



SVD ile PCA Hesaplamak

PCA bolumunde anlatilan yontem temel bileşenlerin hesabinda ozdegerler ve ozvektorler kullandi. Alternatif bir yontem Tekil Deger Ayristirma (Singular Value Decomposition -SVD-) uzereinden bu hesabi yapmaktir. SVD icin *Lineer Cebir Ders 29'a* bakabilirsiniz. Peki ne zaman klasik PCA ne zaman SVD uzereinden PCA kullanmali? Bir cevap belki mevcut kutuphanelerde SVD kodlamasinin daha iyi olmasi, ayristirmanin ozvektor / deger hesabindan daha hizli isleyebilmesi [6].

Ayrica birazdan gorecegimiz gibi SVD, kovaryans matrisi uzereinde degil, A 'nin kendisi uzereinde isletilir, bu hem kovaryans hesaplama asamasini atlamamizi, hem de kovaryans hesabi sirasinda ortaya cikabilecek numerik puruzlerden korunmamizi saglar (cok ufak degerlerin kovaryans hesabini bozabilecegi literatürde bahsedilmektedir).

PCA ve SVD baglantisina geelim:

Biliyoruz ki SVD bir matrisi su sekilde ayristirir

$$A = USV^T$$

U matrisi $n \times n$ dikgen (orthogonal), V ise $m \times m$ dikgen. S 'in sadece kosegeni uzereinde degerler var ve bu σ_j degerleri A 'nin tekil degerleri (singular values) olarak biliniyor.

Simdi A yerine AA^T koyalim, yani A 'nin kovaryans matrisinin SVD ayristirmasini yapalim, acaba elimize ne gececek?

$$\begin{aligned} AA^T &= (USV^T)(USV^T)^T \\ &= (USV^T)(VS^T U^T) \\ &= USS^T U^T \end{aligned}$$

S bir kosegen matrisi, o zaman SS^T matrisi de kosegen, tek farkla kosegen uzereinde artik σ_j^2 degerleri var. Bu normal.

SS^T yerine Λ sembolunu kullanalim, ve denklemleri iki taraftan (ve sagdan) U ile carparsak (unutmayalim U ortanormal bir matris ve $U^T U = I$),

$$AA^T U = U \Lambda U^T U$$

$$AA^T U = U \Lambda$$

Son ifadeye yakindan bakalim, U 'nun tek bir kolonuna, u_k diyelim, odaklanacak olursak, ustteki ifadede bu sadece kolona yonelik nasil bir esitlik cikartabilirdik? Soyle cikartabilirdik,

$$(AA^T)u_k = \sigma_k^2 u_k$$

Bu ifade tanidik geliyor mu? Ozdeger / ozvektor klasik yapısına eristik. Ustteki esitlik sadece ve sadece eger u_k , AA^T 'nin ozvektoru ve σ_k^2 onun ozdegeri ise gecerlidir. Bu esitligi tum U kolonlari icin uygulayabilecegimize gore demek ki U 'nun kolonlarinda AA^T 'nin ozvektorleri vardir, ve AA^T 'nin ozdegerleri A 'nin tekil degerlerinin karesidir.

Bu muthis bir bulus. Demek ki AA^T 'nin ozvektorlerini hesaplamak icin A uzereinde SVD uygulayarak U 'yu bulmamiz yeterli, kovaryans matrisini hesaplamak bile gerekmiyor! AA^T ozdegerleri uzereinde buyukluk karsilastirmasi icin ise A 'nin tekil degerlerine bakmak yeterli!

Ornek

İlk PCA yazısındaki örneğe dönelim, ve özvektörleri SVD üzerinden hesaplayalım. Önce eig usulu ile

```
In [2]: from pandas import *
data = read_csv("testSet.txt", sep="\t", header=None)
means = mean(data, axis=0)
meanless_data = data - means
cov_mat = cov(meanless_data, rowvar=0)
eigs,eigv = linalg.eig(mat(cov_mat))
eigv
```

```
Out[2]: matrix([[ -0.85389096, -0.52045195],
                [ 0.52045195, -0.85389096]])
```

Şimdi de SVD usulu ile

```
In [3]: U,s,Vt = svd(meanless_data.T,full_matrices=False)
U
```

```
Out[3]: array([[ -0.52045195, -0.85389096],
               [-0.85389096,  0.52045195]])
```

```
In [4]: np.dot(U.T,U)
```

```
Out[4]: array([[ 1.,  0.],
               [ 0.,  1.]])
```

Görüldüğü gibi aynı özvektörleri bulduk.

New York Times Yazıları Analizi

Şimdi daha ilginç bir örneğe bakalım. Bir araştırmacı belli yıllar arasındaki NY Times makalelerinde her yazıda hangi kelimenin kaç kere çıktığının verisini toplamış [1,2,3], bu veri 4000'ün aşkın kelime, her satır (yazı) için bir boyut (kolon) olarak kaydedilmiştir. Bu veri nytimes.csv üzerinde ek bir normalize işleminden sonra (bkz norm.R [4] ve nytimes4.csv), onun üzerinde boyut indirgeme yapabiliriz.

Veri setinde her yazı ayrıca ek olarak sanat (arts) ve müzik (music) olarak etiketlenmiş, ama biz PCA kullanarak bu etiketlere hiç bakmadan, verinin boyutlarını azaltarak acaba verinin "ayrılabilir" hale indirgenip indirgenemediğine bakalım. Sonra etiketleri veri üstüne koyup sonucun doğruluğunu kontrol edeceğiz.

Bakmak derken veriyi (en önemli) iki boyuta indirgeyip sonucu grafikleyeceğiz. İlla 2 olması gerekmez tabii, 10 boyuta indirgeyip (ki 4000'ün aşkın boyuttan sonra bu hala muhtemelen bir kazanım) geri kalanlar üzerinde mesela bir kümeleme algoritması kullanabilirdik.

Ana veriyi yükleyip birkaç satırını ve kolonlarını göstereyim.

```
In [12]: from pandas import *
import numpy.linalg as lin
nyt = read_csv ("nytimes.csv")
labels = nyt['class.labels']
nyt.ix[:8,102:107]
```

Out[12]:

	after	afternoon	afterward	again	against
0	1	0	0	0	0
1	1	1	0	0	0
2	1	0	0	1	2
3	3	0	0	0	0
4	0	1	0	0	0
5	0	0	0	1	2
6	7	0	0	0	1
7	0	0	0	0	0
8	0	0	0	0	0

Yuklemeyi yapip sadece etiketleri aldik ve onlari bir kenara koyduk. Simdi onemli bir normalizasyon islemi gerekiyor - ki bu isleme ters dokuman-frekans agirliklandirmasi (inverse document-frequency weighting -IDF-) ismi veriliyor - her dokumanda aşırı fazla ortaya cikan kelimelerin onemi ozellikle azaltiliyor, ki diger kelimelerin etkisi artabilsin.

IDF kodlamasi alttaki gibidir. Once class.labels kolonunu atariz. Sonra "herhangi bir deger iceren" her hucrenin 1 digerlerinin 0 olmasi icin kullanim DataFame uzerinde astype(bools) isletme numarasini kullaniriz, boylece asiri buyuk degerler bile sadece 1 olacaktir. Bazi diger islemler sonrasi her satiri kendi icinde tekrar normalize etmek icin o satirdaki tum degerlerin karesinin toplamının karekokunu aliriz ve satirdaki tum degerler bu karekok ile bolunur. Buna oklitsel (euclidian) normalizasyon denebilir.

Not: Oklitsel norm alirken toplamın hemen ardından çok ufak bir $1e-16$ degeri eklememize dikkat cekelim, bunu toplamın sıfır olma durumu için yapıyoruz, ki sonra sıfırla bölerken NaN sonucundan kacinalim.

```
In [13]: nyt = nyt.drop(['class.labels'],axis=1)
freq = nyt.astype(bool).sum(axis=0)
freq = freq.replace(0,1)
w = np.log(float(nyt.shape[0])/freq)
nyt = nyt.apply(lambda x: x*w,axis=1)
nyt = nyt.apply(lambda x: x / np.sqrt(np.sum(np.square(x))+1e-16), axis=1)
```

```
In [6]: nyt=nyt.ix[:,1:] # ilk kolonu atladik
        nyt.ix[:8,102:107]
```

Out[6]:

	afterward	again	against	age	agent
0	0	0.000000	0.000000	0.051085	0
1	0	0.000000	0.000000	0.000000	0
2	0	0.021393	0.045869	0.000000	0
3	0	0.000000	0.000000	0.000000	0
4	0	0.000000	0.000000	0.000000	0
5	0	0.024476	0.052480	0.000000	0
6	0	0.000000	0.008536	0.000000	0
7	0	0.000000	0.000000	0.000000	0
8	0	0.000000	0.000000	0.000000	0

SVD yapalim

```
In [7]: nyt = nyt - nyt.mean(0)
        u,s,v = lin.svd(nyt.T,full_matrices=False)
```

```
In [8]: s[:10]
```

```
Out[8]: array([ 1.41676764,  1.37161893,  1.31840061,  1.24567955,  1.20596873,
                1.18624932,  1.15118771,  1.13820504,  1.1138296 ,  1.10424634])
```

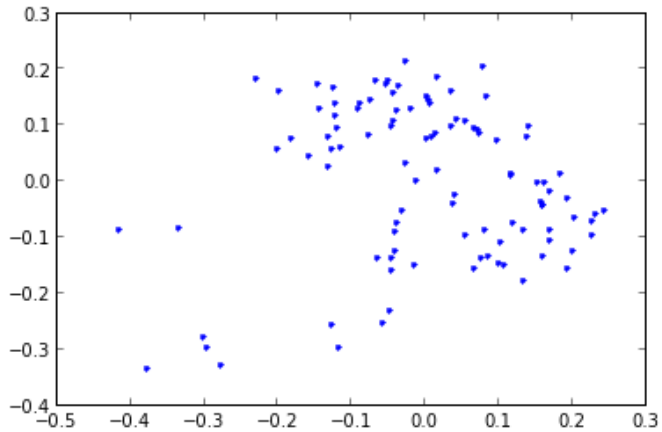
```
In [9]: u.shape
```

```
Out[9]: (4430, 102)
```

SVD'nin verdigi u icinden iki ozvektoru seciyoruz (en bastakiler, cunku Numpy SVD kodu bu ozvektorleri zaten siralanmis halde dondurur), ve veriyi bu yeni kordinata izdusumluyoruz.

```
In [10]: proj = np.dot(nyt, u[:, :2])
proj.shape
plot(proj[:,0],proj[:,1],'.')
```

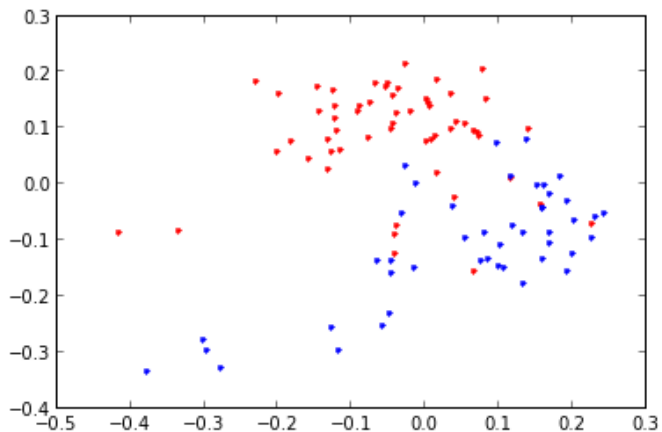
Out[10]: [



Simdi ayni veriyi bir de etiket bilgisini devreye sokarak cizdirelim. Sanat kirmizi muzik mavi olacak.

```
In [11]: arts =proj[labels == 'art']
music =proj[labels == 'music']
plot(arts[:,0],arts[:,1], 'r.')
plot(music[:,0],music[:,1], 'b.')
```

Out[11]: [



Goruldugu gibi veride ortaya cikan / ozvektorlerin kesfettigi dogal ayirim, hakikaten dogruymus.

Metotun ne yaptigina dikkat, bir suru boyutu bir kenara atmamiza ragmen geri kalan en onemli 2 boyut uzzerinden net bir ayirim ortaya cikartabiliyoruz. Bu PCA yonteminin iyi bir is becerdigini gosteriyor, ve kelime sayilarinin makalelerin icerigi hakkında ipucu icerdigini ispatliyor.

[1] Cosma Rohilla Shalizi, Advanced Data Analysis from an Elementary Point of View

[2] <http://www ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC2008T19>

[3] <http://www.stat.cmu.edu/~cshalizi/490/pca/>

[5] <http://www.math.nyu.edu/faculty/goodman/teaching/RPME/notes/Section3.pdf>

[6] Lineer Cebir notlarımızda SVD türetilmesine bakınca özdeğer/vektor mantığına atıf yapıldığını görebiliriz ve akla şu gelebilir; "özdeğer / vektor rutini işletmekten kurtulalım dedik, SVD yapıyoruz, ama onun içinde de özdeğer/vektor hesabı var". Fakat bunu belirtmek gerekir ki SVD numerik hesabını yapmanın tek yöntemi özdeğer/vektor yöntemi değildir. Mesela Numpy Linalg kutuphanesi içindeki SVD, LAPACK dgesdd rutinini kullanır ve bu rutin iç kodlamasında QR, ve bir tür böl / istila et (divide and conquer) algoritması işletmektedir.