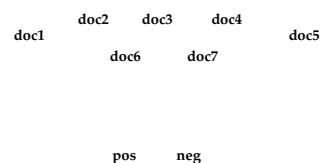


# From Latent Features to Word Embeddings

Lecture 11: Text Analytics for Big Data  
Mark Keane, Insight/CSI, UCD

## Intuitive Idea



Imagine Model-X that maps word features of  
docs to sentiment features

doc1 doc2 doc3 doc4 doc5  
doc6 doc7

happy	nice	sad	pooh	tree	evil
			pos	neg	

**Model-X** tells if docs are positive or negative;  
based on the word-feature in the docs...

doc1	doc2	doc3	doc4	doc5
doc6	doc7			

happy	nice	sad	pooh	tree	evil
			pos	neg	

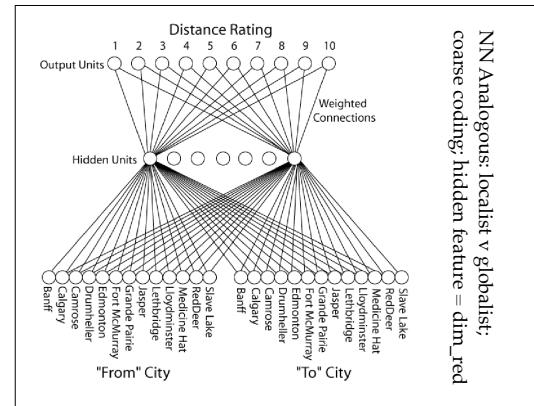
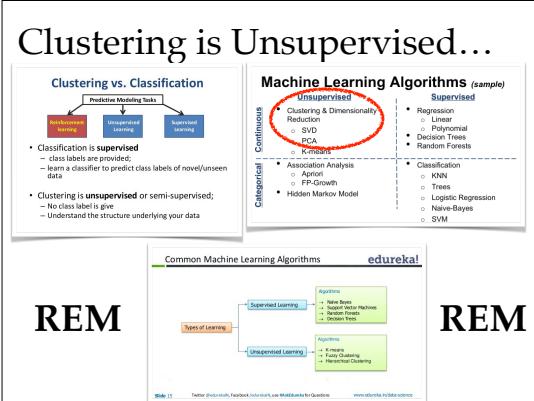
But, **Model-X** does not tell us anything deeper about these feature-dimensions; latent/hidden features

feel-good			feel-bad	
happy	nice	sad	pooh	tree
				evil
			pos	neg

Latent features like a feel-good or a feel-bad feature...

feel-good			feel-bad	
happy	nice	sad	pooh	tree
				evil
			pos	neg

Methods exist to find such latent features;  
basically, factor analytic techniques that are unsupervised ML



## Buzz Words & Terminology

◆ Discovering Latent/Hidden Features, Factors

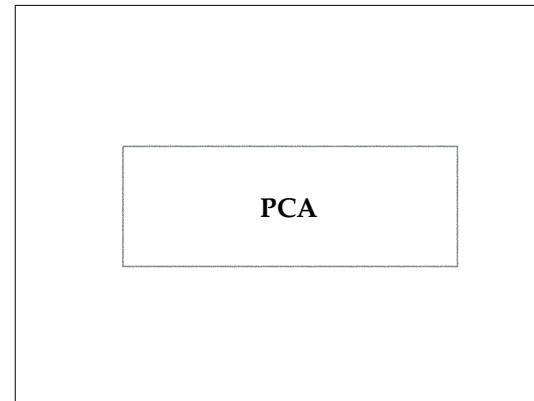
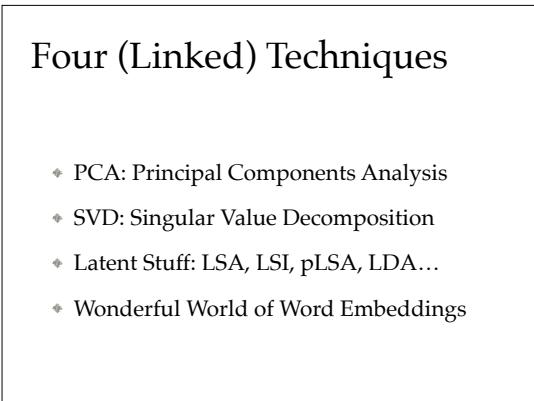
◆ Dimensionality Reduction

◆ Word embeddings and Deep Learning

feature = dimension = variable

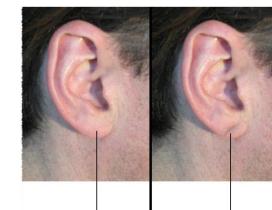
factor is linear combination of variables

factor interchangeable with component



## PCA: The Idea:

Several low-level features may combine (as a factor / latent feature) to predict some outcome (e.g., smoking, ear-lobe folds, high-BMI may represent genetic sub-group that predicts early-stroke risk)



feature = dimension = variable

factor is linear combination of variables

## PCA: The Idea:

- In multivariate analysis, many *variables* (ie word-features) may be correlated (i.e., feel-good words) in their effect on some outcome/prediction (good/bad)
- PCA is a dimension-reduction tool that reduces a large set of *variables* to a smaller set of *factors* that still contains most of the information in the large set
- PCA mathematically transforms a (larger) number of (possibly correlated *variables*) into a (smaller) number of uncorrelated *factors* aka Principal Components

<ftp://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf>

doc1 doc2 doc3 doc4 doc5  
doc6 doc7

feel-good feel-bad tree

happy	nice	sad	pooh	evil	tree
			pos	neg	

feel-good words will be correlated in their effects on the predicted outcome (pos), and feel-bad words in their effects on predicted outcome (neg); BUT the effects of feel-good and feel-bad words will be uncorrelated (because they are a different class of variables). Components/Factors = FEEL-GOOD/FEEL-BAD

## PCA: The Idea: II

- 1st Principal-Component accounts for most of the variance, the second the next next most and so on
- Components summarise the total variance of data
- NB: NO guarantee Components are interpretable
- Like Factor Analysis but actually different...

<ftp://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf>



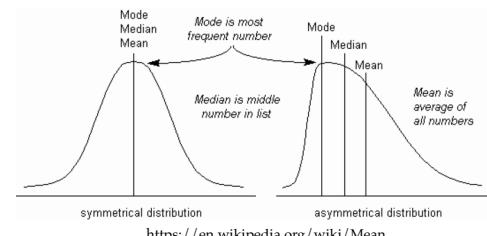
## PCA: Some Stats Nuts & Bolts

- Means & Deviations from it
- Standard Deviations & Variance
- Co-variance & Correlation
- The Matrix Step
- Eigenvectors & Eigenvalues

<ftp://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf>

## I Really Means It !

Mean (arithmetic) is the central value of a discrete set of numbers  
Got by the sum of the values divided by the number of values



<https://en.wikipedia.org/wiki/Mean>





## Oh...and...eigen-whaaaa?

Will all be done using matrix multiplication  
and the like...  
finding factors to “explain” the variance in  
these matrices  
...matrices produced by multiplying original  
by transposed ones...

$$A^T \cdot A \text{ and } A \cdot A^T$$

## Eigenvectors & Eigenvalues

- Given an input matrix  $X$ , an eigenvector of the matrix is a non-zero vector  $v$  that satisfies the equation:

$$Xv = \lambda v$$

where the corresponding number  $\lambda$  is called an eigenvalue.

- Eigendecomposition is the factorisation of a matrix into its eigenvalues and eigenvectors.

$$X = \begin{pmatrix} 1.0000 & 0.5000 & 0.3330 & 0.2500 \\ 0.5000 & 1.0000 & 0.6667 & 0.5000 \\ 0.3333 & 0.6667 & 1.0000 & 0.7500 \\ 0.2500 & 0.5000 & 0.7500 & 1.0000 \end{pmatrix}$$

$$A = \begin{pmatrix} 0.5390 & 0 & 0 & 0 \\ 0 & 0.8483 & 0 & 0 \\ 0 & 0 & 0.4078 & 0 \\ 0 & 0 & 0 & 0.2077 \end{pmatrix}$$

4 x 4 symmetric

Eigenvalues and values exist in pairs:  
every eigenvector has a corresponding  
eigenvalue.

$$V = \begin{pmatrix} 0.3477 & -0.8105 & -0.1451 & 0.0698 \\ -0.3522 & 0.3822 & -0.1876 & 0.7189 \\ -0.5619 & 0.3039 & 0.0487 & -0.7026 \\ -0.5089 & 0.4661 & -0.5015 & 0.5216 \end{pmatrix}$$

4

From D. Greene's course...

- Each eigenvector has a direction - i.e. it is a dimension.
- Eigenvectors of symmetric matrices are orthogonal to each other - i.e. they point in completely different directions.
- Each eigenvalue is a number indicating how much variance there is in the data in that direction.
- The eigenvector with the largest eigenvalue has the most variance, making it a useful dimension for projection, and so on...
- Principal Components (PCs): New dimensions constructed as linear combinations of the original features, which are uncorrelated with one another. Constructed from eigenvectors.

**DIVERSION  
END**

## PCA: The Idea: REM

- In multivariate analysis, many *variables* (ie word-features) may be correlated (i.e., feel-good words) in their effect on some outcome / prediction (good/bad)
- PCA is a dimension-reduction tool that reduces a large set of *variables* to a smaller set of *factors* that still contains most of the information in the large set
- PCA mathematically transforms a (larger) number of (possibly) correlated *variables* into a (smaller) number of uncorrelated *factors* aka Principal Components

<ftp://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf>

doc1 doc2 doc3 doc4 doc5  
doc6 doc7  
**REM**

feel-good feel-bad tree

happy	nice	sad	pooh	evil	tree
			<b>pos</b>	<b>neg</b>	

**feel-good words** will be correlated in their effects on the predicted outcome (**pos**), and **feel-bad words** in their effects on predicted outcome (**neg**); BUT the effects of **feel-good** and **feel-bad words** will be uncorrelated (because they are a different class of variables). Components/factors = FEEL-GOOD/FEEL-BAD

## PCA: The Idea: II REM

- 1st Principal-Component accounts for most of the variance, the 2nd the next most and so on
- Components summarise total variance of data
- No guarantee Components are interpretable
- Components reduce dimensionality of data

[ftp://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf](http://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf)

A principal component analysis can be considered as a rotation of the axes of the original variables onto axes parallel to new orthogonal axes, called principal axes, such that the new axes coincide with the maximum variance directions of the original observations. Consider the line or axis passing through the center of the data set of points in Figure 1. Project the original data points onto this line. Let  $y_m$  be the perpendicular distance of the point  $(x_m, y_m)$  from the axis defined by the direction  $Y_1$ . This axis has the property that the sum of the squared projected points  $y_m^2, m = 1, \dots, n$ , is greater than the variance of the points when projected onto any other axis passing through  $(\bar{x}_1, \bar{y}_1)$ . Any line parallel to  $Y_1$  also has the same property. It is however convenient geometrically to use the first representation.

The property of maximum variation of the projected points defines the first principal axis; it is the line or direction with maximum variance of the projected values of the original data points. The projected values correspond to the *principal component scores*. The first principal axis is often called the *axis of best fit* and the square of the SD of the perpendicular deviations of the original data points from the line is a minimum. Some properties of the principal axis are determined with the property that they are orthogonal to the previous principal axis and that they minimize the variation of the projected points subject to these constraints.

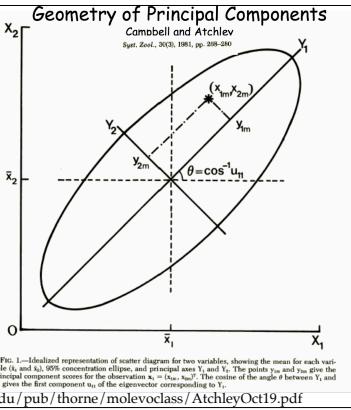


Fig. 1.—Idealized representation of scatter diagram for two variables, showing the mean for each variable ( $\bar{x}_1$  and  $\bar{x}_2$ ), 95% concentration ellipse, and principal axes  $Y_1$  and  $Y_2$ . The points  $y_{1m}$  and  $y_{2m}$  give the perpendicular distances of the point  $(x_m, y_m)$  from the axis  $Y_1$ . The angle  $\theta = \cos^{-1} l_{11}$  is the angle between  $Y_1$  and  $X_1$  given the first component  $u_{11}$  of the eigenvector corresponding to  $Y_1$ .

## PCA: The Maths Bit

PCA is performed on a square, symmetric matrix:

- ♦ **Covariance Matrix:** based on comparing the variance between two variables
- ♦ **Correlation Matrix:** based on correlating two variables (e.g., Pearson's correlation)
- ♦ NB: these are related to one another: CorrMatrix is scaled version of CovMatrix

[ftp://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf](http://statgen.ncsu.edu/pub/thorne/molevoclass/AtchleyOct19.pdf)

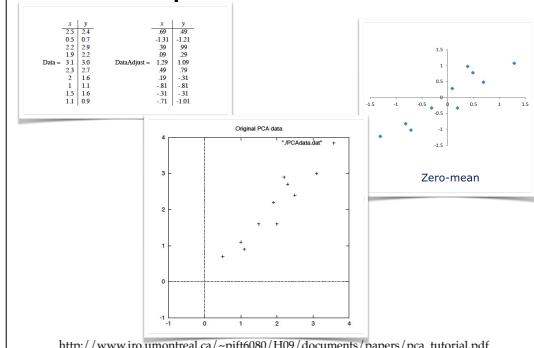
## PCA: Steps...

- ♦ Step1: Subtract means from raw data
- ♦ Step2: Do COV / CORR matrix for variables
- ♦ Step3: Find matrix eigenvectors+eigenvalues
- ♦ Step4: Order eigenvector(s) from high->low
- ♦ Step5\*: Apply & Re-plot data

\* Optional; eigenvectors adopted as new representation of data

[http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca\\_tutorial.pdf](http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca_tutorial.pdf)

## PCA: Step 1: Subtract Means



[http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca\\_tutorial.pdf](http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca_tutorial.pdf)

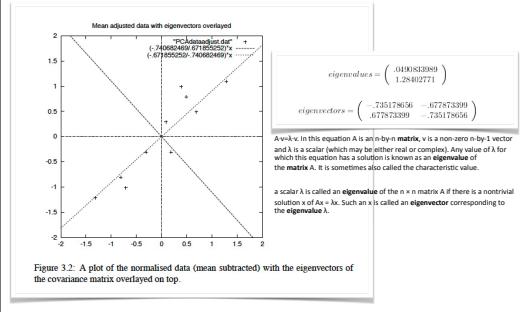
## PCA: Step 2: Covariance Matrix

$$\text{cov} = \begin{pmatrix} .616555556 & .615444444 \\ .615444444 & .716555556 \end{pmatrix}$$

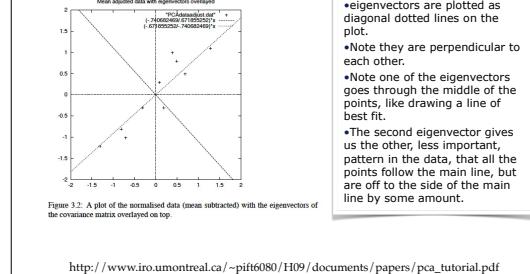
- As the non-diagonal elements in this covariance matrix are positive, we should expect that both the x and y variable increase together
- Aka positively, not negatively correlated

[http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca\\_tutorial.pdf](http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca_tutorial.pdf)

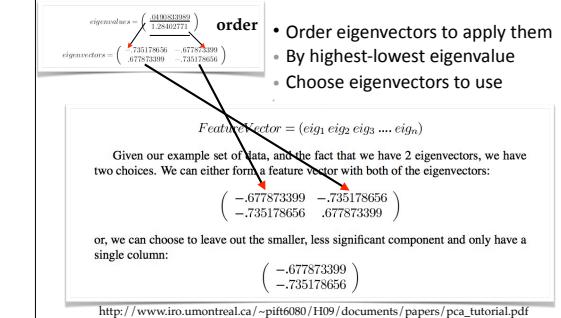
## PCA: Step 3: Get Eigen...Stuff



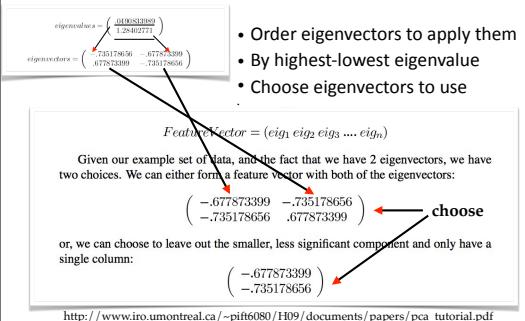
## PCA: Step 3: Get Eigen...Stuff



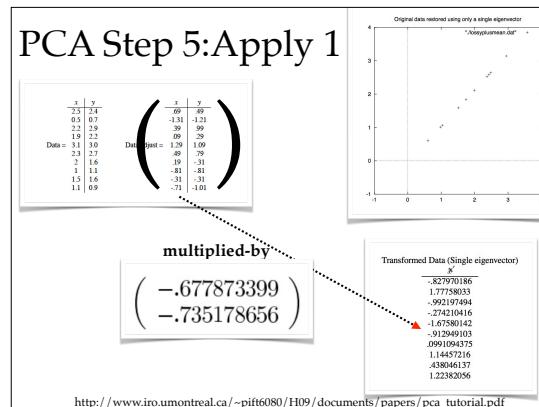
## PCA: Step 4: Order



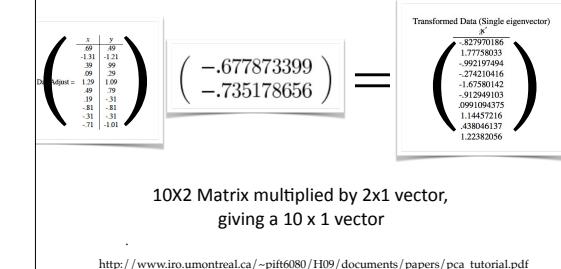
## PCA: Step 4: Order

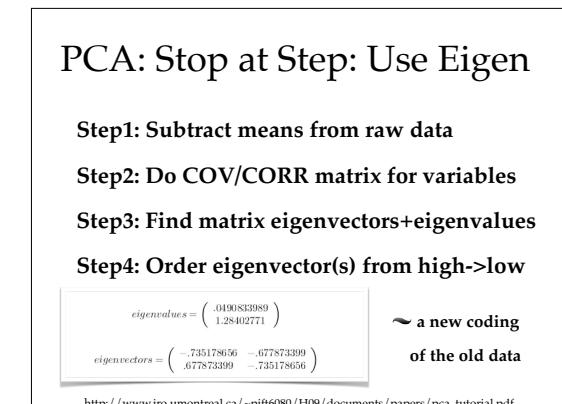
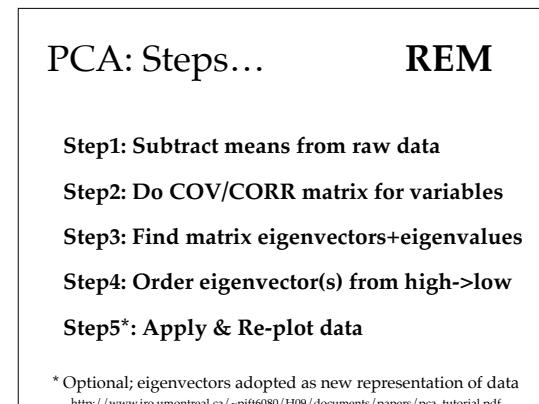
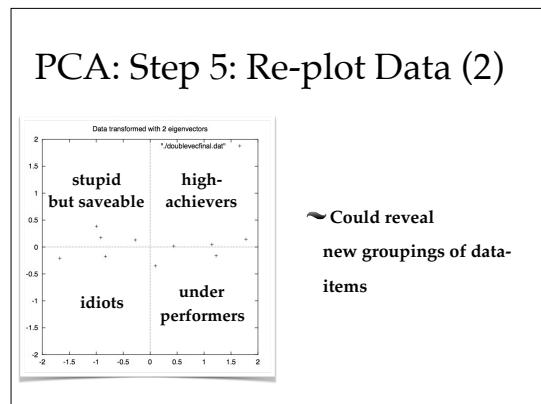
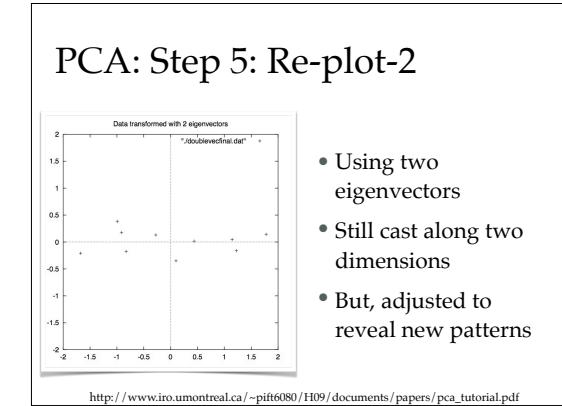
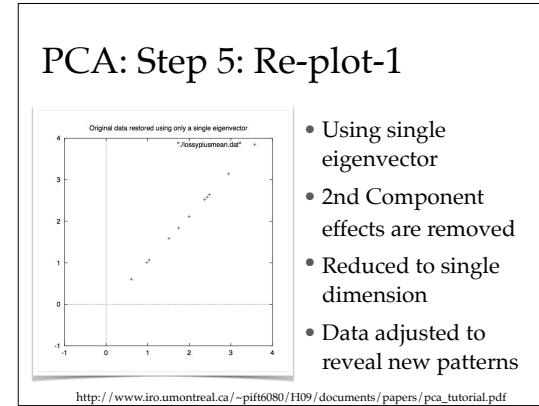
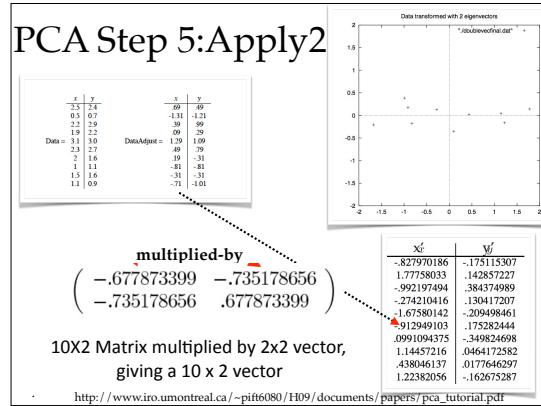


## PCA Step 5: Apply 1



## PCA Step 5: Apply 1





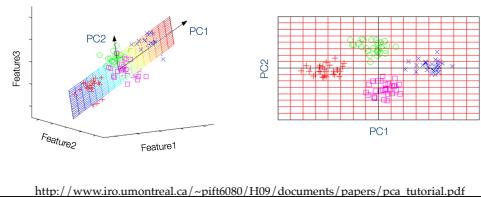
## PCA: Uses of ...

- PCA can be used as a method for feature extraction; finding a set of features to use in your model
- PCA can speed up an ML algorithm by reducing the number of input features/variables
- PCA can aid visualisation; reducing 4-D to 2-D
- Like Factor Analysis but actually different; most think they are the same

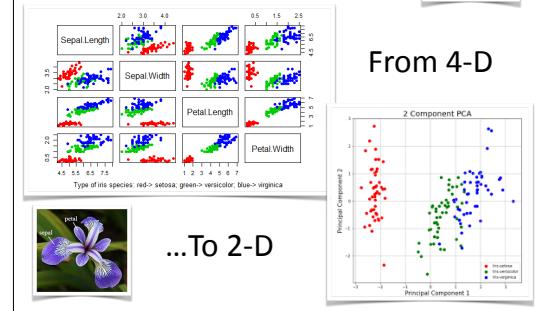
[http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca\\_tutorial.pdf](http://www.iro.umontreal.ca/~pift6080/H09/documents/papers/pca_tutorial.pdf)

## PCA: Prog Example

- PCA can aid visualisation; 4-D to 2-D
- Prog Eg using Irises Database



## PCA: Prog Example

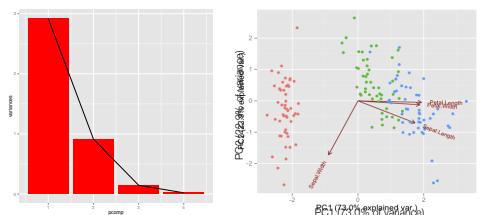


## PCA: Prog Example

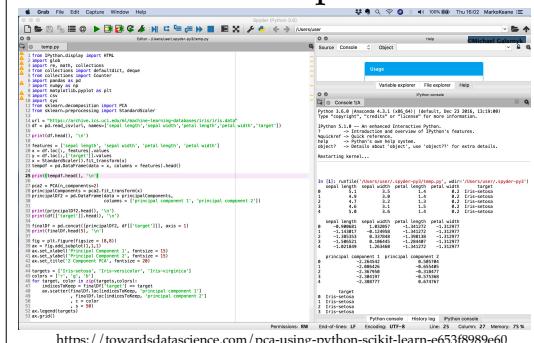
We generally select the  $k$  PCs with the highest variance.

**Example:** Apply PCA to *iris* data set, and examine amount of variance in each PC.

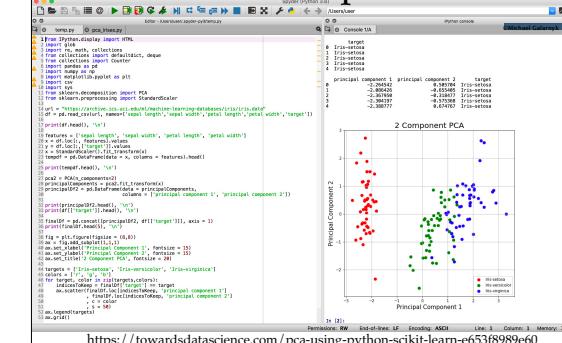
The first two PCs account for ~96% of the variance in the data.



## PCA: Irises Example



## PCA: Irises Example





## SVD: Maths

- SVD decomposes a matrix into constituent parts (factorisation)
- $M$  is the  $m \times n$  matrix we want to decompose
- $U$  is  $m \times m$  matrix
- $\Sigma$  (sigma) is  $m \times n$  diagonal matrix
- $V^*$  is  $n \times n$  matrix

[https://en.wikipedia.org/wiki/Singular\\_value\\_decomposition](https://en.wikipedia.org/wiki/Singular_value_decomposition)

$$M_{m \times n} = U_{m \times m} \Sigma_{m \times n} V^*_{n \times n}$$

$$U_{m \times m} U^*_{m \times m} = I_m$$

$$V^*_{n \times n} V_{n \times n} = I_n$$

## SVD: Maths

- Any matrix can be approximated by one of a lower rank; rank is akin to no of dimensions/factors required to account for data
- $U$  and  $V$  are computed from multiplying original matrix by its transpositions,  $M \cdot M^T$  and  $M^T \cdot M$ , respectively; then getting eigenvectors and eigenvalues of the resultant matrices
- Singular values (in  $\Sigma$ ) are square roots of the eigenvalues of the cross-product matrices ( $U, V$ )
- Singular values can be applied to column/row variables to determine relative “importance” identifying factors (i.e., how much they account for variance)

[http://web.mit.edu/be.400/www/SVD/Singular\\_Value\\_Decomposition.htm](http://web.mit.edu/be.400/www/SVD/Singular_Value_Decomposition.htm)

## SVD: Maths:

$U$  is computed from  $A \cdot A^T$  to get eigenvectors and eigenvalues of resultant matrices

$$A = \begin{bmatrix} 2 & 4 \\ 1 & 3 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$A \cdot A^T = \begin{bmatrix} 2 & 4 \\ 1 & 3 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 2 & 1 & 0 & 0 \\ 1 & 3 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 20 & 14 & 0 & 0 \\ 14 & 10 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} = W$$

$$\Sigma = \begin{bmatrix} 5.47 & 0 \\ 0 & 0.37 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$U = \begin{bmatrix} 0.82 & -0.58 & 0 & 0 \\ 0.58 & 0.82 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

nb: eigen.step not shown

$S$  is the square root of the eigenvalues from  $A \cdot A^T$  or  $A^T \cdot A$

## SVD: Maths:

$$A_{n \times p} = U_{n \times n} S_{n \times p} V^T_{p \times p}$$

$V$  is computed from  $A^T \cdot A$  to get eigenvectors and eigenvalues of resultant matrices

$$A = \begin{bmatrix} 2 & 4 \\ 1 & 3 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$A^T \cdot A = \begin{bmatrix} 2 & 1 & 0 & 0 \\ 4 & 3 & 0 & 0 \end{bmatrix} \begin{bmatrix} 2 & 4 \\ 1 & 3 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 4 \\ 1 & 3 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$V = \begin{bmatrix} 0.40 & -0.91 \\ 0.91 & 0.40 \end{bmatrix}$$

$$S = \begin{bmatrix} 5.47 & 0 \\ 0 & 0.37 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

nb: eigen.step not shown

$S$  is the square root of the eigenvalues from  $A \cdot A^T$  or  $A^T \cdot A$

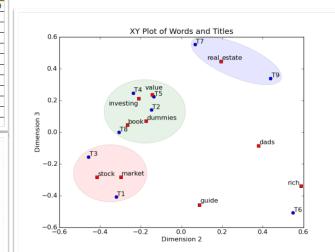
## SVD: Maths: $\Sigma$ : Singular Values

- Square of the singular values are eigenvalues
- Convey relative importance of different eigenvectors

$$\Sigma = \begin{bmatrix} 5.47 & 0 \\ 0 & 0.37 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

## SVD: Eg...see later

Index Words	T1	T2	T3	T4	T5	T6	T7	T8	T9
stock	1	1	1	1	1	1	1	1	1
book	1	1	1	1	1	1	1	1	1
dads	1	1	1	1	1	1	1	1	1
dummies	1	1	1	1	1	1	1	1	1
estate	1	1	1	1	1	1	1	1	1
guide	1	1	1	1	1	1	1	1	1
market	1	1	1	1	1	1	1	1	1
rich	1	1	1	1	1	1	1	1	1
real	1	1	1	1	1	1	1	1	1
value	1	1	1	1	1	1	1	1	1



## Latent Semantic Indexing/Analysis

### Lets Distinguish...

- ◆ LSI (LSA) as a technique for processing term-doc matrices
- ◆ LSI (Latent Semantic Indexing) was the version of this used to build new indices for retrieval systems, typically
- ◆ LSA (Latent Semantic Analysis) was the version of this used to process terms-docs (e.g., for paraphrase id)
- ◆ LSA-Theory (Plato's Problem) was an attempted to build a new LSA-based-vector-representation for words/paras to use for semantic tasks (Landauer & Dumais, 1997)

### LSI/LSA = SVD(Term-Doc Matrix)

- ◆ LSI (LSA) is SVD on a Term-Doc VSM to find latent/hidden associations between words, usually using co-occurrence
- ◆ If two words occur together frequently then they are related to some common topic (meaning by the company it keeps)
- ◆ 'goal', 'net', 'post' will have associated clusters of words about the soccer topic, even though "soccer" may not be explicitly mentioned, the latent-factor could be (interpreted as) SOCCER
- ◆ LSI uses a weighted term-document matrix, performing a SVD on the matrix, and using the resultant matrix to identify the concept-groupings contained in the text.

### LSI/A is based on VSMs

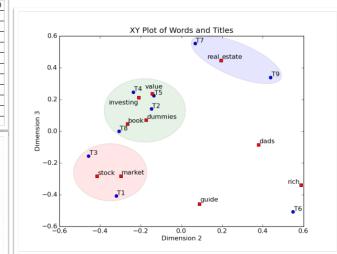
- ◆ LSI/A begins by constructing a term-document matrix, A, to identify the occurrences of the  $m$  unique terms within a collection of  $n$  documents.
- ◆ In a term-document matrix, each term/word is a row, and each document/text-item is a column
- ◆ This matrix is usually very large and very sparse; aka something you want to reduce
- ◆ SVD can give you a new matrix to encode terms/docs

### LSI/LSA

- ◆ Gives you a new set of vectors for the terms and documents; re-coding them in terms of the factors (eigenvector and eigenvalues found)
- ◆ These new vectors can be clustered; to show new, deeper groupings or combine for D-reduction
- ◆ Also, new recoding of words may constitute a latent model of meaning (LSA-Theory)

### From This to This...

Index Words	T1	T2	T3	T4	T5	T6	T7	T8	T9
book	0.15	-0.27	0.04						
dads	0.24	0.38	-0.09						
dummies	0.13	-0.17	0.07						
estate	0.18	0.19	0.45						
guide	0.22	0.09	-0.46						
investing	0.74	-0.21	0.21						
market	0.18	-0.3	-0.28						
real	0.18	0.19	0.45						
rich	0.36	0.59	-0.34						
stock	0.25	-0.42	-0.28						
value	0.12	-0.14	0.23						









## LSA Theory: Bang !

300-dimension word-vectors tells you an awful lot about a word's meaning, it predicts:

- ❖ Associations: night->day (and reaction times)
- ❖ Synonyms: punch -> strike
- ❖ Measure of word and sentence coherence:
  - ❖ "football, score, keeper" V "football, egg, vomit"
- ❖ New vectors for words...KNOWLEDGE !

## AND...

- ❖ LSA gives rise to probabilistic LSA (pLSA); Bayesian approach to considering the hidden factors of the term:document matrix
- ❖ pLSA leads to **Latent Dirichlet Allocation (LDA)**, which comes to be used heavily for topic modelling
- ❖ Any given word within a doc could belong to a topic/factor and the set of topics can be used to define the document; just as the set of factors characterise a given doc; LDA finds these topics, n.b., familiar issues around knowing  $k$  no. of topics in advance
- ❖ Leads to other method: Expectation Maximisation, Non-Negative Matrix Factorisation, all heavily used in Text Analytics

## NB: LSA:Dimensional Reduction

### Real-World Example

Finally, I applied LDA to a set of Sarah Palin's emails a little while ago (see [here](#) for the blog post, or [here](#) for an app that allows you to browse through the emails by the LDA-learned categories), so let's give a brief recap. Here are some of the topics that the algorithm learned:

- **Trig/Family/Inspiration:** family, web, mall, god, son, from, congratulations, children, life, child, down, trig, baby, birth, love, you, syndrome, very, special, bless, old, husband, years, thank, best, ...
- **Wildlife/BP Corrosion:** game, fish, moose, wildlife, hunting, bears, polar, bear, subsistence, management, area, board, hunt, wolves, control, department, year, use, wolf, habitat, hunters, caribou, program, derby, fishing, ...
- **Energy/Fuel/Oil/Mining:** energy, fuel, costs, oil, alaskans, prices, cost, nome, now, high, being, home, public, power, mine, crisis, price, resource, need, community, fairbanks, rebate, use, mining, villages, ...
- **Gas:** gas, oil, pipeline, agia, project, natural, north, producers, companies, tax, company, energy, development, slope, production, resources, line, gasoline, transcanada, said, billion, plan, administration, million, industry, ...
- **Education/Waste:** school, waste, education, students, schools, million, read, email, market, policy, student, year, high, news, states, program, first, report, business, management, bulletin, information, reports, 2008, quarter, ...
- **Presidential Campaign/Elections:** mail, web, from, thank, you, box, mccain, sarah, very, good, great, john, hope, president, sincerely, wasilla, work, keep, make, add, family, republican, support, doing, p.o. ....

<http://blog.echen.me/2011/08/22/introduction-to-latent-dirichlet-allocation/>

Latent Stuff  
**Word Embeddings**

## Word Embeddings...

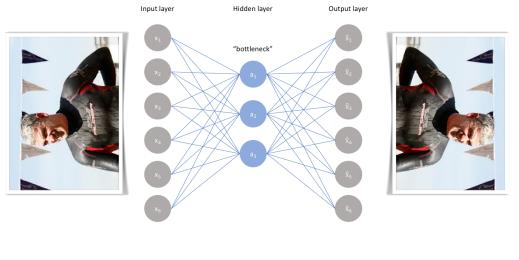
- ❖ LSA variant...computed differently...Google enter the fray...
- ❖ Creates multi-dimensional word vectors in a computationally efficient manner (for all); Vectors manifest semantic relations
- ❖ KING:QUEEN::PRINCE?: (using vectors offsets)
- ❖ A number of current competitors; explosion since 2015
  - ❖ word2vec
  - ❖ GloVe
  - ❖ Others....

## Embeddings: The Idea

- ❖ Co-occurrence analysis: Counting / Prediction
- ❖ *Counting*; basically, analyses a co-occurrence matrix using SVD
- ❖ *Prediction*: uses NNs to learn word-vectors based on co-occurrence, local-context sampling
- ❖ *Both*: GloVe (from Stanford)

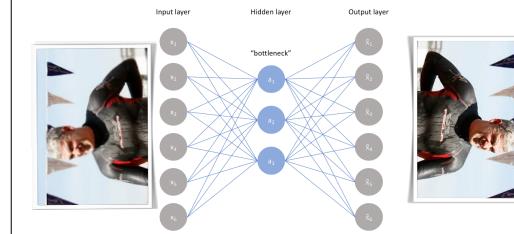
## Embeddings: Prediction

NNs learn patterns: e.g., autoencoder



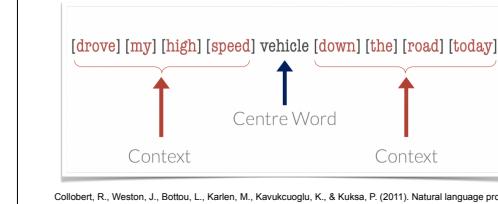
## Embeddings: Prediction

NNs can learn to complete patterns...



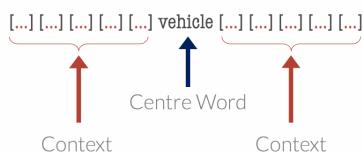
## Embeddings: Predictions

So, can we learn to predict a word from the context-words around it or predict the contexts from the word ?



## Embeddings: Predictions

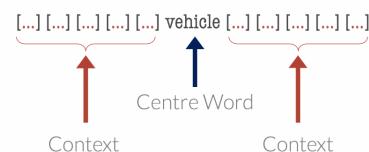
So, can we learn to predict a word from the context-words around it or predict the contexts from the word ? Generalised...



Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., & Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug), 2493-2537.

## Embeddings: Predictions

Slide a window over the text, or sample windows, noting word-contexts for each word to learn a vector



Collobert, R., & Weston, J. (2008, July). A unified architecture for natural language processing: Deep neural networks with multi-task learning. In *Proceedings of the 25th international conference on Machine learning* (pp. 160-167). ACM.

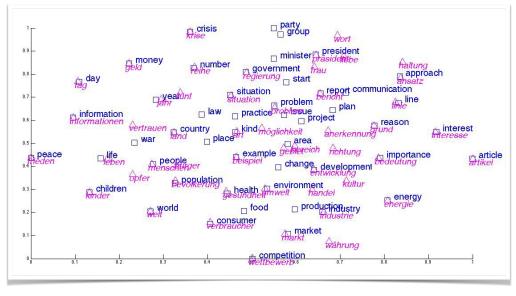
## Embeddings: Prediction

- ◆ Collobert & Weston (2008) using deep learning methods; archane, hard to compute
- ◆ Mikilova et al (2013) word2vec brings it to the masses (softmax, CBOW, Skip-gram)
- ◆ Vectors are better than ones from LSA-Count
- ◆ Area explode with many different techniques

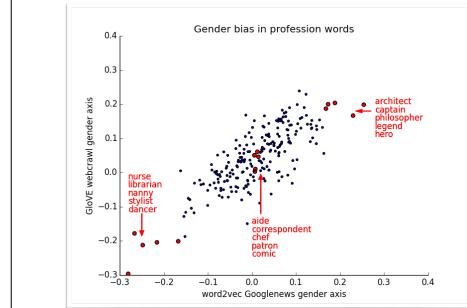
Mikolov, T., Chen, K., Corrado, G. and Dean J. (2013). Efficient estimation of word representations in vector space. CoRR, abs/1301.3781.

Goldberg, Y., & Levy, O. (2014). word2vec Explained: deriving Mikolov et al.'s negative-sampling word-embedding method. arXiv preprint arXiv:1402.3722.

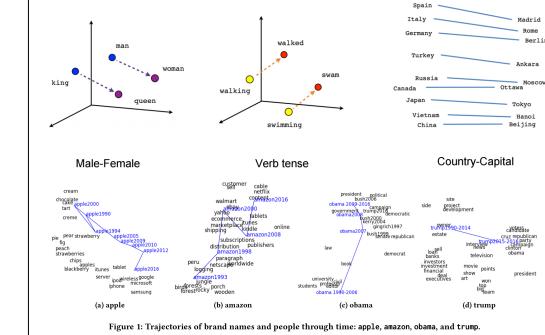
## Using vectors: Languages



## Using vectors: Sexism..



## Using vectors: Trajectories..



## Embeddings: Benefits

- ◆ Uses real world text-sets (all of wikipedia); so avoids pre-processing : *thnx, thanks, plz, thanx*, get similar vectors
- ◆ Normalised vectors mean that transforms in the space are meaningful" man->woman, boy->girl;
- ◆ GloVe makes all more formal

<https://www.shanelynn.ie/get-busy-with-word-embeddings-introduction/>

## Embeddings: Our Work

- ◆ Leavy, Pine & Keane (2017). Government report on clerical sexual abuse, Ryan Report.
- ◆ Patterns of "sacking" but not called that: "dispensation", "vows", and other indirections (e.g., latin)
- ◆ word2vec generate synonyms that enabled automated query expansion for search

Leavy, S., Pine, E., & Keane, M. (2017). Mining the Cultural Memory of Irish Industrial Schools Using Word Embedding and Text Classification. In DH-2017.

## Embeddings: Our Work...

1987	reagan	koch	soviet	iran.comra	nativelova	yuppie	walkman
1988	reagan	reagan	soviet	iran.comra	nativelova	yuppie	mp3 player
1989	bush	koch	soviet	iran.comra	nativelova	yuppie	walkman
1990	bush	dinkins	soviet	iran.comra	nativelova	yuppie	headphones
1991	bush	dinkins	soviet	iran.comra	nativelova	yuppie	cassette player
1992	bush	dinkins	soviet	iran.comra	nativelova	yuppie	boombox
1993	clinton	dinkins	russian	iran.comra	nativelova	yuppie	cd.player
1994	clinton	mr giuliani	russian	iran.comra	sanchez.vicario	yuppie	walkman
1995	clinton	mr giuliani	russian	iran.comra	sanchez.vicario	yuppie	mobile phone
1996	clinton	giuliani	russian	whitehouse	graf	yuppie	mobile phone
1997	clinton	giuliani	russian	iran.comra	hingga	yuppie	headphones
1998	clinton	giuliani	russian	iran.comra	hingga	yuppie	headphones
1999	clinton	mayer giuliani	russian	white house	hingga	yuppie	buttons
2000	clinton	giuliani	russian	white house	hingga	yuppie	bracelet
2001	clinton	giuliani	russian	iran.comra	hingga	yuppie	mp3 player
2002	bush	bloomberg	russian	white house	hingga	gen.x	walkman
2003	bush	bloomberg	russian	white house	agassi	hipsters	mobile phones
2004	bush	bloomberg	russian	white house	agassi	hipsters	car buds
2005	bush	bloomberg	north korean	iran.comra	rodrick	geek	headset
2006	bush	bloomberg	iranian	white house	hingga	teens	headset
2007	bush	bloomberg	iranian	capitol hill	federer	dads	pod

Table 1: Examples of words from 1987 and their analogues over time. Each column corresponds to a single point in vector space, and each row shows the word closest to that point in a given year.

Szymanski, T. (2017). Temporal word analogies: Identifying lexical replacement with diachronic word embeddings. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers) (Vol. 2, pp. 448-453).

## Embeddings: Some Issues

- ◆ Note, people typically use the vectors to do sthm-else; like idea of transfer learning
- ◆ Different runs produce different answers as the text is sampled in various ways...
- ◆ But, algorithmic stability + new metrics

Wendlandt, L., Kummerfeld, J. K., & Mihalcea, R. (2018). Factors Influencing the Surprising Instability of Word Embeddings. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers) (Vol. 1, pp. 2092-2102).

## Embeddings: Why Its Hot !

- ◆ Delivers very good “semantic” vectors
- ◆ Easy of use by anyone; eg. recomputing LSA tests in psych showing better results
- ◆ Many things can be treated as sentence: places I go to each day, genes...

## Why Its Hot !: Other Sequences

Session: Personalized Social Web II

UMAP'18, July 8–11, 2018, Singapore

### Predict Demographic Information Using Word2vec on Spatial Trajectories

Adir Solomon, Ariel Bar, Chen Yanai, Bracha Shapira, Lior Rokach  
Department of Software and Information Systems Engineering  
Telemek Innovation Laboratories at BGU  
Ben-Gurion University of the Negev, P.O. Box 84105, Beer-Sheva, Israel  
(adirsolo, arieli, chenyan)@post.bgu.ac.il, {bshapira, liorl}@bgu.ac.il

## Conclusions

- ◆ We stared by considering how we can get at hidden-variables by dimensional reduction
- ◆ Continuous evolution from early work though PCA to SVD to word2vec...
- ◆ At the beginning of a very exciting step...in capturing the meaning of words