

DIMENSIONALITY REDUCTION:

PCA AND SVD

Paul Burkard

10/28/2016

I. DIMENSIONALITY REDUCTION

II. FEATURE SELECTION

III. FEATURE EXTRACTION - PCA AND SVD

HANDS-ON: FEATURE EXTRACTION WITH PCA

- ▶ What is Dimensionality Reduction?
 - ▶ Why might we want to do it?
 - ▶ What are the broad types?
- ▶ What are feature selection and feature extraction?
- ▶ What is the goal of PCA? SVD?

I. DIMENSIONALITY REDUCTION

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dim reduction</i>	<i>clustering</i>

Q: *What is **dimensionality reduction**?*

A: *A set of techniques for **reducing the size** (in terms of features, records, and/or bytes) **of the dataset** under examination.*

*In general, the idea is to regard the dataset as a matrix and to **decompose the matrix** into simpler, meaningful pieces.*

*Dimensionality reduction is frequently performed as a **pre-processing step** before another learning algorithm is applied.*

Q: *What are the motivations for dimensionality reduction?*

A: *The number of features in our dataset can be difficult to manage, or even misleading (eg, if the relationships are actually simpler than they appear).*

Q: *What is the goal of dimensionality reduction?*

A:

- *reduce computational expense*
- *reduce susceptibility to overfitting*
- *reduce noise in the dataset*
- *enhance our intuition*
- *reduce multicollinearity*

Q: *What are some applications of dimensionality reduction?*

A:

- topic models (document clustering)*
- image recognition/computer vision*
- recommender systems*

Q: *How is dimensionality reduction performed?*

A: *There are two approaches: **feature selection** and **feature extraction**.*

feature selection – *selecting a subset of features using an external criterion (filter) or the learning algorithm accuracy itself (wrapper)*

feature extraction – *mapping the features to a lower dimensional space*

The goal of feature selection is to select out the best possible subset of features for model-building from the original available features.

Feature selection is important, but typically when people say dimensionality reduction, they are referring to feature extraction.

The goal of feature extraction is to create a new set of coordinates that simplify the representation of the data.

II. FEATURE SELECTION

The goal of feature selection is to select out the best possible subset of features for model-building from the original available features.

This problem can be thought of as a search through the space of all possible feature subsets for the optimal subset.

For even a moderate number of features, this comprehensive search becomes impossible, so we need a heuristic approach.

Q: *How do we perform **feature selection**?*

A: *By making use of **wrappers**, **filters**, or **embedded methods***

wrappers** – potential feature subsets are compared based on the success of **built models** projected via **cross-validation

***filters** – feature subsets are determined based on some simple prescribed metric over the features*

***embedded** – feature selection happens within the model-building itself*

Wrappers use some criteria for trying out various feature subsets and build models of the desired type with each subset.

The models' performance are all estimated via cross-validation and the feature subset chosen that yields the greatest performance.

Ex: Stepwise Regression – Starting from zero features, test out adding each feature alone and training the model. Keep the feature that leads to the best model. Continue adding features until no more improvement

Filters use a prescribed general metric for determining which features to include in the model.

Filters are often far less computationally expensive than wrappers as the inclusion criteria for features is easy to compute.

Ex:

- Information gain*
- Correlation Coefficient, etc.*

Embedded methods have the model-building itself incorporate the feature selection.

These methods often strike a good balance between the strengths and weaknesses of wrappers and filters.

Ex: A good example is the LASSO Regression, in which the L_1 regularized regression tends to zero out coefficients and implicitly choose features to be excluded.

III. FEATURE EXTRACTION

*The goal of **feature extraction** is to create a **new set of coordinates** that simplify the representation of the data.*

*Typically we do this by using **matrix factorizations** to map the features to a lower-dimensional space that **minimizes information loss**.*

*Two prominent examples of such matrix factorization methods are **Principal Component Analysis (PCA)** and **Singular Value Decomposition (SVD)***

INTRO TO DATA SCIENCE

PCA

Principal component analysis is a dimension reduction technique that can be used on a matrix of any dimensions.

*This procedure produces a new **basis**, each of whose components retain as much **variance** from the original data as possible.*

*The PCA of a matrix A boils down to the **eigenvalue decomposition** of the **covariance matrix** of A .*

What is *variance*? $s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{(n - 1)}$

Variance is the average of the squared distance between the mean of a dataset and a point in the dataset.

*In other words, it is a measure of the spread in the dataset. Recall that the square root of the variance is the **standard deviation**.*

What is *covariance*?

A measure of how much 2 random variables change together.

Variance:

$$s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{(n-1)}$$

$$\text{var}(X) = \frac{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}{(n-1)}$$

Covariance:

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{(n-1)}$$

The covariance matrix C of a matrix A is always square:

$$C = \begin{bmatrix} E[(X_1 - \mu_1)(X_1 - \mu_1)] & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \cdots & E[(X_1 - \mu_1)(X_n - \mu_n)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[(X_2 - \mu_2)(X_2 - \mu_2)] & \cdots & E[(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - \mu_n)(X_1 - \mu_1)] & E[(X_n - \mu_n)(X_2 - \mu_2)] & \cdots & E[(X_n - \mu_n)(X_n - \mu_n)] \end{bmatrix}.$$

off-diagonal elements C_{ij} give the covariance between X_i, X_j ($i \neq j$)

diagonal elements C_{ii} give the variance of X_i

The eigenvalue decomposition of a square matrix A is given by:

$$A = W \Lambda W^{-1} = W \Lambda W^T$$

*The columns of W are the **eigenvectors** of A , and the values in Λ are the associated **eigenvalues** of A .*

For an eigenvector v of A and its eigenvalue λ , we have the important relation:

$$Av = \lambda v$$

NOTE

This relationship defines what it means to be an eigenvector of A .

*The eigenvectors form a **basis** of the vector space on which A acts (eg, they are orthogonal).*

*That means they're a **transformed coordinate space**.*

*Furthermore the basis elements are ordered by their eigenvalues (from largest to smallest), and these **eigenvalues represent the amount of variance explained by each basis element**.*

Q: *So what goes into a PCA?*

A: *Covariance Matrix of the data matrix X*

- *Fact: $\text{cov}(X) \sim X^T X = A$*
- *Perform eigenvalue decomposition:*
 - $A = W \Lambda W^T$
 - *The columns of W are **eigenvectors of A***
 - *Each eigenvector has an associated eigenvalue*
 - *The diagonal of Λ are **eigenvalues of A***
 - *The eigenvalues are **ordered** in non-decreasing fashion*

Let's think about dimensions:

- *\mathbf{X} is a data matrix of n instances and m features*
 - *$\dim(\mathbf{X}) = n \times m$*
- *$\text{cov}(\mathbf{X}) \sim \mathbf{X}^T \mathbf{X} = \mathbf{A} \longrightarrow \dim(\mathbf{A}) = (m \times n) \times (n \times m) = m \times m$*
- *Eigenvalue decomposition: $\mathbf{A} = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^T$*
 - *$\dim(\mathbf{A}) = m \times m \longrightarrow \dim(\mathbf{W}) = m \times p, \dim(\mathbf{\Lambda}) = p \times p, \dim(\mathbf{W}^T) = p \times m$*
 - *What is p ?*
 - *The **rank** of \mathbf{A}*
 - *Essentially, the # of independent dimensions in \mathbf{X}*
 - *e.g.: data on a perfect line has rank 1, plane has rank 2, etc*

Q: *So what comes out of a PCA?*

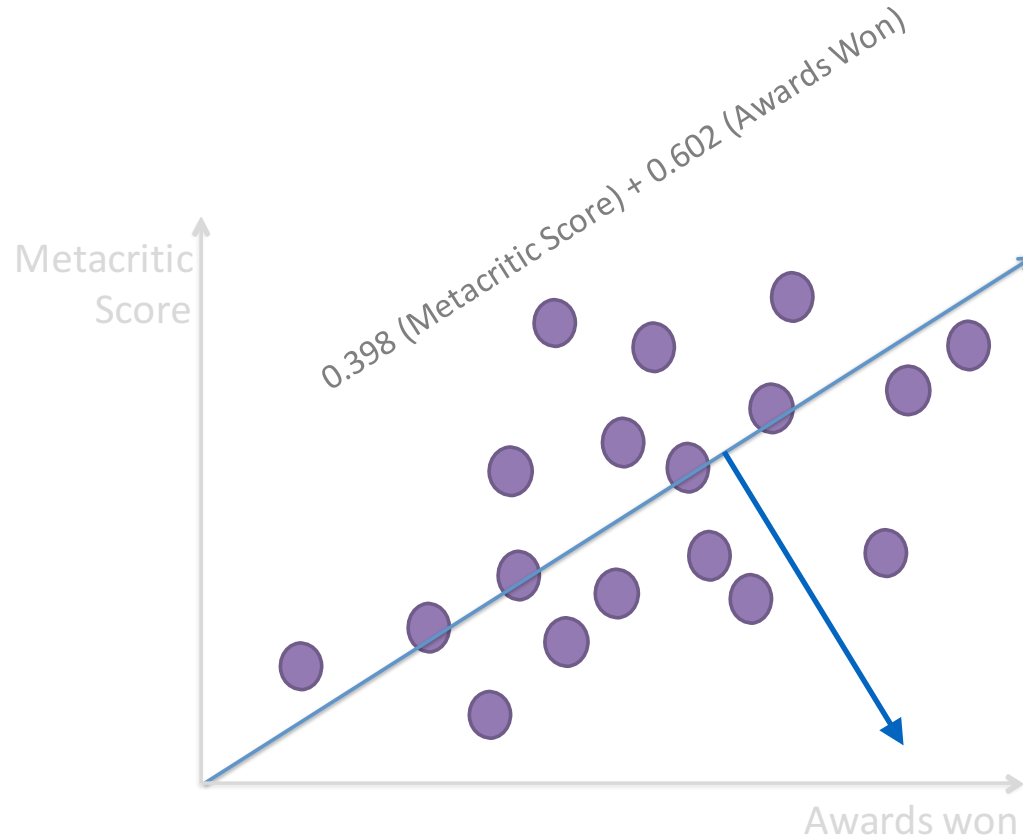
A: *Eigenvectors and eigenvalues.*

- *Eigenvectors are **linear combinations** of the original feature vectors*
- *Each **eigenvector** represents a feature in our new transformed feature space, and are called **principal components***
- *The eigenvalues represent a measure of the amount of variance explained by each corresponding eigenvector (“new feature”)*
- *We can choose only the first k (whatever we like) of our “new features” from the eigenvector space and work with them as our new data knowing we’ll have minimal data loss for a feature space of that size*

Q: *So what comes out of a PCA?*

A: *A transformed coordinate space.*

- *The eigenvectors (principal components) are our new axes.*
- *Our data matrix in this transformed space T is:*
 - $\mathbf{T} = \mathbf{XW}$
- *The features in \mathbf{T} are perfectly independent (orthogonal)!*
- *We keep only the first k (whatever you like) features in T , and know we've retained maximal information!*
- $\mathbf{T}_k = \mathbf{XW}_k$, where \mathbf{W}_k is the first k columns of \mathbf{W}
- $\dim(\mathbf{T}_k) = (n \times m) \times (m \times k) = n \times k$



$$W = \begin{bmatrix} 0.602 & 0.551 \\ 0.398 & -0.834 \end{bmatrix}$$

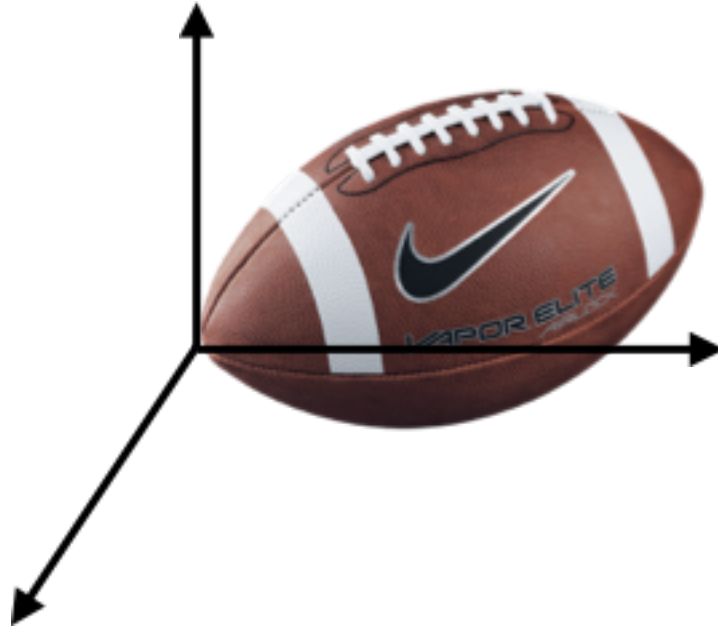
$$T = XW$$

$$t_1 = 0.602x_1 + 0.398x_2$$

$$t_2 = 0.551x_1 - 0.834x_2$$

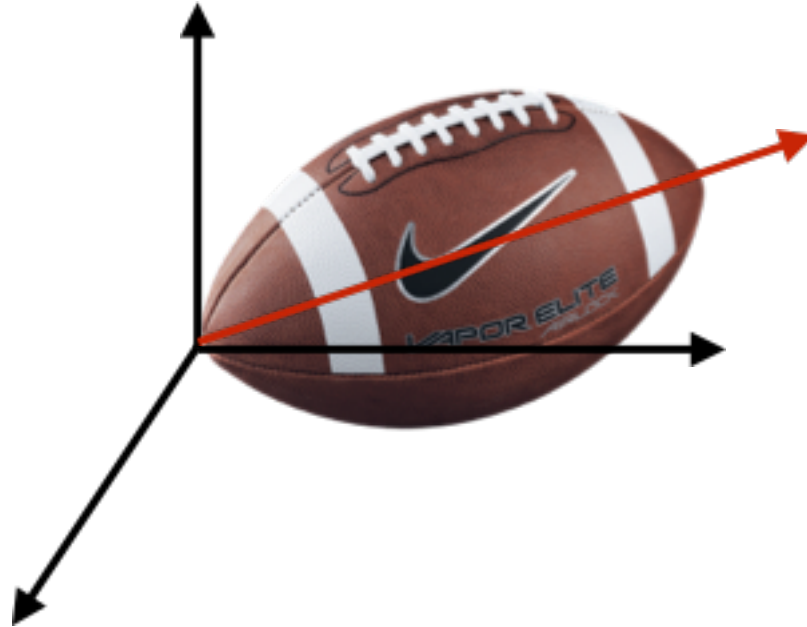
Let's try an example: Mapping 3D to 2D

Suppose we have a dataset representing random points taken inside a football:



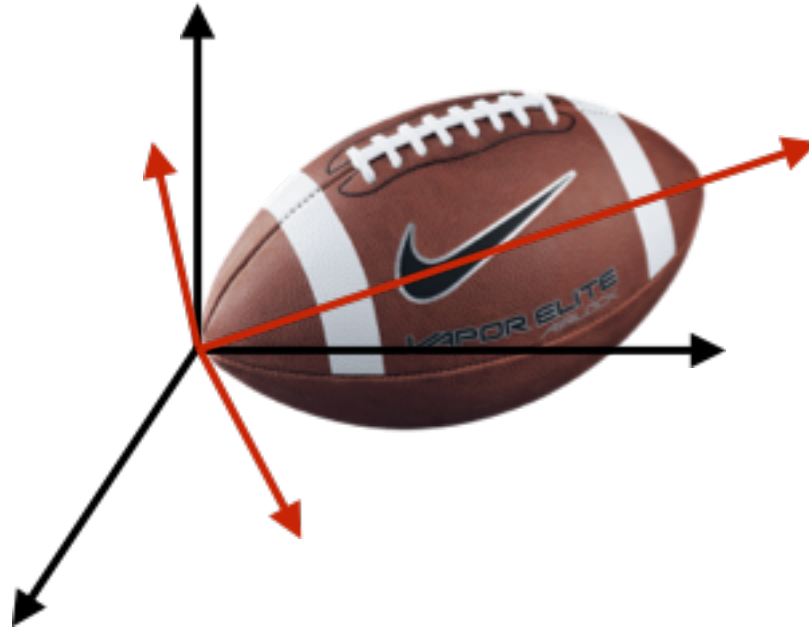
What is the 1st Principal Component Axis?

*The axis along which the data **varies the most***

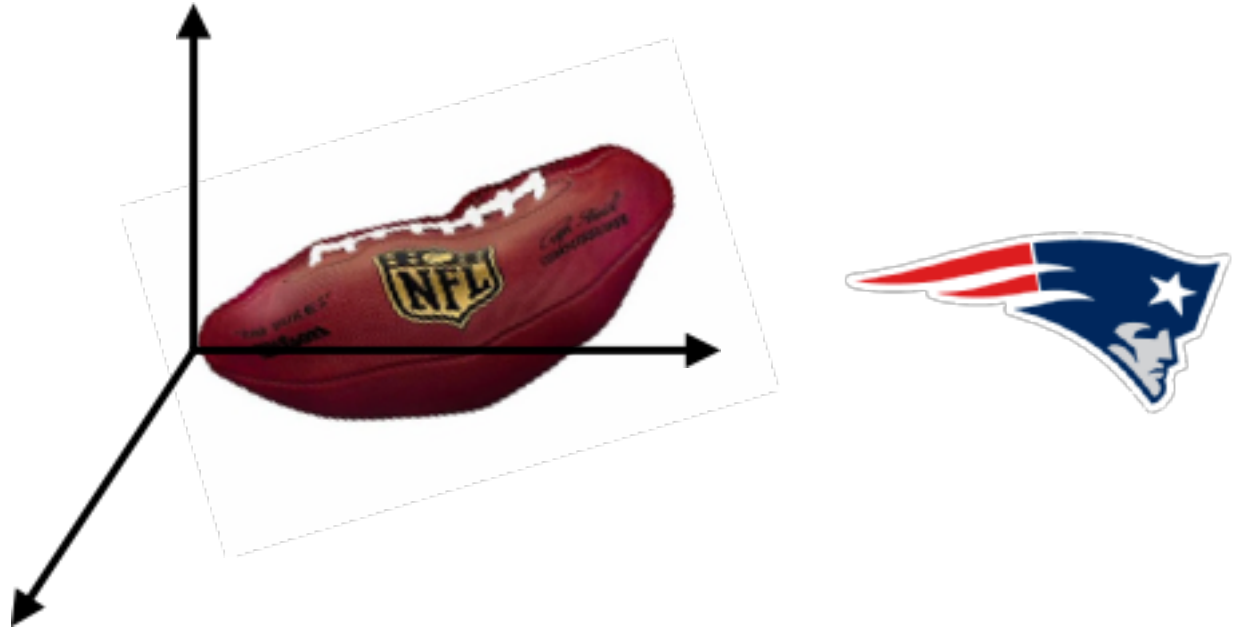


What is the 2nd Principal Component Axis? The 3rd?

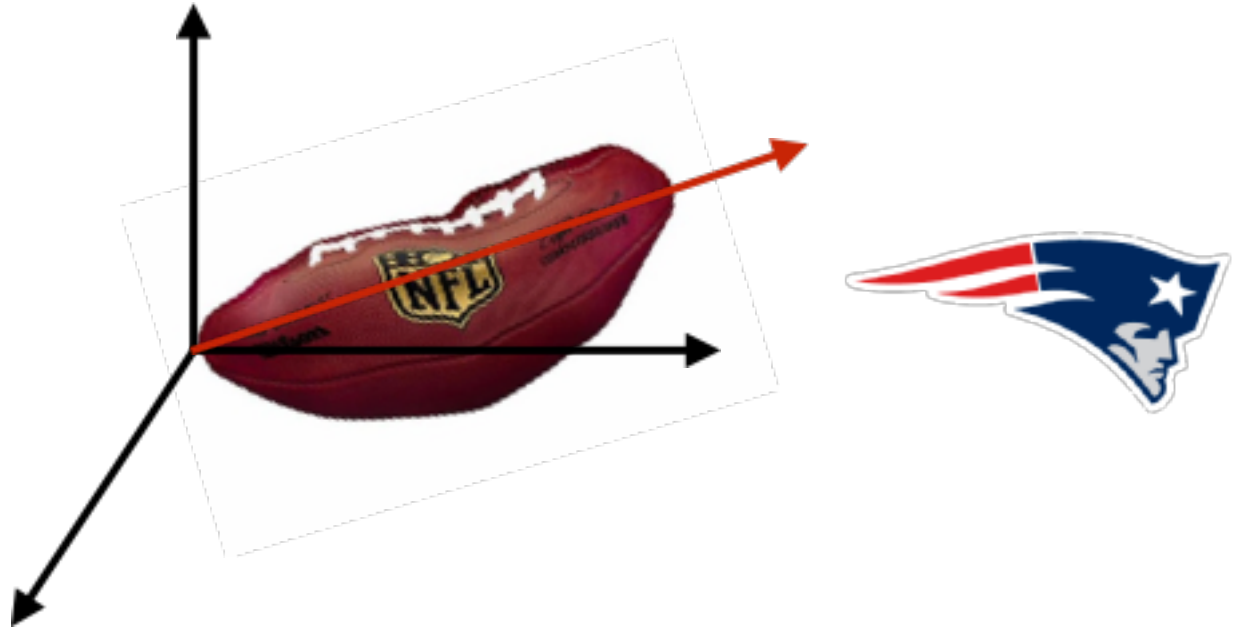
We have many choices right? The football is symmetric along the rest



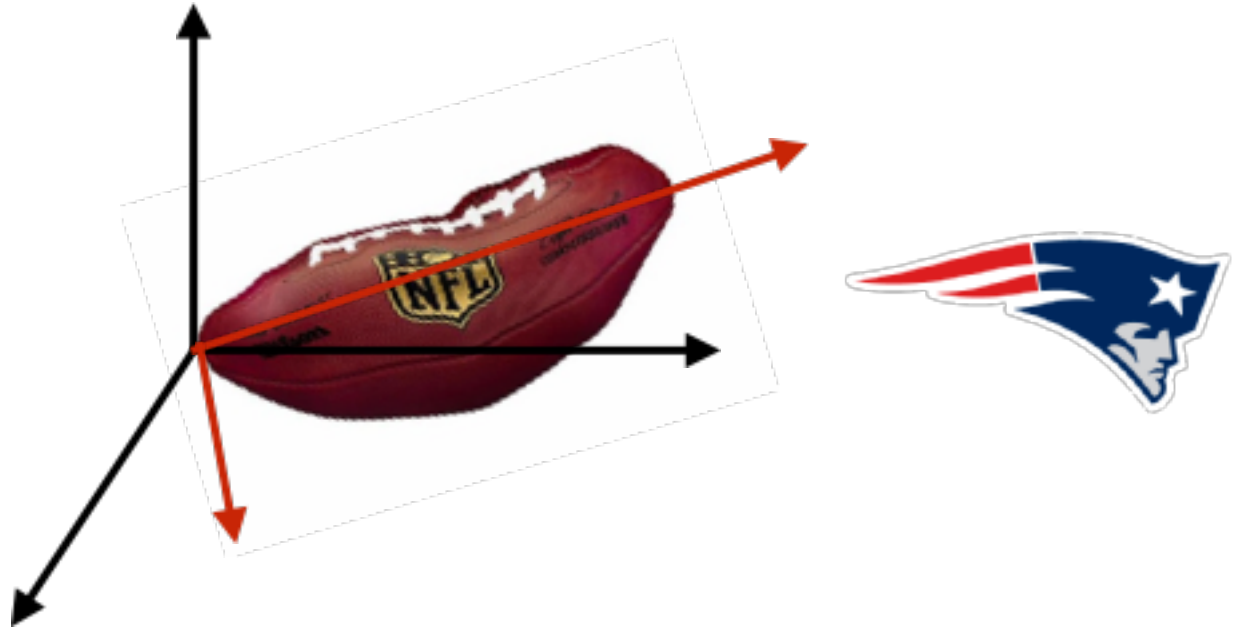
Alright what if the New England Patriots com along and decide to flatten our football a little bit...



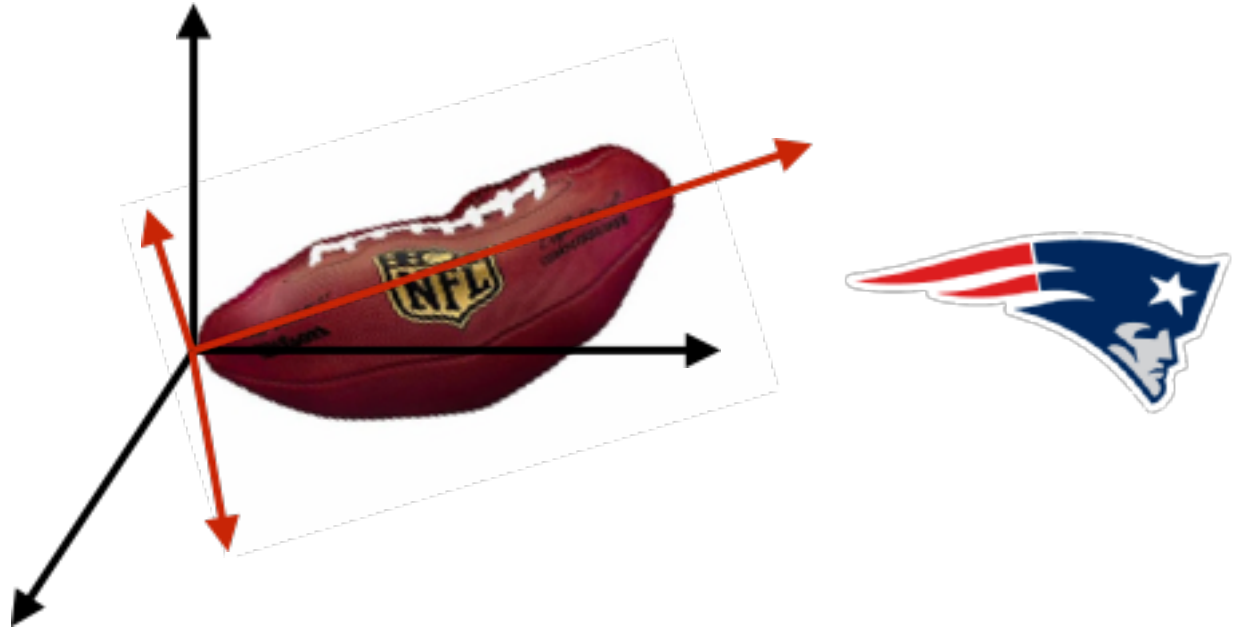
The first axis is still basically the same...



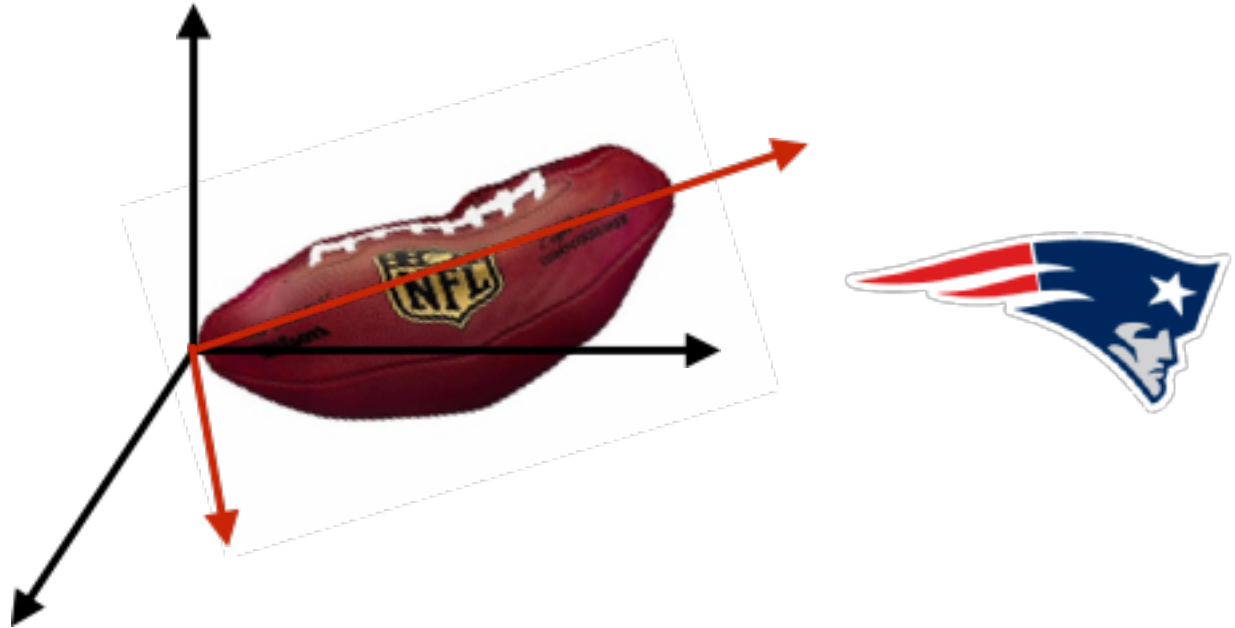
*But now what's the 2nd axis (PC) aka axis with 2nd most variance?
It's coming at us, out of the page, because the ball is squashed!*



The 3rd PC must be orthogonal to the first 2. It's the remaining, "flattened", dimension (up through the laces)!



*If we want, we can retain only the first 2 PC axes going forward.
This retains as **much variance in our data as possible** in 2D.*



What's happening?

*By keeping the **first k (2D) Principal Components**, we're **projecting** our original datapoints in n -dimensional (3D) feature space onto a **new k -dimensional (2D) feature space**.*

*This k -dimensional space **retains as much original variance (information)** from the original data as possible —> **Dimensionality Reduction***

*The **eigenvalue decomposition of the covariance matrix** is the math that finds these projections.*

Q: *So what comes out of a PCA?*

A: *Eigenvectors and eigenvalues.*

- *Eigenvectors are **linear combinations** of the original feature vectors*
- *Each **eigenvector** represents a feature in our new transformed feature space, and are called **principal components***
- *The eigenvalues represent a measure of the amount of variance explained by each corresponding eigenvector (“new feature”)*
- *We can choose only the first k (whatever we like) of our “new features” from the eigenvector space and work with them as our new data knowing we’ll have minimal data loss for a feature space of that size*

Some things to be aware of:

- **PCA is not scale invariant**
 - *Make sure to scale and mean-center (subtract mean to make the mean of all features 0) before performing a PCA*
 - *sklearn can take care of this for you*
- **Feature Importance:**
 - *Can compute explicit proportion of variance explained for each feature t_i :*

$$\text{Explained Variance}_{t_i} = \frac{\lambda_i}{\sum_i \lambda_i} \quad \sum_i \text{Explained Variance}_{t_i} = 1$$

INTRO TO DATA SCIENCE

USING A PCA

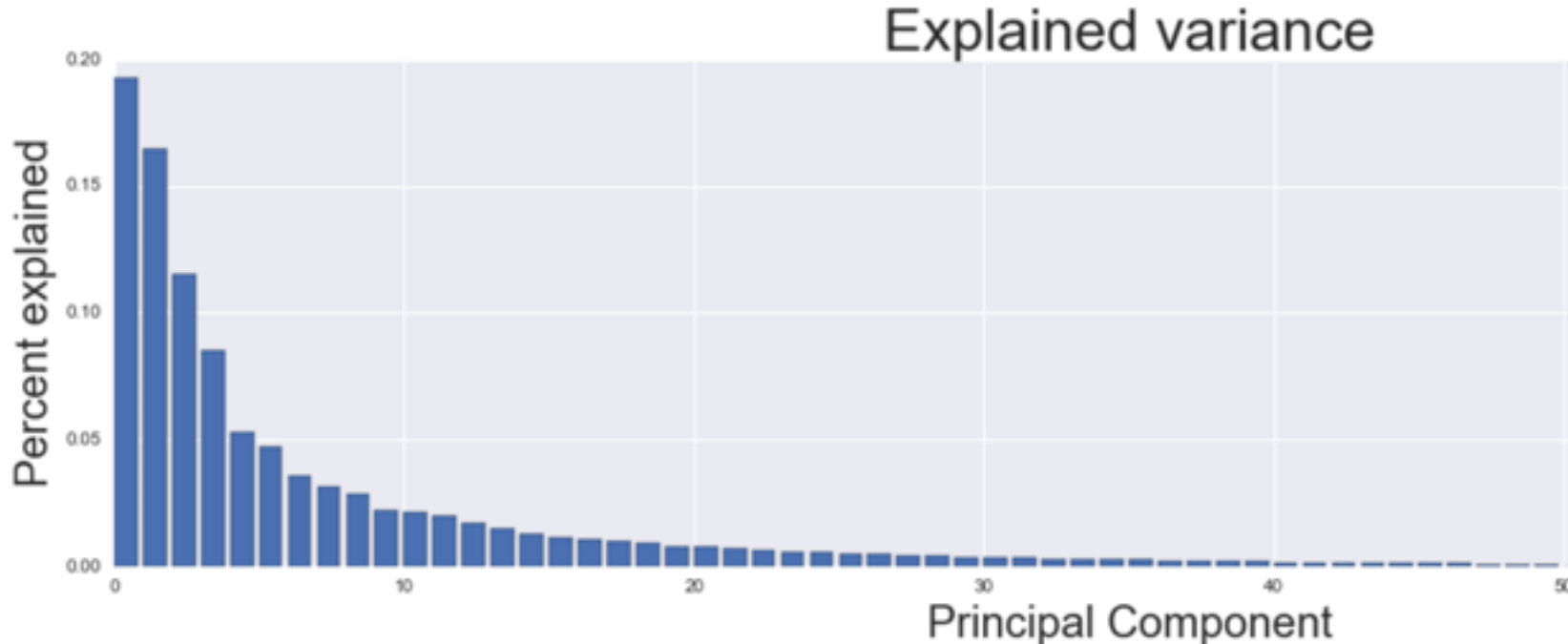
Q: *Okay, we understand PCA, what do we do with it?*

A: *Any ML that you want!*

- *We now have a transformed data matrix T_k ($k \leq m$)*
- *We can **visualize, regress, classify, cluster**, whatever!*
- *Common Workflow:*
 - *Use PCA to reduce to 2 dimensions T_2*
 - ***Explore/Visualize** your data in this space, best 2D space we can get!*
 - *Appear to be clusters? Relationships? etc*
 - ***Build Supervised Models** in the PCA space (of dimension k)*

Q: How to *choose* k (number of features in T to keep)?

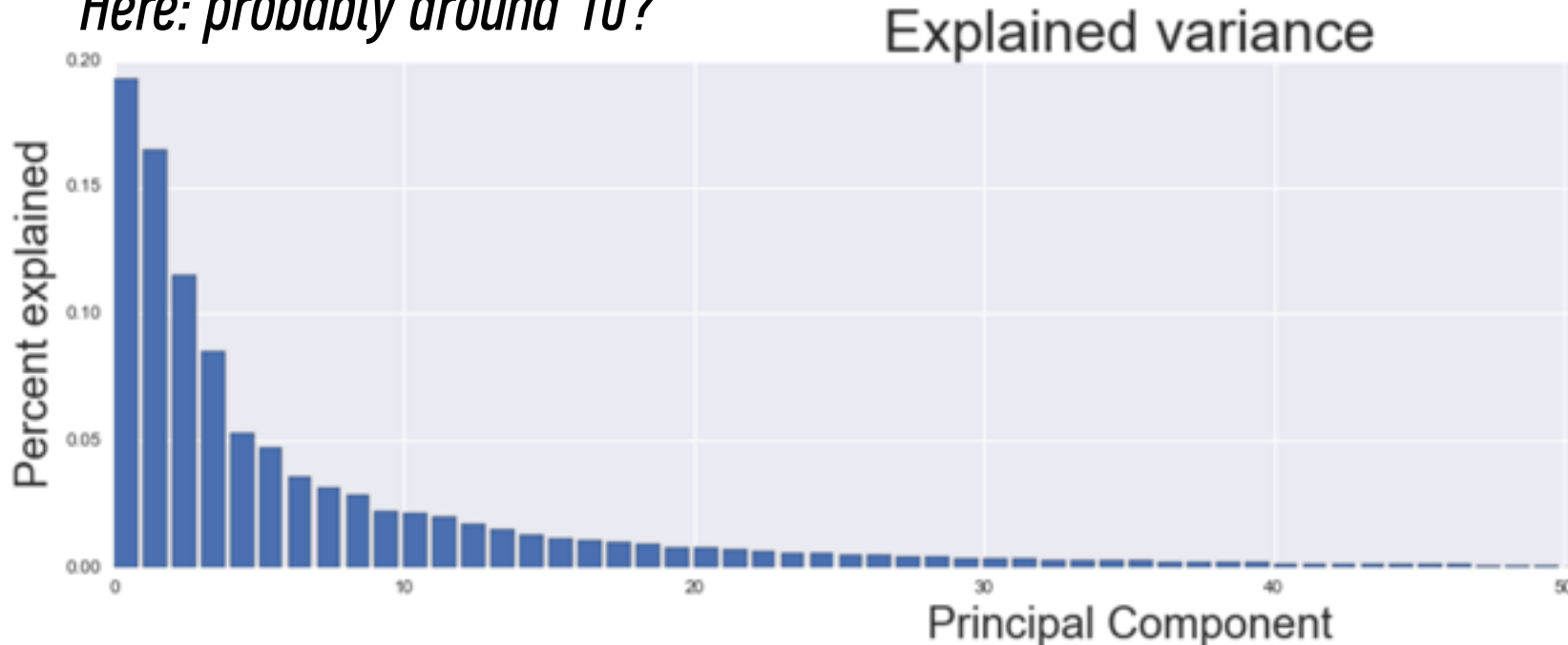
A: Look at the *variance explained* of each feature (plot it)!



Q: *How many Principal Components would you retain here?*

A: *Depends on context, but look for the “elbow”/“knee” in the curve*

Here: probably around 10?



INTRO TO DATA SCIENCE

SVD

*Like PCA, **SVD** is a **matrix factorization** technique that can be used to **reduce the dimensionality** of a data matrix.*

NOTE

PCA and SVD are related by a relatively simple transformation.

Thus, SVD can be (and is) used to compute a PCA more efficiently numerically.

*Like PCA, this procedure produces a new **basis**, each of whose components **retain as much variance from the original data as possible**.*

Consider a matrix M with m rows and n features.

The singular value decomposition of M is given by:

$$\begin{array}{ccccc} M & = & U & \Sigma & V^T \\ \text{(m x n)} & & \text{(m x r)} & \text{(r x r)} & \text{(r x n)} \end{array}$$

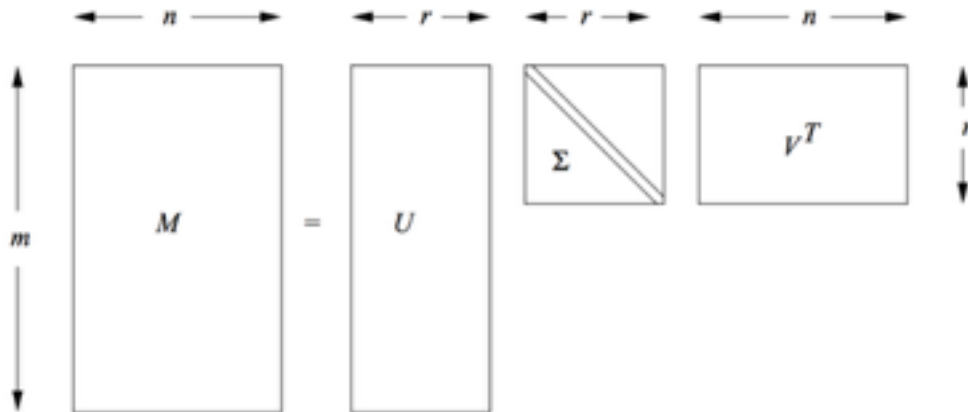
*st. U, V are **orthogonal matrices** and Σ is a **diagonal matrix**.*

$$\rightarrow UU^T = I_m, \quad VV^T = I_n \rightarrow \Sigma_{ij} = 0 \quad (i \neq j)$$

The singular value decomposition of M is given by:

$$M = U \Sigma V^T$$

$(m \times n)$ $(m \times r)$ $(r \times r)$ $(r \times n)$



The nonzero entries of Σ are the singular values of M . These are real, nonnegative, and rank-ordered (decreasing from left to right).

NOTE

The number of singular values is equal to the *rank* of M .

The rank of a matrix measures its *non-degeneracy*.

NOTE

PCA factorizes the covariance matrix of X ($X^T X$).

SVD directly factorizes the data matrix X .

Ratings of movies by users:

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

there are two “concepts” underlying the movies:

science-fiction and romance

Ratings of movies by users:

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix} = \begin{bmatrix} .14 & 0 \\ .42 & 0 \\ .56 & 0 \\ .70 & 0 \\ 0 & .60 \\ 0 & .75 \\ 0 & .30 \end{bmatrix} \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \begin{bmatrix} .58 & .58 & .58 & 0 & 0 \\ 0 & 0 & 0 & .71 & .71 \end{bmatrix}$$

M
 U
 Σ
 V^T

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2

$$\begin{matrix}
 \begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 0 & 0 & 2 & 2 \end{bmatrix} & = & \begin{bmatrix} .14 & 0 \\ .42 & 0 \\ .56 & 0 \\ .70 & 0 \\ 0 & .60 \\ 0 & .75 \\ 0 & .30 \end{bmatrix} & \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} & \begin{bmatrix} .58 & .58 & .58 & 0 & 0 \\ 0 & 0 & 0 & .71 & .71 \end{bmatrix} \\
 M & & U & \Sigma & V^T
 \end{matrix}$$

M: people → movies

U: people → concepts

V: concepts → movies

Σ: the strength of each of the concepts

	Matrix	Alien	Star Wars	Casablanca	Titanic
Joe	1	1	1	0	0
Jim	3	3	3	0	0
John	4	4	4	0	0
Jack	5	5	5	0	0
Jill	0	0	0	4	4
Jenny	0	0	0	5	5
Jane	0	0	0	2	2



$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} = M'$$

$$\begin{bmatrix} .13 & .02 & -.01 \\ .41 & .07 & -.03 \\ .55 & .09 & -.04 \\ .68 & .11 & -.05 \\ .15 & -.59 & .65 \\ .07 & -.73 & -.67 \\ .07 & -.29 & .32 \end{bmatrix} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \begin{bmatrix} .56 & .59 & .56 & .09 & .09 \\ .12 & -.02 & .12 & -.69 & -.69 \\ .40 & -.80 & .40 & .09 & .09 \end{bmatrix}$$

$U \quad \Sigma \quad V^T$

How to reduce dimensions?

Drop Low Singular Values \rightarrow eliminate corresponding rows of U and V

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} =$$

M'



$$\begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix}$$

Σ

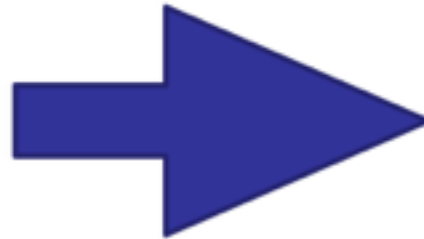
$$\begin{bmatrix} .13 & .02 & -.01 \\ .41 & .07 & -.03 \\ .55 & .09 & -.04 \\ .68 & .11 & -.05 \\ .15 & -.59 & .65 \\ .07 & -.73 & -.67 \\ .07 & -.29 & .32 \end{bmatrix} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \begin{bmatrix} .56 & .59 & .56 & .09 & .09 \\ .12 & -.02 & .12 & -.69 & -.69 \\ .40 & -.80 & .40 & .09 & .09 \end{bmatrix}$$

$U \qquad \qquad \Sigma \qquad \qquad V^T$

How to reduce dimensions?
Drop Low Singular Values

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 3 & 3 & 3 & 0 & 0 \\ 4 & 4 & 4 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 2 & 0 & 4 & 4 \\ 0 & 0 & 0 & 5 & 5 \\ 0 & 1 & 0 & 2 & 2 \end{bmatrix} =$$

M'



$$\begin{bmatrix} .13 & .02 \\ .41 & .07 \\ .55 & .09 \\ .68 & .11 \\ .15 & -.59 \\ .07 & -.73 \\ .07 & -.29 \end{bmatrix} \begin{bmatrix} 12.4 & 0 \\ 0 & 9.5 \end{bmatrix} \begin{bmatrix} .56 & .59 & .56 & .09 & .09 \\ .12 & -.02 & .12 & -.69 & -.69 \end{bmatrix}$$

$$\begin{bmatrix} .13 & .02 & -.01 \\ .41 & .07 & -.03 \\ .55 & .09 & -.04 \\ .68 & .11 & -.05 \\ .15 & -.59 & .65 \\ .07 & -.73 & -.67 \\ .07 & -.29 & .32 \end{bmatrix} \begin{bmatrix} 12.4 & 0 & 0 \\ 0 & 9.5 & 0 \\ 0 & 0 & 1.3 \end{bmatrix} \begin{bmatrix} .56 & .59 & .56 & .09 & .09 \\ .12 & -.02 & .12 & -.69 & -.69 \\ .40 & -.80 & .40 & .09 & .09 \end{bmatrix}$$

 U
 Σ
 V^T

$$= \begin{bmatrix} 0.93 & 0.95 & 0.93 & .014 & .014 \\ 2.93 & 2.99 & 2.93 & .000 & .000 \\ 3.92 & 4.01 & 3.92 & .026 & .026 \\ 4.84 & 4.96 & 4.84 & .040 & .040 \\ 0.37 & 1.21 & 0.37 & 4.04 & 4.04 \\ 0.35 & 0.65 & 0.35 & 4.87 & 4.87 \\ 0.16 & 0.57 & 0.16 & 1.98 & 1.98 \end{bmatrix}$$

- *SVD is immensely valuable.*
- *We'll return to it next week for text data!*

INTRO TO DATA SCIENCE

HANDS-ON: PCA