

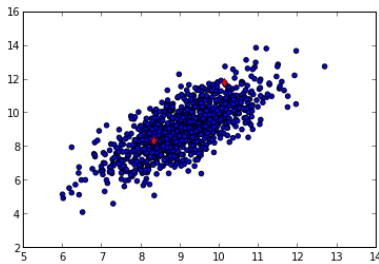
Temel Bileşen Analizi (Principal Component Analysis -PCA-)

PCA yöntemi boyut azaltan yöntemlerden biri, takip edilmeden (unsupervised) işleyebilir. Ana fikir veri noktalarının izdusumunun yapılacağı yonlar bulmaktır ki bu yonlar bağlamında (izdusum sonrası) noktaların arasındaki sayısal varyans (empirical variance) en fazla olsun, yani noktalar grafik bağlamında düşünürsek en "yayılmış" şekilde bulunsunlar. Böylece birbirinden daha uzaklaşan noktaların mesela daha rahat kümelenebileceğini umabiliriz. Bir diğer amaç, hangi değişkenlerin varyansının daha fazla olduğunun görülmesi üzerine, o değişkenlerin daha önemli olabileceğinin anlaşılması. Örnek olarak alttaki grafiğe bakalım,

```
from pandas import *
data = read_csv("testSet.txt", sep="\t", header=None)
print data[:10]
```

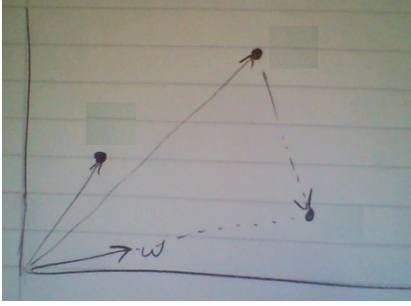
	0	1
0	10.235186	11.321997
1	10.122339	11.810993
2	9.190236	8.904943
3	9.306371	9.847394
4	8.330131	8.340352
5	10.152785	10.123532
6	10.408540	10.821986
7	9.003615	10.039206
8	9.534872	10.096991
9	9.498181	10.825446

```
plt.scatter(data.ix[:,0],data.ix[:,1])
plt.plot(data.ix[1,0],data.ix[1,1], 'rd')
plt.plot(data.ix[4,0],data.ix[4,1], 'rd')
plt.savefig('pca_1.png')
```



PCA ile yapmaya çalıştığımız öyle bir yon bulmak ki, x veri noktalarının tamamının o yone izdusumu yapılıncaya sonuc olacak, "izdusumu yapılmış" z 'nin varyansı en büyük olsun. Bu bir maksimizasyon problemidir. Fakat ondan önce x nedir, z nedir bunlara yakından bakalım.

Veri x ile tüm veri noktaları kastedilir, fakat PCA probleminde genellikle bir "vektörün diğeri üzerine" yapılan izdusumu, "daha optimal bir w yonu bulma", ve "o yone doğru izdusum yapmak" kelimeleri kullanılır. Demek ki veri noktalarını bir vektör olarak görmeliyiz. Eğer üstte kırmızı ile işaretlenen iki noktayı alırsak (bu noktalar verideki 1. ve 4. sıradaki noktalar),



gibi bir görüntüden bahsediyoruz. Hayali bir w kullandık, ve noktalardan biri veri noktası, w üzerine izdüşüm yapılarak yeni bir vektörü / noktayı ortaya çıkartılıyor. Genel olarak ifade edersek, bir nokta için

$$z = w^T x$$

Yapmaya çalıştığımız sayısal varyansı maksimize etmek demistik. Özel bir izdüşüm yönünü referans alırsak, en büyük $\text{Var}(z_1)$ 'i bulacağız. Bu yeni "seyin" beklentisi ve varyansına bakalım,

$$E[w^T x] = w^T E[x] = w^T \mu$$

$$\text{Var}(w^T x) = w^T \text{Var}(x) w = w^T \Sigma w$$

Üstteki sonuçların boyutlarına dikkat: $w^T \mu$ durumunda $1 \times N \cdot N \times 1 = 1 \times 1$, $w^T \Sigma w$ durumunda ise $1 \times N \cdot N \times N \cdot N \times 1 = 1 \times 1$. İki durumda da tek boyutlu skalar değerler elde ettik. Yani w yönündeki izdüşüm bize tek boyutlu bir çizgi verecektir. Bu sonuç aslında çok sasirtici olmasa gerek, tüm veri noktalarını alıp, başlangıcı başnokta 0,0 (origin) noktasında olan vektörlere çevirip aynı yöne işaret edecek şekilde düzenliyoruz, bu vektörleri tekrar nokta olarak düşünürsek, tabii ki aynı yönü gösteriyorlar, bilahere aynı çizgi üzerindeki noktalara dönüşüyorlar. Aynı çizgi üzerinde olmak ne demek? Tek boyuta inmiş olmak demek.

Bastaki amacımıza dönersek, $\text{Var}(z_1)$ 'i maksimize etmek aynı anda $\text{Var}(w_1^T \Sigma w_1)$ 'i maksimize etmek demektir.

Ufak bir sorun $w_1^T \Sigma w_1$ 'i sürekli daha büyük w_1 'lerle sonsuz kadar büyütebilirsiniz. Bize ek bir kısıtlama şartı daha lazım, bu şart $\|w\| = 1$ olabilir, yani w 'nin norm'u 1'den daha büyük olmasın. Böylece optimizasyon w 'yi sürekli büyüte büyüte maksimizasyon yapmayacak, sadece yön bulmak ile ilgilenecek, iyi, zaten biz w 'nin yönü ile ilgileniyoruz. Aradığımız ifadeyi yazalım, ve ek sınırı Lagrange ifadesi olarak ekleyelim, ve yeni bir L ortaya çıkartalım,

$$L(w_1, \lambda) = w_1^T \Sigma w_1 - \lambda(w_1^T w_1 - 1)$$

Niye eksiden sonraki terim o sekilde eklendi? O terim oyle sekilde secildi ki, $\partial L / \partial \lambda = 0$ alinince $w_1^T w_1 = 1$ geri gelsin / ortaya ciksinsin [2, sf 340], bu Lagrange'in dahice bulusu. Bunu kontrol edebilirsiniz, λ 'ya gore turev alirken w_1 sabit olarak yokolur, parantez icindeki ifadeler kalir ve sifira esitlenince orijinal kisiltlama ifadesi geri gelir. Simdi

$$\max_{w_1} L(w_1, \lambda)$$

Turevi w_1 'e gore alirsak, ve sifira esitlersek,

$$2w_1 \Sigma - 2\lambda w_1 = 0$$

$$2w_1 \Sigma = 2\lambda w_1$$

$$\Sigma w_1 = \lambda w_1$$

Ustteki ifade ozdeger, ozvektor ana formulune benzemiyor mu? Evet. Eger w_1 , Σ 'nin ozvektoru ise ve esitligin sagindaki λ ona tekabul eden ozdeger ise, bu esitlik dogru olacaktir.

Peki hangi ozdeger / ozvektor maksimal degeri verir? Unutmayalim, maksimize etmeye calistigimiz sey $w_1^T \Sigma w_1$ idi

Eger $\Sigma w_1 = \lambda w_1$ yerine koyarsak

$$w_1^T \lambda w_1 = \lambda w_1^T w_1 = \lambda$$

Cunku $w_1^T w_1$ 'nin 1 olacagi sartini koymustuk. Neyse, maksimize etmeye calistigimiz deger λ cikti, o zaman en buyuk λ kullanirsak, en maksimal varyansi elde ederiz, bu da en buyuk ozdegerin ta kendisidir.

Demek ki izdusum yapilacak "yon" kovaryans Σ 'nin en buyuk ozdegerine tekabul eden ozvektor olarak secilirse, temel bileşenlerden en onemlisini hemen bulmus olacagiz.

Cebirin geri kalani w_2, w_3 icin devam eder, bu turetmenin detaylarini [1] ve [2] gibi kaynaklarda bulabilirsiniz. Fakat ulasilan sonuc en bileşenlerin onem sirasinin aynen ozdegerlerin buyukluk sirasina tekabul ediyor olmasi.

Ornek

Simdi tum bunlari bir ornek uzerinde gorelim. Iki boyutlu ornek veriyi ustte yuklemistik. Simdi veriyi "sifirda ortalayacagiz" yani her kolon icin o kolonun ortalama degerini tum kolondan cikartacagiz. PCA ile islem yaparken tum degerlerin sifir merkezli olmasi gerekiyor.

Daha sonra ozdegerlerini, vektorlerini hesaplayabilmek icin verinin kovaryansini hesaplayacagiz.

```
import numpy.linalg as lin
from pandas import *
data = read_csv("testSet.txt", sep="\t", header=None)
print data.shape
print data[:10]
```

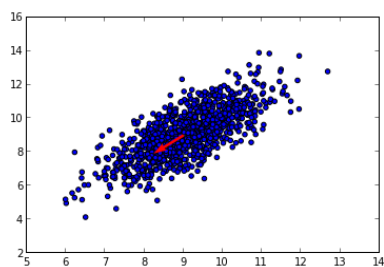
```
means = data.mean()
meanless_data = data - means
cov_mat = np.cov(meanless_data, rowvar=0)
print cov_mat.shape
eigs,eigv = lin.eig(cov_mat)
eig_ind = np.argsort(eigs)
print eig_ind
```

```
(1000, 2)
      0      1
0  10.235186  11.321997
1  10.122339  11.810993
2   9.190236   8.904943
3   9.306371   9.847394
4   8.330131   8.340352
5  10.152785  10.123532
6  10.408540  10.821986
7   9.003615  10.039206
8   9.534872  10.096991
9   9.498181  10.825446
(2, 2)
[0 1]
```

```
print eigs[1],eigv[:,1].T
print eigs[0],eigv[:,0].T
2.89713495618 [-0.52045195 -0.85389096]
0.366513708669 [-0.85389096  0.52045195]
```

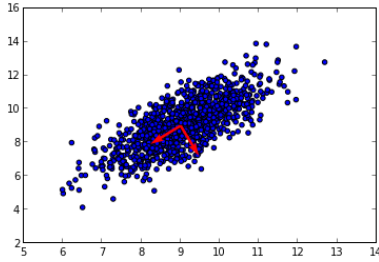
En buyuk olan yonu quiver komutunu kullanarak orijinal veri seti uzerinde gostere-
lim,

```
plt.scatter(data.ix[:,0],data.ix[:,1])
# merkez 9,9, tahminen secildi
plt.quiver(9,9,eigv[1,1],eigv[0,1],scale=10,color='r')
plt.savefig('pca_2.png')
```



Goruldugu gibi bu yon hakikaten dagilimin, veri noktalarinin en cok yayilmis oldugu yon. Demek ki PCA yontemi dogru sonucu buldu. Her iki yonu de ciz-
ersek,

```
plt.scatter(data.ix[:,0],data.ix[:,1])  
plt.quiver(9,9,eigv[1,0],eigv[0,0],scale=10,color='r')  
plt.quiver(9,9,eigv[1,1],eigv[0,1],scale=10,color='r')  
plt.savefig('pca_3.png')
```



Bu ikinci yon birinciye dik olmalidi, ve o da bulundu. Aslinda iki boyut olunca baska secenek kalmiyor, 1. yon sonrasi ikincisi baska bir sey olamazdi, fakat cok daha yuksek boyutlarda en cok yayilimin oldugu ikinci yon de dogru sekilde geri getirilecekti.

SVD ile PCA Hesaplamak

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

Ayrica birazdan gorecegimiz gibi SVD, kovaryans matrisi uzerinde degil, A'nin kendisi uzerinde isletilir, bu hem kovaryans hesaplama asamasini atlamamizi, hem de kovaryans hesabi sirasinda ortaya cikabilecek numerik puruzlerden ko-
runmamizi saglar (cok ufak değerlerin kovaryans hesabini bozabilecegi literaturde bahsedilmektedir).

PCA ve SVD baglantisina gelelim:

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 uzerinde değerler var ve bu σ_j değerleri A'nin tekil değerleri (singular values) olarak biliniyor.

Simdi A yerine AA^T koyalım, yani A 'nin kovaryans matrisinin SVD ayrıştırmasını yapalım, acaba elimize ne geçecek?

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

S bir köşegen matrisi, o zaman SS^T matrisi de köşegen, tek farkla köşegen üzerinde artık σ_j^2 değerleri var. Bu normal.

SS^T yerine Λ sembolünü kullanalım, ve denklemi iki taraftan (ve sağdan) U ile çarparsak (unutmayalım U ortanormal bir matris ve $U^T U = I$),

$$\begin{aligned} AA^T U &= U \Lambda U^T U \\ AA^T U &= U \Lambda \end{aligned}$$

Son ifadeye yakından bakalım, U 'nun tek bir kolonuna, u_k diyelim, odaklanacak olursak, üstteki ifadede bu sadece kolona yönelik nasıl bir eşitlik çıkartabilirdik? Soyle çıkartabilirdik,

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

Bu ifade tanıdık geliyor mu? Özdeğer / özvektor klasik yapısına eristik. Üstteki eşitlik sadece ve sadece eğer u_k , AA^T 'nin özvektörü ve σ^2 onun özdeğeri ise geçerlidir. Bu eşitliği tüm U kolonları için uygulayabileceğimize göre demek ki U 'nun kolonlarında AA^T 'nin özvektörleri vardır, ve AA^T 'nin özdeğerleri A 'nin tekil değerlerinin karesidir.

Bu muthis bir buluş. Demek ki AA^T 'nin özvektörlerini hesaplamak için A üzerinde SVD uygulayarak U 'yu bulmak ise yarar, kovaryans matrisini hesaplamak gerekli değil (bir hesap yerine otekini koymuş olduk $A^T A$ çarpımı yerine yerine AA^T hesabı). AA^T özdeğerleri üzerinde büyüklük karşılaştırması için ise A 'nin tekil değerlerine bakmak yeterli!

Örnek

İlk bölümdeki örneğe donelim, ve özvektörleri SVD üzerinden hesaplatalım.

```
U,s,Vt = svd(meanless_data.T,full_matrices=False)
print U
```

```
[[-0.52045195 -0.85389096]
 [-0.85389096  0.52045195]]
```

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

```
[[ 1.00000000e+00  3.70255042e-17]
 [ 3.70255042e-17  1.00000000e+00]]
```

Goruldugu gibi ayni ozvektorleri bulduk.

New York Times Yazıları Analizi

Simdi daha ilginç bir ornege bakalım. Bir araştırmacı belli yıllar arasındaki NY Times makalelerinde her yazıda hangi kelimenin kaç kere çıktığının verisini toplamış [1,2,3], bu veri 4000 kusur kelime, her satır (yazı) için bir boyut (kolon) olarak kaydedilmiştir. Bu veri nytimes.csv üzerinde ek bir normalize işleminden sonra, 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 kusur boyuttan sonra bu hala muhtemelen bir kazanım) geri kalanlar üzerinde mesela bir kümeleme algoritması kullanabiliriz.

Ana veriyi yükleyip birkaç satırını ve kolonlarını göstereyim.

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

	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

Yüklemeyi yapıp sadece etiketleri aldık ve onları bir kenara koyduk. Şimdi önemli bir normalizasyon işlemi gerekiyor - ki bu işleme ters doküman-frekans ağırlıklandırması (inverse document-frequency weighting -IDF-) ismi veriliyor - her dokümanda aşırı fazla ortaya çıkan kelimelerin önemi özellikle azaltılıyor, ki diğer kelimelerin etkisi artabilsin.

IDF kodlamasi alttaki gibidir. Önce `class.labels` kolonunu atarız. 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 çekelim, bunu toplamın sıfır olma durumu için yapıyoruz, ki sonra sıfırla bölürken NaN sonucundan kaçınalım.

```
nyt2 = nyt.drop('class.labels', axis=1)
freq = nyt2.astype(bool).sum(axis=0)
freq = freq.replace(0,1)
w = np.log(float(nyt2.shape[0])/freq)
nyt2 = nyt2.apply(lambda x: x*w, axis=1)
nyt2 = nyt2.apply(lambda x: x / np.sqrt(np.sum(np.square(x))+1e-16), axis=1)
#nyt2 = nyt2.div(nyt2.sum(axis=0), axis=1)
nyt2=nyt2.ix[:,1:] # ilk kolonu atladi
print nyt2.ix[:8,102:107]
```

	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

Not: Bir diger normalize metodu

```
import pandas as pd

df = pd.DataFrame([[1.,1.,np.nan],
                  [1.,2.,0.],
                  [1.,3.,np.nan]])

print df
print df.div(df.sum(axis=0), axis=1)
```

	0	1	2
0	1	1	NaN
1	1	2	0
2	1	3	NaN

	0	1	2
0	0.333333	0.166667	NaN
1	0.333333	0.333333	NaN
2	0.333333	0.500000	NaN

SVD yapalım


```

nyt3 = nyt2 - nyt2.mean(0)
u,s,v = lin.svd(nyt3.T,full_matrices=False)
print s[:10]

[ 1.41676764  1.37161893  1.31840061  1.24567955  1.20596873  1.18624932
  1.15118771  1.13820504  1.1138296   1.10424634]

print u.shape

(4430, 102)

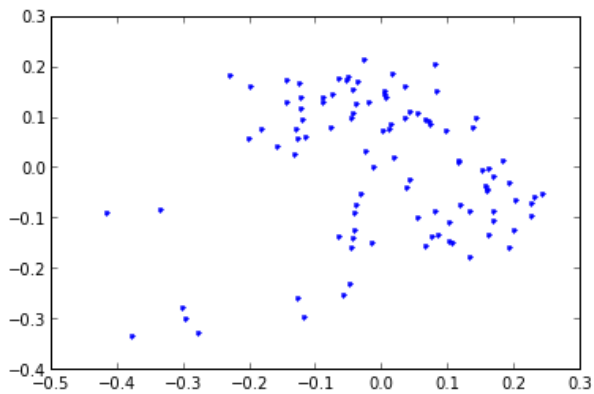
```

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

```

proj = np.dot(nyt, u[:, :2])
proj.shape
plt.plot(proj[:,0],proj[:,1],'.')
plt.savefig('pca_4.png')

```

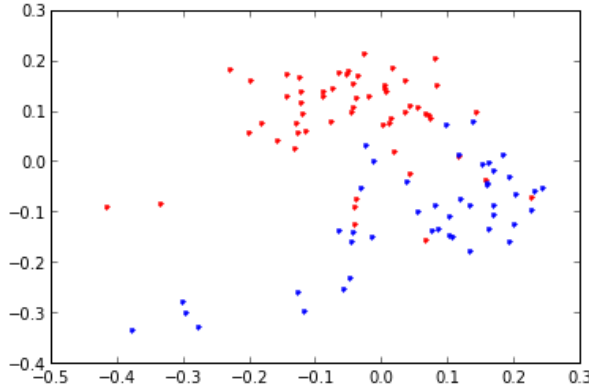


Simdi ayni veriyi bir de etiket bilgisini devreye sokarak cizdirelim. Sanat kirmizi muzik mavi olacak.

```

arts =proj[labels == 'art']
music =proj[labels == 'music']
plt.plot(arts[:,0],arts[:,1],'.r')
plt.plot(music[:,0],music[:,1],'.b')
plt.savefig('pca_5.png')

```



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 uzerinden net bir ayirim ortaya cikartabiliyoruz. Bu PCA yonteminin iyi bir is becerdigini gosteriyor, ve kelime sayilarinin makalelerin icerigi hakkında ipucu icerdigini ispatliyor.

Kaynaklar

[1] Alpaydin, E., Introduction to Machine Learning, 2nd Edition

[2] Strang, G., Linear Algebra and Its Applications, 4th Edition

[3] <http://www.stat.columbia.edu/~fwood/Teaching/w4315/Spring2010/PCA/slides.pdf>

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

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

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

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

[8] Lineer Cebir notlarimizda SVD turetilmesine bakinca ozdeger/vektor mantigina atif yapildigini gorebiliriz ve akla su gelebilir; "ozdeger / vektor rutini isletmekten kurtulalim dedik, SVD yapiyoruz, ama onun icinde de ozdeger/vektor hesabi var". Fakat sunu belirtmek gerekir ki SVD numerik hesabini yapmanin tek yontemi ozdeger/vektor yontemi degildir. Mesela Numpy Linalg kutuphanesi icindeki SVD, LAPACK dgesdd rutinini kullanir ve bu rutin ic kodlamasinda QR, ve bir tur bol / istila et (divide and conquer) algoritmasi isletmektedir.