Uses of PCA

Lyle Ungar

Learning objectives

PCA for feature creation PCR, Semi-supervised learning PCA for visualization Eigenfaces: see the worksheet Eigenwords: word embeddings

PCR: Principal Component Regression

Do a PCA on X to get scores Z and loadings V
 Do OLS regression using Z as features

 $\mathbf{y} = \mathbf{w} \cdot \mathbf{z}$ $\mathbf{w} = (\mathbf{Z}^{\mathsf{T}}\mathbf{Z})^{-1}\mathbf{Z}^{\mathsf{T}}\mathbf{Y}$

For future predictions, use $z = V^T x$ $y = w \cdot V^T x$

PCR

♦ How to find z for a new x?
X = ZV^T
x^T V = z^T V^TV = z^T
z = V^Tx

Semi-supervised learning

Use unlabeled data X₀ to derive new features z = φ(x)
Train y = f(φ(x), w)



Semi-supervised PCR 1a. Do a PCA on a big X_u to get loadings V 1b. Project X (with labels y) to get scores Z 2. Do OLS regression using Z as features $y = w \cdot z$ $w = (Z^TZ)^{-1}Z^TY$

> For future predictions, use $z = V^T x$ $y = w \cdot V^T x$

Semi-Supervised Learning

 Hypothesis: P(c|x) can be more accurately computed using shared structure with P(x)



from Socher and Manning

Semi-Supervised Learning

 Hypothesis: P(c|x) can be more accurately computed using shared structure with P(x)



from Socher and Manning

PCA for visualization

Project original observations, x, into 2 dimensions

- PCA minimizes reconstruction error
- This is one of many ways of trying to make points that were close in *m* dimensions still be close in 2 dimensions

Look at the loadings

• Often prefer sparse loadings

Country well-being

```
# import OECD data
df_OECD = pd.read_csv('http://www.cis.upenn.edu/~cis545/OECD_well-being.csv')
df_OECD
X, country = df OECD.iloc[:,1:].values, df OECD.iloc[:,0].values
```

sc = StandardScaler()
X std= sc.fit transform(X)

```
def plot_points(X, labels):
    plt.scatter(X[:,0],X[:,1])
    for i, txt in enumerate(labels):
        plt.annotate(txt, (X[i,0],X[i,1]))
    plt.show()
```

```
# plot PCA
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(X_std)
X_PCA = pca.transform(X_std)
plot_points(X_PCA,country)
```



OCED PCA loadings

- **Dwellings without basic facilities**
- Housing expenditure
- Rooms per person
- Household net adjusted disposable income Household net financial wealth
- Δ
- Labour market insecurity 5
- **Employment rate** 6
- Long-term unemployment rate
- Personal earnings 8
- 9 Quality of support network 10 Educational attainment
- Student skills 11
- 12 Years in education
- 13 Air pollution
- 14 Water quality15 Stakeholder engagement for developing regulations
- 16 Voter turnout
- Life expectancy 17
- 18 Self-reported health
- 19 Life satisfaction
- 20 Feeling safe walking alone at night
- Homicide rate 21
- 22 Employees working very long hours23 Time devoted to leisure and personal care

OCED PCA loadings

- Dwellings without basic facilities
- Housing expenditure
- 2
- Rooms per person Household net adjusted disposable income Household net financial wealth 3
- 4
- Labour market insecurity 5
- **Employment rate** 6
- Long-term unemployment rate 7
- Personal earnings 8
- 9 Quality of support network 10 Educational attainment
- Student skills 11
- Years in education 12
- 13 Air pollution
- 14 Water quality
- Stakeholder engagement for developing regulations 15
- 16 Voter turnout
- Life expectancy 17
- **18 Self-reported health** 19 Life satisfaction
- Feeling safe walking alone at night 20
- 21 Homicide rate
- 22 Employees working very long hours23 Time devoted to leisure and personal care

All the variance is in the first PC

Eigenwords

Learning objectives Distributional similarity Word embeddings SVD on words

Lyle Ungar University of Pennsylvania

Represent each word by its context

	l ate ham You ate c You ate	hees	e							
		rd Be	fore		Word	Word After				
		ate	chees	e har	n I You	ate c	heese	e han	n I You	
	ate	0	0	0	12	0	1	1	00	
	cheese	1	0	0	00	0	0	0	0 0	
word	ham	1	0	0	0 0	0	0	0	0 0	
		0	0	0	00	1	0	0	0 0	
	You	0	0	0	0 0	2	0	0	0 0	

Distributional Similarity

Hypothesis: Words with similar contexts have similar meanings

Eigenwords:*Word embeddings*

 Project high dimensional context to low dimensional space (SVD/PCA)

Similar words are close in this low dimensional space

I ate ham You ate c You ate	hees	e									
	Word Before ate cheese ham I You					Word After ate cheese ham I Yo					
ate	0	0	0	1	2	0	1	1	0 0		
cheese	1	0	0	0	0	0	0	0	0 0		
ham	1	0	0	0	0	0	0	0	0 0		
	0	0	0	0	0	1	0	0	0 0		
You	0	0	0	0	0	2	0	0	00		

Eigenwords as SVD

- Left singular vectors are eigenwords
 - a vector representing each word "word embeddings"

Right singular vectors times context give eigentokens

- vectors mapping contexts to the latent space
 - I ate ham You ate cheese You ate

	W	ord Be	efore		Word After				
	ate o	cheese	e ham		You	ate cl	neese	e han	n I You
ate	0	0	0	1	2	0	1	1	0 0
cheese	1	0	0	0	0	0	0	0	0 0
ham	1	0	0	0	0	0	0	0	00
	0	0	0	0	0	1	0	0	0 0
You	0	0	0	0	0	2	0	0	0 0

Similar words are close

Nouns and verbs

Pronouns

Numbers

PC 1

Names

Word Sense Disambiguation

 Estimate "state vector" ("contextualized embedding") for a word using right singular vectors

• Similar meanings will again be close.

- The ships dock in the port.
- The port is loaded onto ships and sent to America
- The meat is tender.
- I have tender feelings for her.
- The company will tender an offer.

Use eigenwords/eigentokens in supervised learning

Similar' words have embeddings that are close

- Predict labels for tokens based on their estimated "state vector"
 - Part of speech
 - Named entity type (person, place, thing...)
 - Word sense ("meaning") disambiguation
- Or embed sentences

Word2vec

- Word embeddings, often found by deep learning, are very popular now
 - Word2Vec has similar performance to the simpler eigenwords
- Deep learning (contextualized) versions such as BERT and friends work better
 - To be covered later

What you should know

PCR

• Use principal components from training to find scores on test data

Interpretation: Scores and loadings

• percentage of variance explained

Word embeddings

- Context free (left singular vectors: words)
- Context sensitive (right singular vectors: context)