir sayidan buyukluge gore filtreleme yapalim, ve sonuc icinde bir tane bile True olup olmadigini kontrol edelim,

```
print not any(U.T-la.inv(U) > 1e-15)

True
```

Yeni Bob verisi

```
bob = np.array([5,5,0,0,0,5])
```

O zaman

```
print bob.T.shape
print u.shape
S_k = np.eye(2)*Sigma[:2]
bob_2d = np.dot(np.dot(la.inv(S_k),vt.T),bob.T)
print bob_2d
```

```
(6, )
(4, 2)
[-0.37752201 -0.08020351]
```

Ustte `eye` ve `Sigma` ile ufak bir takla attik, bunun sebebi `svd` cagrisindan gelen `Sigma` sonucunun bir vektor olması ama ustteki islem için kosegen bir "matrise" ihtiyacımız olması. Eğer birim (identity) matrisini alıp onu `Sigma` ile carparsak, bu kosegen matrisi elde ederiz.

Simdi mesela kosinus benzerligi kullanarak bu izdusumlenmis yeni vektorun hangi diger vektorlere benzedigini bulalim.

```
for i,user in enumerate(u):
    print data.index[i],cos_sim(user,bob_2d)
```

```
Ben 0.993397525045
Tom 0.891664622942
John 0.612561691287
Fred 0.977685793579
```

Sonuca gore yeni kullanıcı Bob, en çok Ben ve Fred'e benziyor. Sonuca eristik! Artık bu iki kullanıcının yüksek not verdigi ama Bob'un hiç not vermedigi sezonlari alıp Bob'a tavsiye olarak sunabiliriz.

Movielens 1M Verisi

Bu veri seti 6000 kullanıcı tarafından yaklaşık 4000 tane filme verilen not / derece (rating) verisini iceriyor, 1 milyon tane not verilmiş, yani $4000 * 6000 = 24$ milyon olasılık içinde sadece 1 milyon veri noktası dolu. Bu oldukça seyrek bir matris demektir.

Verinin ham hali diger ders notlarimizi iceren ust dizinlerde var, veriyi SVD ile kullanilir hale getirmek için bu dizindeki `movielens_prep.py` adli script kullanilir. Islem bitince `movielens.csv` adli bir dosya script'te gorulen yere yazilacak. Bu dosyada olmayan derecelendirmeler, verilmemiş notlar bos olacaktır. Bu bosluklari sifirlarsak, seyrek matrisi o noktaları atlar. Ardından bu seyrek matris üzerinde seyrek SVD isletilebilir. Bu normal SVD'den daha hızlı isleyecektir.

Tavsiye kodlamamız için yazının başında anlatılan teknigi kullanacagiz, film verisi üzerinde boyut azaltılması yapılacak, benzer kullanıcı bulunacak, ve herhangi bir yeni kullanıcı / film kombinasyonu için bu diger benzer kullanıcının o filme verdigi not baz alınacak.

Veriyi eğitim ve test olarak iki parçaya böleceğiz. SVD eğitim bölümü üzerinde isletilecek.

Bu bağlamda, önemli bir diger konu eksik veri noktalarının SVD sonuçlarını nasıl etkileyeceği. Sonuçta eksik yerler `nan`, oradan sıfır yapıp ardından seyrek matris kodlaması üzerinden "atlanıyor" olabilir, fakat bu değerler atlanıyor (yani hızlı isleniyor, depolanıyor) olsa bile, onların sıfır olmasının bir anlamı yok mudur? Evet vardır. Not bakımından sıfır da bir not'tur, ve bu sebeple sonuçları istenmeyen biçimde etkileyebilir.

O zaman mevcut veriyi oyle bir degistirelim ki verilmemis notlar, yani sifir degerleri sonucu fazla degistirmesin.

Bunu yapmanin yollarindan biri her film icin bir ortalama not degeri hesaplamak, ve bu ortalama degeri o filme verilen tum not degerlerinden cikartmaktir. Bu isleme "sifir cevresinde merkezlemek" ismi de verilir, hakikaten mesela film j icin ortalama 3 ise, 5 degeri 2, 3 degeri sifir, 2 degeri -1 haline gelecektir. Bu bir ilerlemedir cunku ortalama 3 degeri zaten bizim icin "onemsiz" bir degerdir, tavsiye problemi baglaminda bizim en cok ilgilendigimiz sevilen filmler, ve sevilmeyen filmler. Bu degerler sirasiyla arti ve eksi degerlere donusecekler, ve SVD bu farklilik matematiksel olarak kullanabilme yetenegine sahip.

Altta Pandas `mean` cagrisi ile bu islemin yapildigini goruyoruz, dikkat, Pandas dataframe icinde `nan` degerleri olacaktir, ve Pandas bu degerleri atlamasi gerektigini bilir, yani bu degerler ortalamaya etki etmez. Ardindan merkezleme islemi egitim verisi uzerinde uygulaniyor.

```
import pandas as pd, os
import scipy.sparse as sps
df = pd.read_csv("%s/Downloads/movielens.csv" % os.environ['HOME'], sep=';')
print df.shape
df = df.ix[:,1:] # id kolonunu atla
df = df.ix[:, :3700] # sadece filmleri al
df_train = df.copy().ix[:5000, :]
df_test = df.copy().ix[5001:, :]
df_train[np.isnan(df_train)] = 0.0
movie_avg_rating = np.array(df_train.mean(axis=0))
df_train = df_train - movie_avg_rating
dfs_train = sps.coo_matrix(df_train)

df_train = np.array(df_train)
df_test = np.array(df_test)

print df_train.shape
print df_test.shape

__top_k__ = 10
import scipy.sparse.linalg as slin
import scipy.linalg as la
U, Sigma, V = slin.svds(dfs_train, k=__top_k__)
print U.shape, Sigma.shape, V.shape
Sigma = np.diag(Sigma)

(6040, 3731)
(5001, 3700)
(1039, 3700)
(5001, 10) (10,) (10, 3700)
```

Altta test verisi uzerinde satir satir ilerliyoruz, ve her satir (test kullanicisi) icinde film film ilerliyoruz. "Verilmis bir not" ariyoruz (cogunlukla not verilmemis oluyor cunku), ve buldugumuz zaman artik elimizde test edebilecegimiz bir sey var, o notu "sifirlayip" vektorun geri kalanini azaltilmis boyuta yansitiyoruz, ve sonra o boyuttaki tum diger U vektorleri icinde arama yapiyoruz, en yakin diger

kullaniciyi buluyoruz ve onun bu filme verdigi notu tahminimiz olarak kullanıyoruz.

Altta eger bulunan diger kullanıcı o filme not vermemisse, basitleştirme amaçlı olarak, o filmi atladık. Gerçek dünya şartlarında filme not vermiş ve yakın olan (en yakın olmasa da) ikinci, üçüncü kullanıcılar bulunup onların notu kullanılabilir. Hatta en yakın k tane kullanıcının ortalaması alınabilir (o kullanıcılar kNN gibi bir metotla bulunur belki), vs.

```
def euclid(inA,inB):
    return 1.0/(1.0 + la.norm(inA - inB))

rmse = 0; n = 0
for i,test_row in enumerate(df_test):
    for j, test_val in enumerate(test_row):
        # nan olmayan bir not buluncaya kadar ara
        if np.isnan(test_val): continue
        # bulduk, test satirini tamamen kopyala ve bulunan notu silerek
        # onu nan / sifir haline getir cunku yansitma (projection) oncesi
        # o notu 'bilmiyormus gibi' yapmamiz lazim.
        curr = test_row.copy()
        curr[j] = np.nan
        curr[np.isnan(curr)] = 0.

        proj_row = np.dot(np.dot(la.inv(Sigma),V),curr)

        sims = np.array(map(lambda x: euclid(x, proj_row), U[:,:__top_k__]))
        isim = np.argmax(sims)

        # eger bulunan kullanıcı o filme not vermemisse atla
        if np.isnan(df.ix[isim, j]): continue

        # egitim verisinde notlar sifir etrafında ortalanmış, tekrar
        # normal haline dondur
        est = df_train[isim, j]+movie_avg_rating[j]

        # gercek not
        real = df_test[i, j]

        print i, 'icin en yakin', isim, 'urun',j, 'icin oy', est, 'gercek', real
        rmse += (real-est)**2
        n += 1
        break # her kullanıcı için tek film test et
    if i == 20: break # 20 kullanıcı test et

print "rmse", np.sqrt(rmse / n)

0 icin en yakin 1903 urun 144 icin oy 5.0 gercek 5.0
1 icin en yakin 239 urun 144 icin oy 5.0 gercek 5.0
2 icin en yakin 2045 urun 844 icin oy 4.0 gercek 4.0
3 icin en yakin 4636 urun 0 icin oy 3.0 gercek 4.0
4 icin en yakin 139 urun 845 icin oy 4.0 gercek 5.0
5 icin en yakin 427 urun 1107 icin oy 4.0 gercek 5.0
6 icin en yakin 3620 urun 31 icin oy 4.0 gercek 4.0
```


7 için en yakın 1870 ürün 0 için oy 4.0 gerçek 3.0
8 için en yakın 4816 ürün 106 için oy 5.0 gerçek 5.0
9 için en yakın 3511 ürün 0 için oy 3.0 gerçek 4.0
10 için en yakın 3973 ürün 1212 için oy 5.0 gerçek 4.0
11 için en yakın 2554 ürün 287 için oy 4.0 gerçek 5.0
12 için en yakın 4733 ürün 31 için oy 4.0 gerçek 3.0
13 için en yakın 2339 ürün 9 için oy 4.0 gerçek 3.0
14 için en yakın 3036 ürün 10 için oy 4.0 gerçek 3.0
15 için en yakın 2748 ürün 253 için oy 5.0 gerçek 5.0
16 için en yakın 450 ürün 16 için oy 4.0 gerçek 4.0
17 için en yakın 1133 ürün 9 için oy 5.0 gerçek 2.0
18 için en yakın 3037 ürün 253 için oy 5.0 gerçek 4.0
19 için en yakın 1266 ürün 107 için oy 3.0 gerçek 3.0
20 için en yakın 537 ürün 253 için oy 5.0 gerçek 5.0
rmse 0.975900072949

Sonuc fena değil. Tavsiye programlarında RMSE 0.9 civarı iyi olarak bilinir, Netflix yarışmasında [3] mesela kazanan algoritma RMSE 0.85'e erismistir.

Kaynaklar

[1] <http://www.igvita.com/2007/01/15/svd-recommendation-system-in-ruby/>

[2] Harrington, P., Machine Learning in Action

[3] http://en.wikipedia.org/wiki/Netflix_Prize

[4] <http://stats.stackexchange.com/questions/31096/how-do-i-use-the-svd-in-collaborative-filtering>