Dimensionality Reduction:
CS 445/545
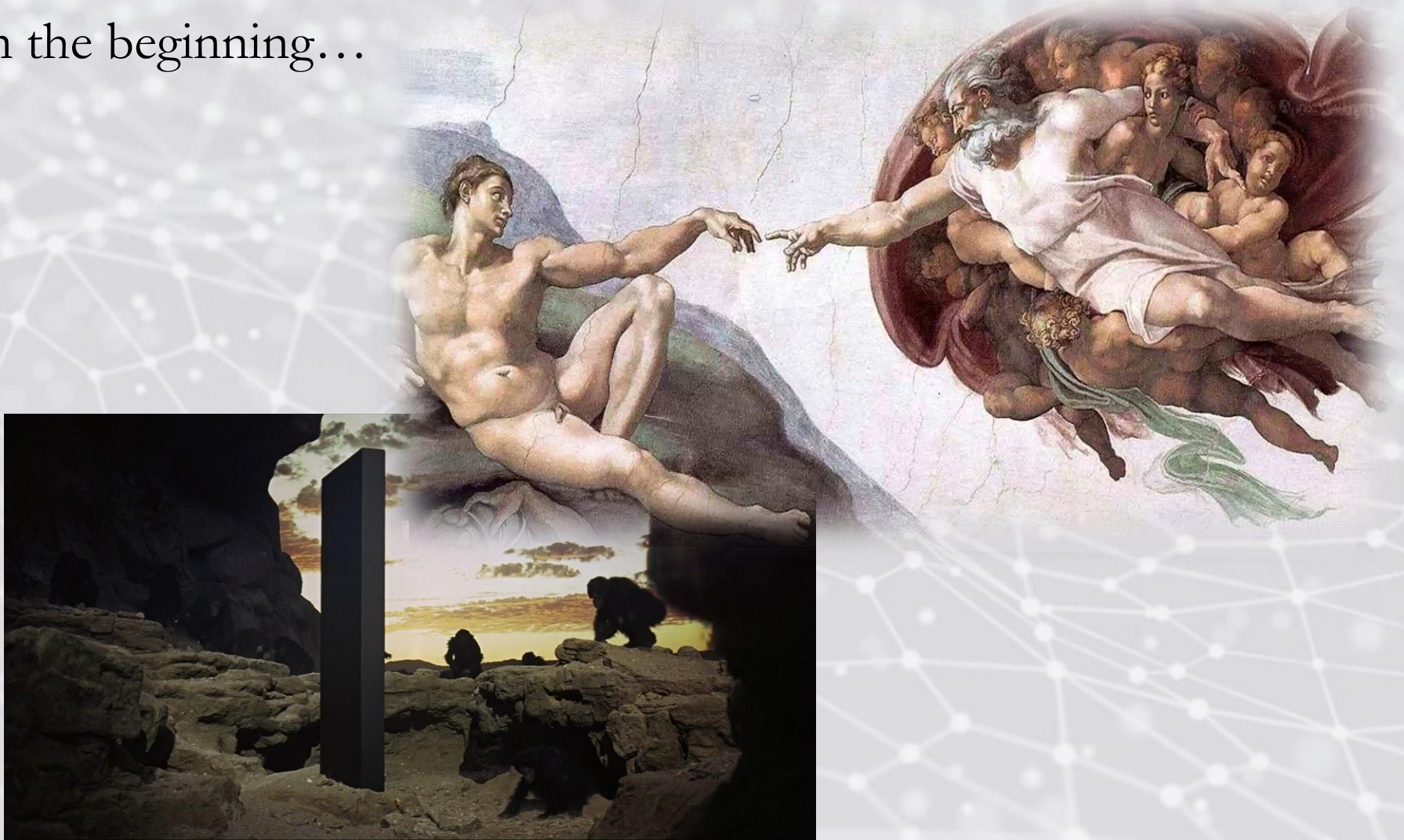
# Outline

- "Big data" and motives for dimensionality reduction

- **LDA** (linear discriminant analysis)

- **PCA** (principal component analysis)
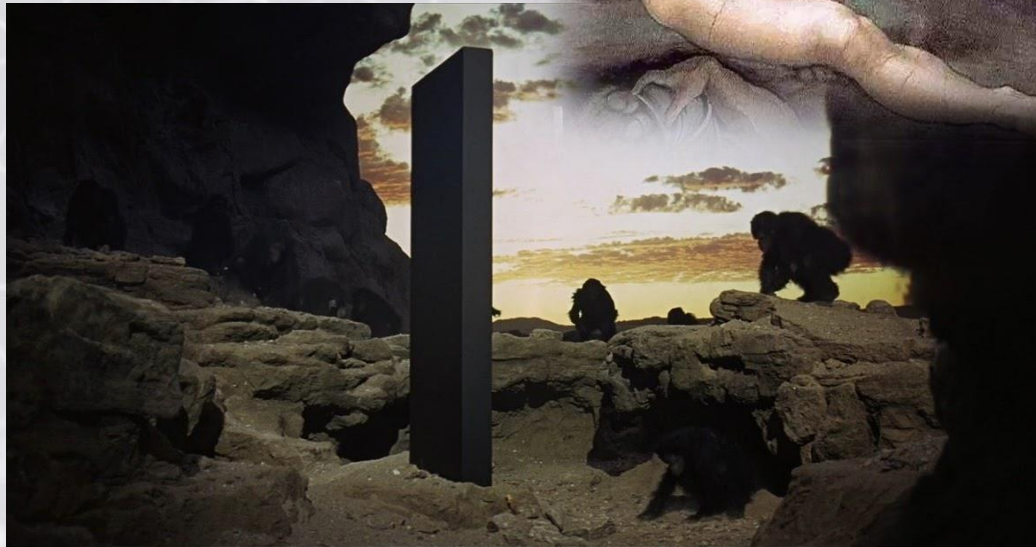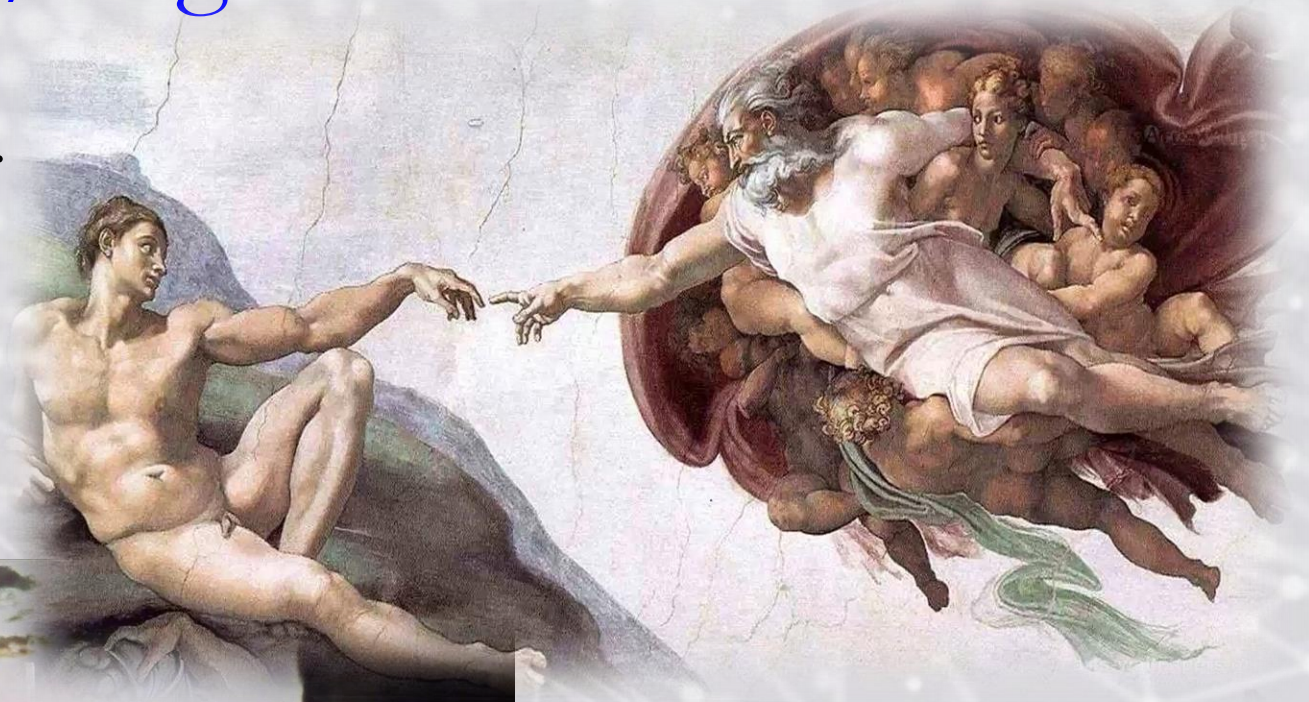
- **SVD** (singular value decomposition)


I WAS INTO BIG DATA

BEFORE BIG DATA WAS BIG

# Big Data/High-Dimensional Data

In the beginning…

# Big Data/High-Dimensional Data

In the beginning…

There was *small data*…

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Big Data/High-Dimensional Data

# Aside: Is the Basis of the Universe Information?



• In the 1940s, "the father of the digital age," Claude Shannon, formalized the notion of *information* through *entropy*.

$$H = -\sum_i p_i \log p_i$$

• The theoretical physicist John Archibald Wheeler later stated that The universe had (3) parts:

 (1) Everything is particles
 (2) Everything is fields
 (3) Everything is information

# Aside: Is the Basis of the Universe Information?

- A study in 2017 revealed substantial evidence that **we live in a holographic universe**.

- In this view, we might be caught inside a giant hologram; the cosmos is a projection, much like a 3D simulation.



holographic phase

Physics of early Universe leaves imprint on cosmic microwave background

Seed formation of stars and galaxies in late time Universe

- If the nature of reality is in fact <u>reducible to information itself</u>, that implies a conscious mind on the receiving end, to interpret and comprehend it.

# Aside: Is the Basis of the Universe Information?

- Wheeler himself believed in a *participatory universe*, where consciousness holds a central role.

- It is possible that information theory may in the future help bridge the gap between general relativity and quantum mechanics, or aid in our understanding of dark matter.





"The universe is a physical system that contains and processes information in a systematic fashion and that can do everything a computer can do." – Seth Lloyd MIT

Susskind: "On the world as hologram" https://www.youtube.com/watch?v=2DIl3Hfh9tY

# Aside: Is the Basis of the Universe Information?

- Just for fun…here is a short conversation with Minsky on the question of whether information is a basic building block of reality.



https://www.closertotruth.com/series/information-fundamental

# Introduction

- Most traditional statistical techniques (e.g. regression/classification) were developed in *low-dimensional* settings (i.e. $n >> p$ where $n$ is the data size and $p$ is the number of features).

- Over the last several decades, new technologies have drastically changed the way that data are collected (see "big data age"). Consequently, it is now commonplace to work with data with a very large number of features (i.e. $p >> n$).

- While $p$ can be extremely large, the number of observations $n$ is often limited due to cost, sample availability, or other considerations.

# Introduction

- Data containing more features than observations are typically referred to as **high-dimensional**.

- Issues pertaining to the *bias-variance tradeoff* and *overfitting* are commonly exacerbated in high dimensions.

- With a large number of features, statistical models (e.g. regression) can become too flexible and hence overfit the data.

- Recall the *curse of dimensionality*, which poses two fundamental, associated problems: (1) "neighborhoods" become very large (this is problematic in particular for kernel and clustering methods), (2) we need a much larger data set to adequately "fill" the space for predictive modeling, etc.

# Interpretability in High Dimensions

- In high-dimensional settings we need to be **cautious about how we interpret our results** – that is to say if they can be reasonably interpreted at all.

- Of course, it is oftentimes adequate, depending on the application, to treat a machine learning model as a *mere* predictive "black box" (e.g. statistical arbitrage, government work).

- Conversely, if we want to say that the features in our model directly impact the outcomes we observe (note: in ML we almost never use the c-word – viz., variables *caused* observed effect) we need to be alert to ***multicollinearity***.

- In high dimensions, it is very likely that some of our model variables are mutually correlated. This means we can never know exactly which variables (if any) are truly predictive of the outcome. Moreover, we can rarely identify the optimal set of features for a given phenomenon of interest.

# Interpretability in High Dimensions

- The "first rule" of data science and ML: **one can always add more and more features to achieve zero classification/predictive error**, a perfect correlation coefficient value, etc.

- In the end, however, ***this is a useless model***.  We always need to report results on an independent test or validation set.

- In 2008, Hinton *et al*, developed a non-linear dimensionality technique known as ***t-SNE*** (*t-distributed stochastic neighbor embedding*) that is particularly well-suited for embedding high-dimensional data into 2 or 3 dimensions, which can be visualized with a scatter plot.

- Specifically, it models each high-dimensional object by a two- or three-dimensional point in such a way that **similar objects are modeled by nearby points** and dissimilar objects are modeled by distant points.

# t-SNE: H-D Data Visualization

- First, t-SNE constructs a probability distribution over pairs of high-dimensional objects in such a way that similar objects have a high probability of being picked, whilst dissimilar points have an extremely small probability of being picked.

- Second, t-SNE defines a similar probability distribution over the points in the low-dimensional map, and it minimizes the **KL divergence** (a standard measure of "distance" between probability distributions) between the two distributions with respect to the locations of the points in the map.

$$D_{\mathrm{KL}}(P\|Q) = \sum_i P(i) \ln \frac{P(i)}{Q(i)}$$

$p(x)$  $q(x)$

$D_{KL}(P\|Q)$

$$D_{\mathrm{KL}}(P\|Q) = -\sum_x p(x) \log q(x) + \sum_x p(x) \log p(x)$$
$$= \quad H(P,Q) \quad - \quad H(P)$$

# t-SNE for MNIST

# t-SNE for Atari! (Deepmind)

# Word2vec (2013)

- **Word2vec** is a group of related models (Google) that are used to produce word embeddings.

- These models are shallow, two-layer neural network that are trained to reconstruct linguistic contexts of words.

- Word2vec takes as its input a large corpus of text and produces a vector space (usually of high dimensions), with each unique word in the corpus being assigned a corresponding vector in the space.

- Word vectors are positioned in the vector space such that **words that share common contexts in the corpus are located in close proximity to one another in the space**.



Load up the word vectors

QUEEN [0.3, 0.9]

KING [0.5, 0.7]

WOMAN [0.3, 0.4]

MAN [0.5, 0.2]



So king + man - woman = queen!

QUEEN

KING

WOMAN

MAN



Country and Capital Vectors Projected by PCA

# t-SNE for word2vec

Word2vec model computed from 6 billion word corpus of news articles

| Type of relationship | Word Pair 1 | | Word Pair 2 | |
|---|---|---|---|---|
| Common capital city | Athens | Greece | Oslo | Norway |
| All capital cities | Astana | Kazakhstan | Harare | Zimbabwe |
| Currency | Angola | kwanza | Iran | rial |
| City-in-state | Chicago | Illinois | Stockton | California |
| Man-Woman | brother | sister | grandson | granddaughter |
| Adjective to adverb | apparent | apparently | rapid | rapidly |
| Opposite | possibly | impossibly | ethical | unethical |
| Comparative | great | greater | tough | tougher |
| Superlative | easy | easiest | lucky | luckiest |
| Present Participle | think | thinking | read | reading |
| Nationality adjective | Switzerland | Swiss | Cambodia | Cambodian |
| Past tense | walking | walked | swimming | swam |
| Plural nouns | mouse | mice | dollar | dollars |
| Plural verbs | work | works | speak | speaks |

"Efficient Estimation of Word Representations in Vector Space" Tomas Mikolov, Kai Chen, Greg Corrado, Jeffrey Dean, Arxiv 2013

# Dimensionality Reduction

- In general: the higher the number of dimensions we have, the more training data we need.

- Additionally, <u>computational cost is generally an explicit function of dimensionality</u>.

- Dimensionality reduction can also **remove noise** in a data set, which can, in turn, significantly improve the results of a learning algorithm.

- These are perhaps the strongest reasons why dimensionality reduction is useful (in addition to improving visualization/interpretability).

In general, there are (3) common ways to perform dimensionality reduction:

(1) **Feature selection** – determine whether the features available are actually useful, i.e. are they correlated with the output variables.

(2) **Feature derivation** – means deriving new features from old ones, generally by applying transforms to the data set that change the coordinate system axes (e.g., by moving or rotating); this is usually achieved through matrix multiplication.

(3) **Clustering** – group together similar data points to see whether this allows fewer features to be used.

# LDA

- We will consider linear discriminant analysis (LDA) in an **unsupervised setting** (that is to say the data has no target label), due to Fisher (1936).



- Consider the case of **two data classes**; we can compute various summary statistics for these data, including μ, the mean of the entire set, as well as $\mu_1$ and $\mu_2$, the means of each class, respectively; and the covariance of each class: $\sum_j \left( x_j - \mu \right)\left( x_j - \mu \right)^T$



**\*Note**: We use the standard notation for column vectors, so that $\mathbf{w}^T\mathbf{w}$ is equivalent to the dot product of $\mathbf{w}$ with itself, whereas $\mathbf{w}\mathbf{w}^T$ <u>yields an nxn matrix</u>, for n-dimensional $\mathbf{w}$.

# LDA

- The principal insight of LDA is that the **covariance matrix** can tell us about the *scatter within a dataset*, which is the amount of spread extant within the data.

- The way to find the *scatter* is to multiply the covariance by $p_c$, **the probability of the given class** (that is, the number of data points there are in the class divided by the total number).

- Adding the values of this for all of the classes gives us a measure called the **within-class scatter of the data set**:

$$S_W = \sum_{classes\ c} \sum_{j \in c} p_c \left( x_j - \mu_c \right) \left( x_j - \mu_c \right)^T$$

# LDA

- **Within-class scatter of the data set:** $S_W = \sum\limits_{classes\ c} \sum\limits_{j \in c} p_c \left( x_j - \mu_c \right) \left( x_j - \mu_c \right)^T$

- If the data is **easy to separate into classes**, then this <u>within-class scatter should be small</u>, so that each class is tightly clustered together.

- Conversely, in order to separate data, we also want the **distance between classes to be large**. This quantity is known as the **between-classes scatter**:

$$S_B = \sum\limits_{classes\ c} \left( \mu_c - \mu \right) \left( \mu_c - \mu \right)^T$$

# LDA



$$S_W = \sum_{\text{classes } c} \sum_{j \in c} p_c \left( x_j - \mu_c \right) \left( x_j - \mu_c \right)^T \qquad S_B = \sum_{\text{classes } c} \left( \mu_c - \mu \right) \left( \mu_c - \mu \right)^T$$

# LDA

- "Good separation" means $S_W$ should be small and $S_B$ should be large; accordingly we wish to make the ratio: **$S_B/S_W$ as large as possible**.

- Regarding dimensionality reduction: we want **$S_B/S_W$** to be large when we reduce the number of dimensions of our data.

- Recall that the **projection** of a data point (i.e. a vector) **x** onto another vector **w** can be written as $z = \mathbf{w}^T\mathbf{x}$.

- So let's consider the between and within scatter quantities projected onto some vector **w**:

why?

$$\sum_{classes\ c}\sum_{j\in c} p_c\left(w^T\left(x_j - \mu_c\right)\right)\left(w^T\left(x_j - \mu_c\right)\right)^T = w^T S_W w$$

$$\sum_{classes\ c}\left(w^T\left(\mu_c - \mu\right)\right)\left(w^T\left(\mu_c - \mu\right)^T\right) = w^T S_B w$$

# LDA

- Thus, the ratio of projected between-class and within-class scatter is:

$$\frac{w^T S_B w}{w^T S_W w}$$

How do we find the maximum value of this expression with respect to **w**?

# LDA

- Thus, the ratio of projected between-class and within-class scatter is:

$$\frac{w^T S_B w}{w^T S_W w}$$

How do we find the maximum value of this expression with respect to **w**?

Take the derivative with respect to **w** and set it equal to zero (i.e. find the "critical points").

# LDA

- Thus, the ratio of projected between-class and within-class scatter is:

$$\frac{w^T S_B w}{w^T S_W w}$$

How do we find the maximum value of this expression with respect to **w**?

Take the derivative with respect to **w** and set it equal to zero (i.e. find the "critical points").

- This yields:

Why?

$$\frac{S_B w \left( w^T S_W w \right) - S_W w \left( w^T S_B w \right)}{\left( w^T S_W w \right)^2} = 0$$

# LDA

- Thus, the ratio of projected between-class and within-class scatter is:

$$\frac{w^T S_B w}{w^T S_W w}$$

How do we find the maximum value of this expression with respect to **w**?

<u>Take the derivative with respect to **w** and set it equal to zero</u> (i.e. find the "critical points").

- This yields:

Why?

$$\frac{S_B w \left( w^T S_W w \right) - S_W w \left( w^T S_B w \right)}{\left( w^T S_W w \right)^2} = 0$$

Calculus!

$$\left( \frac{u}{v} \right)' = \frac{u'v - uv'}{v^2}$$

**FYI**: Matrices have analogous "derivative" rules (see *Matrix Calculus*):
$$\frac{\partial}{\partial \mathbf{x}} \left[ \mathbf{x}^T A \mathbf{x} \right] = \mathbf{x}^T \left( A + A^T \right) \text{, etc.}$$

# LDA

- After some simplification and solving a sub-problem known as the generalized eigenvalue problem, we arrive at Fisher's solution in the 2-class case:

$$w^* = \arg\max\left[\frac{w^T S_B w}{w^T S_W w}\right] = S_W^{-1}\left(\mu_1 - \mu_2\right)$$

Let's recap the LDA framework and what w* represents, now that we have the solution:

(*) LDA seeks to <u>reduce dimensionality while preserving as much of the class discriminatory information as possible</u>.

(*) We seek to obtain a scalar y by projecting each datum $\mathbf{x}$ onto a line: $y = \mathbf{w}^T\mathbf{x}$.

(*) Of all possible lines, <u>w* represents the one that maximizes the separability of the scalars</u>.



FIGURE 6.5  Plot of the iris data showing the three classes *left:* before and right: after LDA has been applied.

# LDA

(*) In the figure: the two classes A and B appear overlapped along both $X_1$ and $X_2$ directions.

(*) However, they are perfectly separated along the discriminant function w*.

(*) Projecting the data onto this discriminant function renders a perfect separation of the two classes.



$$w* = \arg\max \left[ \frac{w^T S_B w}{w^T S_W w} \right] = S_W^{-1} \left( \mu_1 - \mu_2 \right)$$

* Note that there are a variety of different forms of LDA, including the use of Bayes' theorem for posterior classification, and QDA (quadratic discriminant analysis), where the discriminant function is parabolic.

# PCA

- Like LDA, **principal component analysis** (PCA) amounts to computing a transformation of a data set in order to identify a (useful) lower-dimensional set of axes.

*Essential idea*: PCA <u>generates a particular set of coordinate axes</u> (usually in fewer dimensions than the original data) **that capture the maximum variability in the data**; furthermore, <u>these new coordinate axes are **orthogonal**</u> (which is to say they are uncorrelated).

# PCA: (Aside) Gram-Schmidt

- One of the quintessential results in Linear Algebra is the **Gram-Schmidt algorithm** (*NB*: this is the same Gram per the "Gram matrix").

The G-S algorithm takes a basis for, say, an *inner product space* (viz., a vector space equipped with an inner product, such as the *Euclidean space* $R^n$) and returns an ***orthonormal basis*** for the same space.



Gram

# PCA: (Aside) Gram-Schmidt

- One of the quintessential results in Linear Algebra is the **Gram-Schmidt algorithm** (*NB*: this is the same Gram per the "Gram matrix").

The G-S algorithm takes a basis for, say, an *inner product space* (viz., a vector space equipped with an inner product, such as the *Euclidean space* $R^n$) and returns an ***orthonormal basis*** for the same space.

Gram

(*) Given a basis set, each step of the G-S algorithm amounts to <u>iteratively subtracting off the orthogonal projection of the current vector from each of the previous vectors in the process</u>.

(*) In the diagram, the set $\{\mathbf{v_1}, \mathbf{v_2}\}$ is the original basis set;
(1) Normalize $\mathbf{v_1}$:  $e_1 = v_1/\|v_1\|$
(2) Project $\mathbf{v_2}$ onto $\mathbf{v_1}$
(3) Define $\mathbf{u_2}$ as the difference of $\mathbf{v_2}$ and the projection of $\mathbf{v_2}$ onto $\mathbf{v_1}$
(4) Normalize $\mathbf{u_2}$: $e_2 = u_2/\|u_2\|$    (5) the set $\{e_1, e_2\}$ is an **orthonormal basis**.

# PCA: (Aside) Gram-Schmidt

The G-S algorithm takes a basis for, say, an *inner product space* (viz., a vector space equipped with an inner product, such as the *Euclidean space* $R^n$) and returns an **orthonormal basis** for the same space.

$$\mathbf{u}_1 = \mathbf{v}_1,$$

$$\mathbf{u}_2 = \mathbf{v}_2 - \text{proj}_{\mathbf{u}_1}(\mathbf{v}_2),$$

$$\mathbf{u}_3 = \mathbf{v}_3 - \text{proj}_{\mathbf{u}_1}(\mathbf{v}_3) - \text{proj}_{\mathbf{u}_2}(\mathbf{v}_3),$$

$$\mathbf{u}_4 = \mathbf{v}_4 - \text{proj}_{\mathbf{u}_1}(\mathbf{v}_4) - \text{proj}_{\mathbf{u}_2}(\mathbf{v}_4) - \text{proj}_{\mathbf{u}_3}(\mathbf{v}_4),$$

$$\vdots$$

$$\mathbf{u}_k = \mathbf{v}_k - \sum_{j=1}^{k-1} \text{proj}_{\mathbf{u}_j}(\mathbf{v}_k),$$

$$\mathbf{e}_1 = \frac{\mathbf{u}_1}{\|\mathbf{u}_1\|}$$

$$\mathbf{e}_2 = \frac{\mathbf{u}_2}{\|\mathbf{u}_2\|}$$

$$\mathbf{e}_3 = \frac{\mathbf{u}_3}{\|\mathbf{u}_3\|}$$

$$\mathbf{e}_4 = \frac{\mathbf{u}_4}{\|\mathbf{u}_4\|}$$

$$\vdots$$

$$\mathbf{e}_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}.$$

# PCA: (Aside) Gram-Schmidt

The G-S algorithm takes a basis for, say, an *inner product space* (viz., a vector space equipped with an inner product, such as the *Euclidean space* $R^n$) and returns an **orthonormal basis** for the same space.

$$\mathbf{u}_1 = \mathbf{v}_1, \qquad\qquad \mathbf{e}_1 = \frac{\mathbf{u}_1}{\|\mathbf{u}_1\|}$$

$$\mathbf{u}_2 = \mathbf{v}_2 - \mathrm{proj}_{\mathbf{u}_1}(\mathbf{v}_2), \qquad \mathbf{e}_2 = \frac{\mathbf{u}_2}{\|\mathbf{u}_2\|}$$

$$\mathbf{u}_3 = \mathbf{v}_3 - \mathrm{proj}_{\mathbf{u}_1}(\mathbf{v}_3) - \mathrm{proj}_{\mathbf{u}_2}(\mathbf{v}_3), \qquad \mathbf{e}_3 = \frac{\mathbf{u}_3}{\|\mathbf{u}_3\|}$$

$$\mathbf{u}_4 = \mathbf{v}_4 - \mathrm{proj}_{\mathbf{u}_1}(\mathbf{v}_4) - \mathrm{proj}_{\mathbf{u}_2}(\mathbf{v}_4) - \mathrm{proj}_{\mathbf{u}_3}(\mathbf{v}_4), \quad \mathbf{e}_4 = \frac{\mathbf{u}_4}{\|\mathbf{u}_4\|}$$

$$\vdots \qquad\qquad\qquad \vdots$$

$$\mathbf{u}_k = \mathbf{v}_k - \sum_{j=1}^{k-1} \mathrm{proj}_{\mathbf{u}_j}(\mathbf{v}_k), \qquad \mathbf{e}_k = \frac{\mathbf{u}_k}{\|\mathbf{u}_k\|}.$$

(*) How is this useful? G-S algorithm is used to render the **QR factorization** of a matrix.

(*) PCA employs a similar strategy to generate the "principal components" of a data set.

# PCA: (Aside) Gram-Schmidt

$\mathbf{u}_1 = \mathbf{v}_1,$

$\mathbf{u}_2 = \mathbf{v}_2 - \mathrm{proj}_{\mathbf{u}_1}(\mathbf{v}_2),$

$\mathbf{u}_3 = \mathbf{v}_3 - \mathrm{proj}_{\mathbf{u}_1}(\mathbf{v}_3) - \mathrm{proj}_{\mathbf{u}_2}(\mathbf{v}_3),$

$\mathbf{u}_4 = \mathbf{v}_4 - \mathrm{proj}_{\mathbf{u}_1}(\mathbf{v}_4) - \mathrm{proj}_{\mathbf{u}_2}(\mathbf{v}_4) - \mathrm{proj}_{\mathbf{u}_3}(\mathbf{v}_4),$

$\vdots$

$\mathbf{u}_k = \mathbf{v}_k - \sum_{j=1}^{k-1} \mathrm{proj}_{\mathbf{u}_j}(\mathbf{v}_k),$

$\mathbf{e}_1 = \dfrac{\mathbf{u}_1}{\|\mathbf{u}_1\|}$

$\mathbf{e}_2 = \dfrac{\mathbf{u}_2}{\|\mathbf{u}_2\|}$

$\mathbf{e}_3 = \dfrac{\mathbf{u}_3}{\|\mathbf{u}_3\|}$

$\mathbf{e}_4 = \dfrac{\mathbf{u}_4}{\|\mathbf{u}_4\|}$

$\vdots$

$\mathbf{e}_k = \dfrac{\mathbf{u}_k}{\|\mathbf{u}_k\|}.$

Apply the Gram-Schmidt process to the following basis.

$$\begin{array}{cccc} & \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 \\ B = & \{(1,1,0), & (1,2,0), & (0,1,2)\} \end{array}$$

Sol: $\mathbf{v}_1 = \mathbf{u}_1 = (1,1,0)$

$\mathbf{v}_2 = \mathbf{u}_2 - \dfrac{\mathbf{u}_2 \cdot \mathbf{v}_1}{\mathbf{v}_1 \cdot \mathbf{v}_1}\mathbf{v}_1 = (1,2,0) - \dfrac{3}{2}(1,1,0) = (-\dfrac{1}{2}, \dfrac{1}{2}, 0)$

$\mathbf{v}_3 = \mathbf{u}_3 - \dfrac{\mathbf{u}_3 \cdot \mathbf{v}_1}{\mathbf{v}_1 \cdot \mathbf{v}_1}\mathbf{v}_1 - \dfrac{\mathbf{u}_3 \cdot \mathbf{v}_2}{\mathbf{v}_2 \cdot \mathbf{v}_2}\mathbf{v}_2$

$= (0,1,2) - \dfrac{1}{2}(1,1,0) - \dfrac{1/2}{1/2}(-\dfrac{1}{2}, \dfrac{1}{2}, 0) = (0,0,2)$

# PCA

PCA generates a particular set of coordinate axes **that capture the maximum variability in the data**; furthermore, these new coordinate axes are **orthogonal.**

The figure shows two versions of the same data set.



- In the **first image**, the data are arranged in an ellipse that runs at 45° axes; while in **the second,** the axes have been moved so that the data now runs along the x-axis and is centered on the origin.

- Key idea: the potential for dimensionality reduction rests in the fact that the $y$ dimension now does not demonstrate much variability – and so it might be possible to ignore it and simply use the $x$ axis values alone for learning, etc.

(*) In fact, applying this dimensionality reduction often has the nice effect of removing some of the noise in the data.

# PCA

Q: How do we choose the new axes?

A: With PCA, the *principal component* is <u>the direction in the data with the largest variation</u>.

- The algorithm first **centers the data** (by subtracting off the mean), and then **chooses the direction with the largest variation** and places an axis in that direction.

- For subsequent steps, the algorithm then looks at the variation that remains and **finds another axis that is orthogonal to the first and covers as much of the remaining variation as possible**.

  Rinse and repeat…

# PCA

Q: How do we choose the new axes?

A: With PCA, the *principal component* is <u>the direction in the data with the largest variation</u>.

- The algorithm first **centers the data** (by subtracting off the mean), and then **chooses the direction with the largest variation** and places an axis in that direction.

- For subsequent steps, the algorithm then looks at the variation that remains and **finds another axis that is orthogonal to the first and covers as much of the remaining variation as possible**.

  Rinse and repeat…

(*) The end result is that all the variation is along the axes of the coordinate set, and so the **<u>covariance matrix of the transformed data is diagonal</u>** (since each new variable is uncorrelated with every variable except itself).

(*) Because some of the axes generated in this process have very little variation, **we can typically remove them** without drastically affecting the variability in the data (note the implicit assumption: variation in data features equates to useful information for classification/inference).

# PCA

Let's formally work out the PCA algorithm.

- Suppose we have a data matrix **X** of dimension **n by m** (*n* is the number of training instances, *m* is the dimension of each datum).

**<u>Goal for PCA</u>**: rotate the data so that we render a coordinate system so the new axes are uncorrelated (i.e. orthogonal) and can be ranked according to maximum variation.

# PCA

Let's formally work out the PCA algorithm.

- Suppose we have a data matrix **X** of dimension **n by d** (*n* is the number of training instances, *d* is the dimension of each datum).

**Goal for PCA**: rotate the data so that we render a coordinate system so the new axes are uncorrelated (i.e. orthogonal) and can be ranked according to maximum variation.

Using standard techniques from linear algebra, we can express the <u>rotation of the data matrix X</u> as:

$$Y = P^T X$$

Where P is a "rotation matrix" with the natural property that $P^T = P^{-1}$ (why?).

$$\begin{pmatrix} \cos 45° & -\sin 45° \\ \sin 45° & \cos 45° \end{pmatrix}$$

**Counter Clockwise** $\begin{pmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{pmatrix}$

my graph

45 degrees

**Clockwise** $\begin{pmatrix} \cos \phi & \sin \phi \\ -\sin \phi & \cos \phi \end{pmatrix}$

Simple rotation matrix example

# PCA

$$Y = P^T X$$

- Recall that for PCA, <u>we want the covariance of the transformed matrix $Y$ to be diagonal</u>.

- That is to say, we want to find a matrix $P$ where:

$$\text{cov}(Y) = \text{cov}(P^T X) = \begin{pmatrix} \lambda_1 & 0 & \ldots & 0 \\ 0 & \lambda_2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix}$$

# PCA

$$Y = P^T X$$

- Recall that for PCA, <u>we want the covariance of the transformed matrix $Y$ to be diagonal</u>.

- That is to say, we want to find a matrix $P$ where:

$$\text{cov}(Y) = \text{cov}(P^T X) = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix}$$

- Let's relate cov(Y) and cov(X), so that we can more easily solve for P.

$$\text{cov}(Y) = E\left[ YY^T \right]$$

(By definition of covariance; recall the data is centered)

# PCA

$$Y = P^T X$$

- Recall that for PCA, <u>we want the covariance of the transformed matrix $Y$ to be diagonal</u>.

- That is to say, we want to find a matrix $P$ where:

$$\mathrm{cov}(Y) = \mathrm{cov}(P^T X) = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix}$$

- Let's relate cov(Y) and cov(X), so that we can more easily solve for P.

$$\mathrm{cov}(Y) = E\left[YY^T\right]$$

$$= E\left[\left(P^T X\right)\left(P^T X\right)^T\right]$$

# PCA

$$Y = P^T X$$

- Recall that for PCA, <u>we want the covariance of the transformed matrix $Y$ to be diagonal</u>.

- That is to say, we want to find a matrix $P$ where:

$$\text{cov}(Y) = \text{cov}(P^T X) = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix}$$

- Let's relate cov(Y) and cov(X), so that we can more easily solve for P.

$$\text{cov}(Y) = E\left[ YY^T \right]$$

$$= E\left[ \left( P^T X \right)\left( P^T X \right)^T \right]$$

$$= E\left[ \left( P^T X \right)\left( X^T P \right) \right] \quad \text{Why?}$$

# PCA

$$Y = P^T X$$

- Recall that for PCA, <u>we want the covariance of the transformed matrix $Y$ to be diagonal</u>.

- That is to say, we want to find a matrix $P$ where:

$$\mathrm{cov}(Y) = \mathrm{cov}(P^T X) = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix}$$

- Let's relate cov(Y) and cov(X), so that we can more easily solve for P.

$$\mathrm{cov}(Y) = E\left[YY^T\right]$$

$$= E\left[\left(P^T X\right)\left(P^T X\right)^T\right]$$

$$= E\left[\left(P^T X\right)\left(X^T P\right)\right]$$

$$= P^T E\left[XX^T\right] P$$

$$= P^T \mathrm{cov}(X) P$$

Why?

# PCA

- In summary:  $\mathrm{cov}(Y) = P^T \, \mathrm{cov}\left(X\right) P$

- This tells us that:

$$P\,\mathrm{cov}(Y) = PP^T \, \mathrm{cov}\left(X\right) P = \mathrm{cov}\left(X\right) P$$

Why?

# PCA

- In summary:  $\mathrm{cov}(Y) = P^T \,\mathrm{cov}(X)P$

- This tells us that:

$$P\,\mathrm{cov}(Y) = PP^T \,\mathrm{cov}(X)P = \mathrm{cov}(X)P$$

$$P^T = P^{-1}$$

(\*) Remember that **cov(Y) is a diagonal matrix**; if we write the matrix P as a set of column vectors: $P = [\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_N}]$, then we have:

$$P\,\mathrm{cov}(Y) = [\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_N}]\begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix} = [\lambda_1\mathbf{p_1}, \lambda_2\mathbf{p_2}, ..., \lambda_N\mathbf{p_N}]$$

# PCA

- In summary: $\mathrm{cov}(Y) = P^T \, \mathrm{cov}(X) P$ and

$$P \, \mathrm{cov}(Y) = [\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_N}] \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix} = [\lambda_1 \mathbf{p_1}, \lambda_2 \mathbf{p_2}, ..., \lambda_N \mathbf{p_N}]$$

(*) If we write $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, ..., \lambda_N)^T$ *and* $\mathbf{Z} = \mathrm{cov}(\mathbf{X})$, we then arrive at the fundamental equation:

# PCA

- In summary: $\mathrm{cov}(Y) = P^T \mathrm{cov}(X) P$ and

$$P\,\mathrm{cov}(Y) = [\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_N}] \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix} = [\lambda_1\mathbf{p_1}, \lambda_2\mathbf{p_2}, ..., \lambda_N\mathbf{p_N}]$$

(*) If we write $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, ..., \lambda_N)^T$ *and* $\mathbf{Z} = \mathrm{cov}(\mathbf{X})$, we then arrive at the fundamental equation:

$$\boxed{\boldsymbol{\lambda}\mathbf{p_i} = \mathbf{Z}\mathbf{p_i} \ \textit{for each } \mathbf{p_i}}$$

$$P\,\mathrm{cov}(Y) = \mathrm{cov}(X)P$$

# PCA

- In summary: $\operatorname{cov}(Y) = P^T \operatorname{cov}(X) P$ and

$$P\operatorname{cov}(Y) = [\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_N}] \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix} = [\lambda_1 \mathbf{p_1}, \lambda_2 \mathbf{p_2}, ..., \lambda_N \mathbf{p_N}]$$

(*) If we write $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, ..., \lambda_N)^T$ *and* $\mathbf{Z} = \operatorname{cov}(\mathbf{X})$, we then arrive at the fundamental equation:

$$\boldsymbol{\lambda} \mathbf{p_i} = \mathbf{Z} \mathbf{p_i} \ for \ each \ \mathbf{p_i}$$

This is the *i*th column of cov(X)P

This is the *i*th column of Pcov(Y)

# PCA

- In summary: $\mathrm{cov}(Y) = P^T \mathrm{cov}(X) P$ and

$$P\,\mathrm{cov}(Y) = [\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_N}] \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix} = [\lambda_1 \mathbf{p_1}, \lambda_2 \mathbf{p_2}, ..., \lambda_N \mathbf{p_N}]$$

(*) If we write $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, ..., \lambda_N)^T$ *and* $\mathbf{Z} = \mathrm{cov}(\mathbf{X})$, we then arrive at the fundamental
equation:

$$\boldsymbol{\lambda}\mathbf{p_i} = \mathbf{Z}\mathbf{p_i} \; \textit{for each} \; \mathbf{p_i}$$

This is the *i*th column of cov(X)P

This is the *i*th column of Pcov(Y)

(*) Do you recognize anything "special" about the structure of these vectors?

# PCA

- In summary: $\mathrm{cov}(Y) = P^T \, \mathrm{cov}(X) P$ and

$$P\,\mathrm{cov}(Y) = [\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_N}] \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix} = [\lambda_1 \mathbf{p_1}, \lambda_2 \mathbf{p_2}, ..., \lambda_N \mathbf{p_N}]$$

(*) If we write $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, ..., \lambda_N)^T$ *and* $\mathbf{Z} = \mathrm{cov}(\mathbf{X})$, we then arrive at the fundamental equation:

$$\boldsymbol{\lambda}\mathbf{p_i} = \mathbf{Z}\mathbf{p_i} \; for \; each \; \mathbf{p_i}$$

The $\mathbf{p_i}$ vectors are **<u>eigenvectors</u>** of $\mathbf{Z}$=cov($\mathbf{X}$)!

# PCA

- In summary: we need to solve $\lambda \mathbf{p_i} = \mathbf{Z} \mathbf{p_i}$ *for each* $\mathbf{p_i}$, where each $\mathbf{p_i}$ vector is an eigenvector of $\mathbf{Z} = \mathrm{cov}(\mathbf{X})$.

(\*) **Useful fact (1)**: all eigenvectors of a square <u>symmetric matrix A are orthogonal to each other</u>.

# PCA

- In summary: we need to solve $\lambda \mathbf{p_i} = \mathbf{Z} \mathbf{p_i}$ *for each* $\mathbf{p_i}$, where each $\mathbf{p_i}$ vector is an eigenvector of $\mathbf{Z}$=cov($\mathbf{X}$).

(*) **Useful fact (1)**: all eigenvectors of a square <u>symmetric matrix A are orthogonal to each other</u>.

(*) **Useful fact (2)**: <u>$\mathbf{Z}$=cov($\mathbf{X}$) is a square symmetric matrix</u> (as is true for *any* covariance matrix).

# PCA

- In summary: we need to solve $\lambda \mathbf{p_i} = \mathbf{Z} \mathbf{p_i}$ *for each* $\mathbf{p_i}$, where each $\mathbf{p_i}$ vector is an eigenvector of $\mathbf{Z} = \text{cov}(\mathbf{X})$.

(*) **Useful fact (1)**: all eigenvectors of a square <u>symmetric matrix A are orthogonal to each other</u>.

(*) **Useful fact (2)**: <u>$\mathbf{Z} = \text{cov}(\mathbf{X})$ is a square symmetric matrix</u> (as is true for *any* covariance matrix).

(*) If we put (1) and (2) together, this means that the solution to the PCA problem (i.e. solving for the matrix P) boils down to determining the ***eigendecomposition*** (also called the *spectral decomposition*) – which is guaranteed to exist – of $\text{cov}(\mathbf{X})$.

$$\mathbf{Z} = \text{cov}(\mathbf{X}) = \mathbf{E}\mathbf{D}\mathbf{E}^{T}$$

Where $\mathbf{D}$ is the diagonal matrix of eigenvalues for $\mathbf{Z}$, and $\mathbf{E}$ is the corresponding (orthogonal) matrix of eigenvectors.

# PCA

$$\mathbf{Z} = \mathrm{cov}(\mathbf{X}) = \mathbf{EDE}^{T}$$

- Note: In the eigendecomposition for cov($\mathbf{X}$), the <u>dimensions with large eigenvalues have lots of variation and are therefore useful dimensions</u>.

- In order to perform a **dimensionality reduction** on our data set, we can therefore <u>throw away dimensions for which the eigenvalues are very small</u> (usually smaller than some chosen parameter).

# PCA

- Here is the **<u>PCA algorithm</u>**:

(1) Write $N$ data points $x_i = (x_{1i}, x_{2i}, \dots, x_{Mi})$ as row vectors.

(2) Put these vectors into the data matrix $\mathbf{X}$ (of size $N$ x $M$).

(3) Center the data by subtracting off the mean of each column, place into matrix $\mathbf{B}$.

(4) Computer the covariance matrix: $\mathbf{C} = \dfrac{1}{N}\mathbf{BB}^T$

(5) Computer the *eigenvalues* and *eigenvectors* of $\mathbf{C}$, so: $\mathbf{C} = \mathbf{VDV}^T$
where $\mathbf{D}$ is the diagonal matrix of eigenvalues; V is the matrix of corresponding eigenvectors.

(6) <u>Sort of the columns of D</u> into order of decreasing eigenvalues, and apply the same order to the columns of V.

(7) Reject those with eigenvalues less than some given threshold, leaving $L$ dimensions in the data.

# PCA for MNIST



MNIST Dataset reduced to 2 Components using PCA

# PCA vs. LDA



**PCA:**
component axes that maximize the variance

**LDA:**
maximizing the component axes for class-separation

bad projection

good projection: separates classes well

# Extending PCA

Q: What strong assumptions did we make about the surface for the directions of maximum variation with PCA?

# Extending PCA

Q: What strong assumptions did we make about the surface for the directions of maximum variation with PCA?

A: We assumed these surfaces of maximum variation are **straight lines** (this is a strong assumption!)

Q: How can are "break" the linear restriction for PCA?

# Extending PCA

Q: What strong assumptions did we make about the surface for the directions of maximum variation with PCA?

A: We assumed these surfaces of maximum variation are **straight lines** (this is a strong assumption!)

Q: How can are break the linear restriction for PCA?

A: **"Kernelize" PCA**!



**Fig. 1.** Basic idea of kernel PCA: by using a nonlinear kernel function $k$ instead of the standard dot product, we implicitly perform PCA in a possibly high–dimensional space $F$ which is nonlinearly related to input space. The dotted lines are contour lines of constant feature value.

# Kernel PCA

(*) All we have to do is express the covariance matrix **C** (recall this was the covariance of the data matrix X after centering) in terms of a kernel transformation:

$$\mathbf{C} = \frac{1}{N} \sum_{i=1}^{N} \Phi(x_n) \cdot \Phi(x_n)^T$$

(*) Next we compute the *eigendecomposition* of **C** and use the eigenvectors with the largest associated eigenvalues for PCA.

(*) Recall (from SVM lecture) that by using a kernel function we implicitly perform a dot product in a larger dimensional feature space (this is the crux of the **kernel trick**), with the upshot of enhanced expressiveness.

# PCA & Auto-encoders

• Note that PCA is intimately connected with MLPs.

(*) An MLP can perform (non-linear) PCA using what is called an *auto-associator* (more commonly: **auto-encoder**).

(*) If we train the MLP **where the output equals the input**, we are asking the network to <u>learn a data "reconstruction" process</u>; we therefore train to *minimize the reconstruction error*.

(*) Usually the hidden layers are smaller in dimension than the output/input layers so that they form a **compression "bottleneck"**.

(*) The activations at the hidden layers (i.e. the feature vectors) encode a *dimensionality reduction* of the data.



"bottleneck" hidden layer

input layer

output layer
(reconstruction of input layer)

all layers are fully connected but not drawn

# PCA & Auto-encoders: Image Denoising



The image shows how a "denoising" autoencoder may be used to generate correct input from corrupted input.



corrupt input    cleaned input

# SVD

- Of the vast array of different matrix factorizations used in applied math, the **singular value decomposition** (SVD) is one of the most common and useful (we show its conceptual connection with PCA in the subsequent slides).

- Recall that a _symmetric matrix A_ admits of an eigendecomposition. Notably, if A is **not** symmetric, then there is no guarantee that it has an eigendecomposition – so not every matrix can be factored in this way.

When it came to our previous discussion of PCA, we had $A=XX^T$.

Now, is $A$ symmetric?

# SVD

- Of the vast array of different matrix factorizations used in applied math, the **singular value decomposition** (SVD) is one of the most common and useful (we show its conceptual connection with PCA in the subsequent slides).

- Recall that a *symmetric matrix $\boldsymbol{A}$* admits of an eigendecomposition. Notably, if $\mathbf{A}$ is **not** symmetric, then there is no guarantee that it has an eigendecomposition – so not every matrix can be factored in this way.

When it came to our previous discussion of PCA, we had $\mathbf{A}=\mathbf{X}\mathbf{X}^{\mathrm{T}}$.

Now, is $A$ symmetric? **Yes**: $\mathbf{A}^{\mathrm{T}}=(\mathbf{X}\mathbf{X}^{\mathrm{T}})^{\mathrm{T}}=(\mathbf{X}^{\mathrm{T}})^{\mathrm{T}}\mathbf{X}^{\mathrm{T}}=\mathbf{X}\mathbf{X}^{\mathrm{T}}=\mathbf{A}$.

The PCA algorithm was crucially reliant on the fact that $\mathbf{A}$ was symmetric and thus it has an eigendecomposition.

# SVD

- Conversely, *every* matrix has a singular value decomposition!

Definition: For an *m x n* matrix **A**, the *singular values* of **A** are the square roots of the eigenvalues of $\mathbf{A}^T\mathbf{A}$. They are denoted:

$$\sigma_1, ..., \sigma_n$$

It is conventional to arrange the singular values in decreasing order, whence: $\sigma_1 \geq \sigma_2 \geq ... \geq \sigma_n$

**Example**:

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

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

has eigenvalues $\lambda_1 = 3$ and $\lambda_2 = 1$. Consequently, the singular values of **A** are:

$$\sigma_1 = \sqrt{\lambda_1} = \sqrt{3}$$
$$\sigma_2 = \sqrt{\lambda_2} = 1$$

# SVD

- Definition: Let A be an *m x n* matrix with singular values, $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_r > 0$ and $\sigma_{r+1} = \sigma_{r+2} = \cdots = \sigma_n = 0$. Then there exist an *m x n* orthogonal matrix **U**, and *n x n* orthogonal matrix **V**, and an *m x n diagonal* matrix $\Sigma$ of the form:

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

Note: the columns of **U** are called *left singular vectors* of **A**, and the columns of **V** are called *right singular vectors* of **A**. The matrices **U** and **V** are <u>not</u> uniquely determined by **A**.

(*) *NB*: **rank(A) = r**.

# SVD

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



Example:

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

$$A^T A = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

# SVD

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



Example:

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

$$A^T A = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$
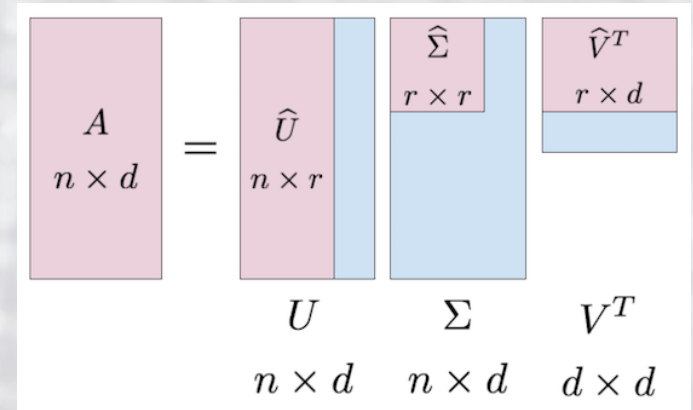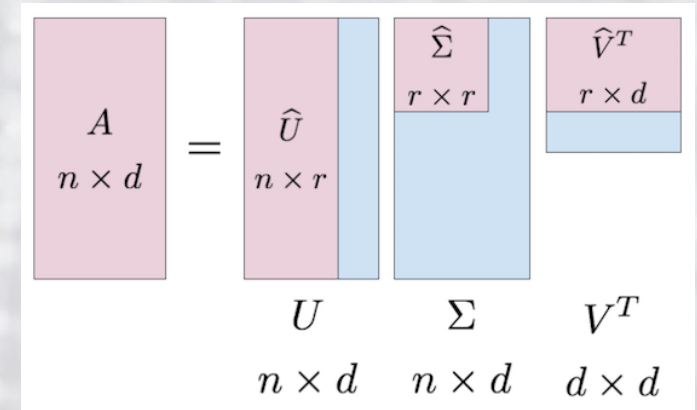
$$eigenvalues\left(A^T A\right): \lambda_1 = 2, \lambda_2 = 1, \lambda_3 = 0$$

# SVD

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



$$\begin{matrix} A \\ n \times d \end{matrix} = \begin{matrix} \widehat{U} \\ n \times r \end{matrix} \quad \begin{matrix} \widehat{\Sigma} \\ r \times r \end{matrix} \quad \begin{matrix} \widehat{V}^T \\ r \times d \end{matrix}$$

$$\begin{matrix} U \\ n \times d \end{matrix} \quad \begin{matrix} \Sigma \\ n \times d \end{matrix} \quad \begin{matrix} V^T \\ d \times d \end{matrix}$$

<u>Example</u>:

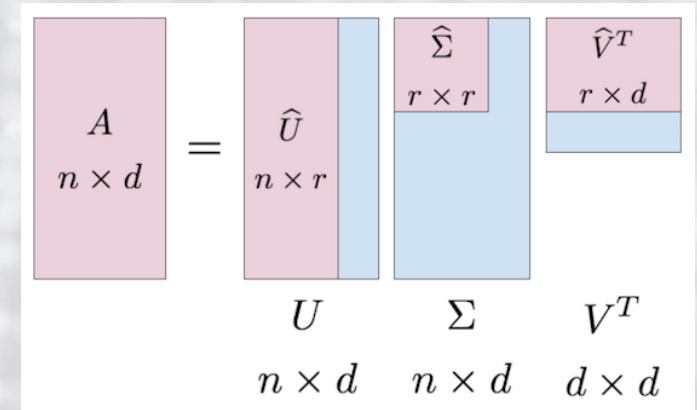$$A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$A^T A = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$eigenvalues\left(A^T A\right) : \lambda_1 = 2, \lambda_2 = 1, \lambda_3 = 0 \qquad eigenvectors\left(A^T A\right) : \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}$$

These vectors are orthogonal, so now we normalize them:

$$V = \begin{bmatrix} 1/\sqrt{2} & 0 & -1/\sqrt{2} \\ 1/\sqrt{2} & 0 & 1/\sqrt{2} \\ 0 & 1 & 0 \end{bmatrix}, \Sigma = \begin{bmatrix} \sqrt{2} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

# SVD



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

Example:

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

$eigenvalues\left(A^T A\right): \lambda_1 = 2, \lambda_2 = 1, \lambda_3 = 0$

$eigenvectors\left(A^T A\right): \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}$

$$V = \begin{bmatrix} 1/\sqrt{2} & 0 & -1/\sqrt{2} \\ 1/\sqrt{2} & 0 & 1/\sqrt{2} \\ 0 & 1 & 0 \end{bmatrix}, \Sigma = \begin{bmatrix} \sqrt{2} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

To find **U** we compute:

$$u_1 = \frac{1}{\sigma_1} A v_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad u_2 = \frac{1}{\sigma_2} A v_2 = \frac{1}{1} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

# SVD

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



Example:

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \sqrt{2} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \\ -1/\sqrt{2} & 1/\sqrt{2} & 0 \end{bmatrix} = U\Sigma V^T$$
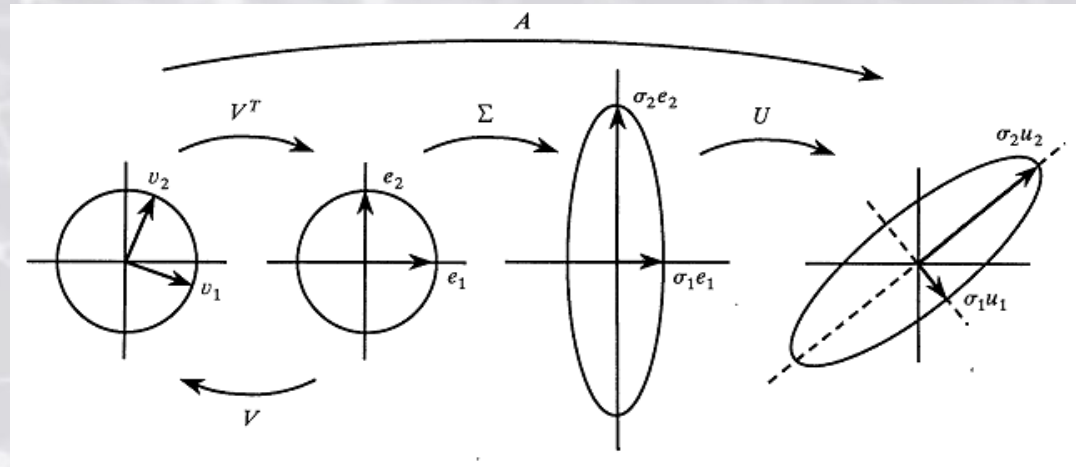
# SVD

Geometric Interpretation: In general, $\mathbf{\Sigma}$ can be regarded as a scaling matrix, and $\mathbf{U}$, $\mathbf{V}$ can be viewed as rotation matrices.

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

Thus the expression $\mathbf{U\Sigma V}$ can be intuitively interpreted as a **composition of three successive geometrical transformations**: a rotation or reflection, a scaling and another rotation or reflection.



As shown in the figure, the singular values can be interpreted as the **semiaxes of an ellipse in 2D**. This concept can be generalized to *n*-dimensional *Euclidean space*, with the singular values of any $n \times n$ square matrix being viewed as the semiaxes of an *n*-dimensional ellipsoid.

As in PCA, these coordinate axes provide a natural framework for determining a dimensionality reduction scheme that captures maximal variation.

# SVD: Outer Product Form

• SVD factorization yields a useful method for "low rank" approximations/dimensionality reduction of data.

Theorem: For a given SVD decomposition of an *m x n* matrix **A**, we can express **A** in the so-called **outer product form**:

$$\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + ... + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$$

Where $\sigma_1 \geq \sigma_2 \geq … \geq \sigma_r > 0$ denote the singular values of **A**; **u** and **v** are the corresponding *left singular* and *right singular vectors*.

(*) Note that the **condition number** of a matrix **A** is defined as the ratio of the largest and the smallest singular values of **A**. Matrices with large condition numbers are called **ill-conditioned** (this has a significant impact on the stability of many different kinds of numerical algorithms in linear algebra).

$$cond(A) = \frac{\sigma_{max}}{\sigma_{min}}$$

# SVD: Outer Product Form

$$\mathbf{A} = \sigma_1\mathbf{u}_1\mathbf{v}_1^T + \ldots + \sigma_r\mathbf{u}_r\mathbf{v}_r^T$$
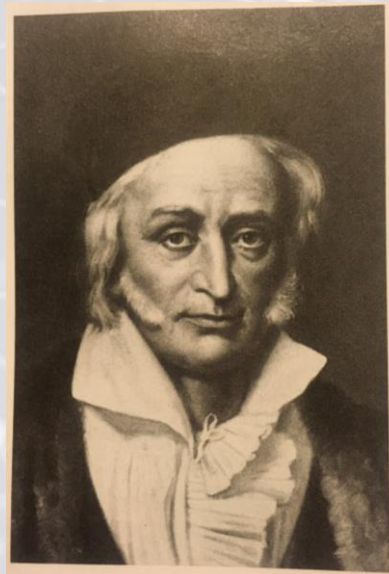
Example:

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \sqrt{2} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 0 & 0 & 1 \\ -1/\sqrt{2} & 1/\sqrt{2} & 0 \end{bmatrix} = U\Sigma V^T$$

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \sqrt{2}\begin{bmatrix} 1 \\ 0 \end{bmatrix}\begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} & 0 \end{bmatrix} + 1\begin{bmatrix} 0 \\ 1 \end{bmatrix}\begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$$

# SVD: Outer Product Form for Image Compression

• Consider the task of compressing a grayscale image of dimension *340 x 280*; each pixel is in the range [0, 255].

• We can store this image in a *340 x 280* dimension matrix, but transmitting and manipulating these 95,200 numbers is very expensive.

• Let's use <u>SVD for efficient image compression</u>. Recall that the small singular values in the SVD of a matrix correspond with "less informative" data features.
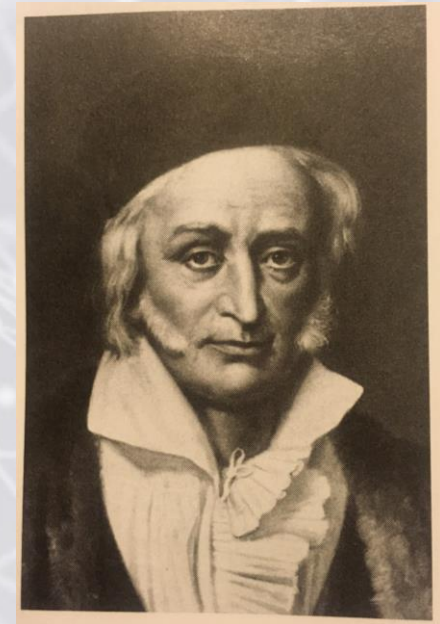
# SVD: Outer Product Form for Image Compression

- Suppose we have the SVD of A expressed in outer product form:

$$\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \ldots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$$

- For the original *340 x 280* image shown, we have r = 280 (why?).

- Define: $\mathbf{A}_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \ldots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T, \;\; k \leq r$

  as the *k-rank approximation to **A***.

# SVD: Outer Product Form for Image Compression
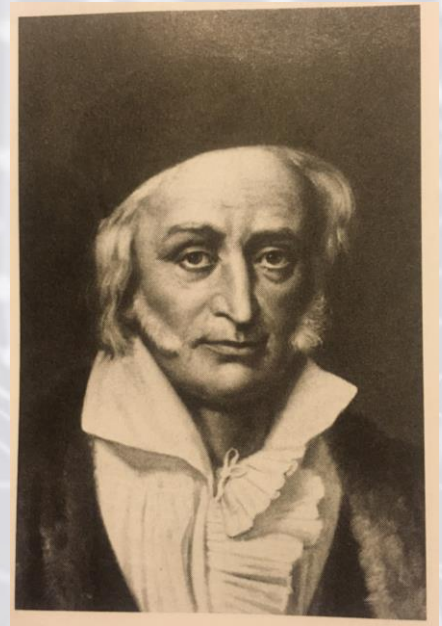
- Suppose we have the SVD of A expressed in outer product form:

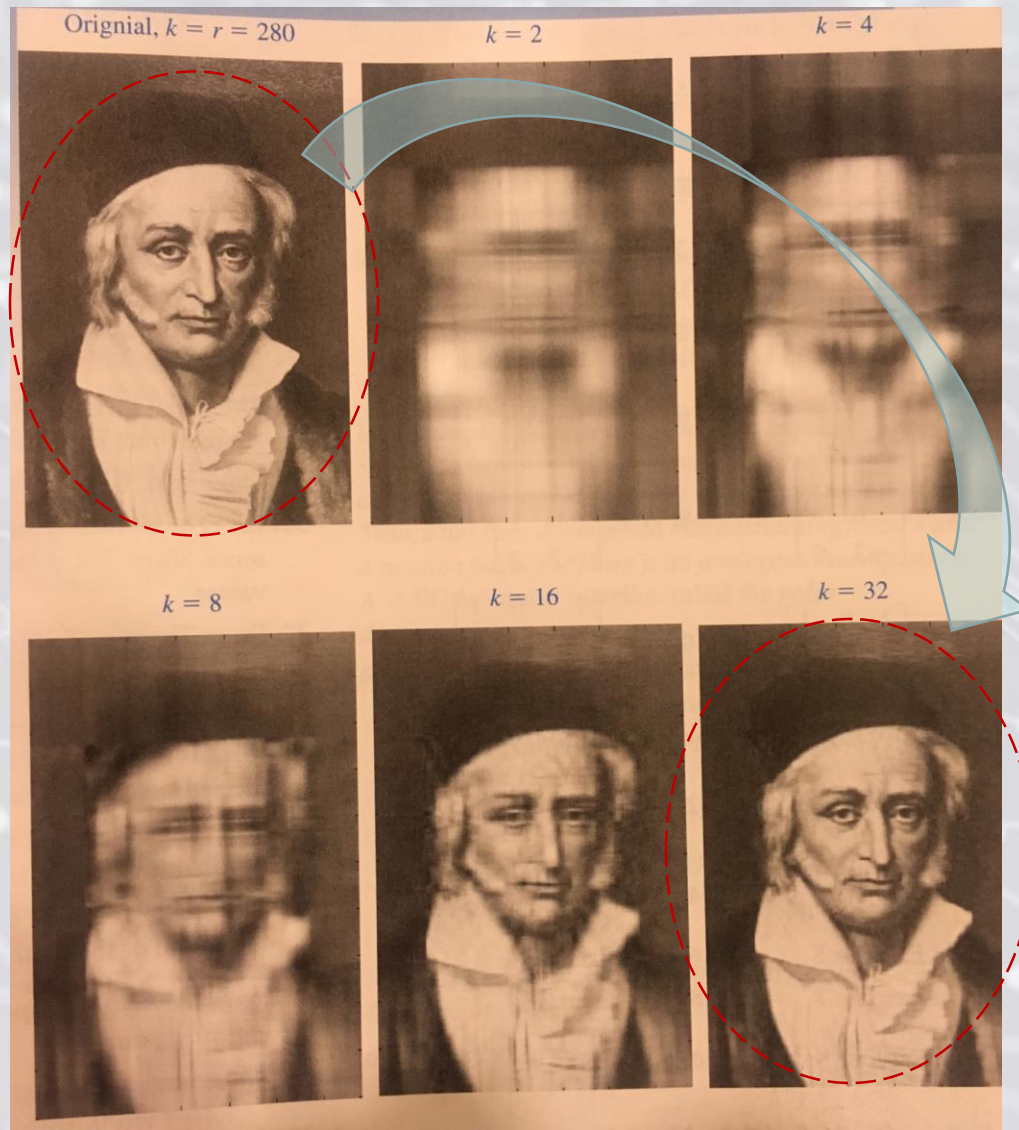$$\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \ldots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$$

- For the original *340 x 280* image shown, we have r = 280 (why?).

- Define: $\mathbf{A}_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \ldots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T, \ \ k \le r$

    as the *k-rank approximation to **A***.

(*) If for example, we use a **k = 20** *rank approximation* for **A** (i.e. we use the largest 20 singular values), the storage/ computational overhead is <u>reduced from 95,200 numbers to 12,420</u>!
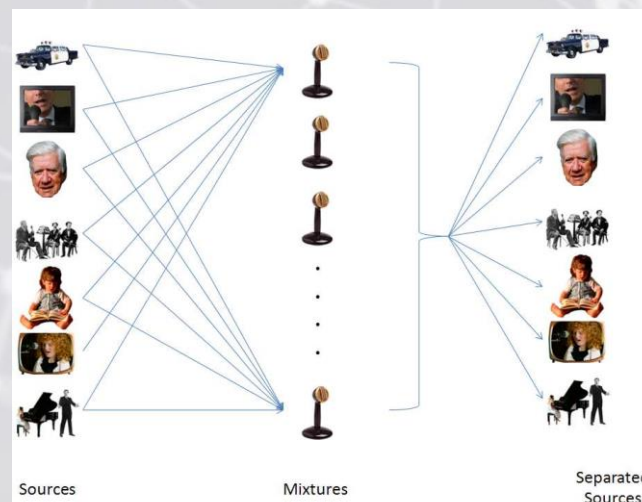
# SVD: Outer Product Form for Image Compression



$$\mathbf{A}_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + ... + \sigma_k \mathbf{u}_k \mathbf{v}_k^T, \quad k = 32$$
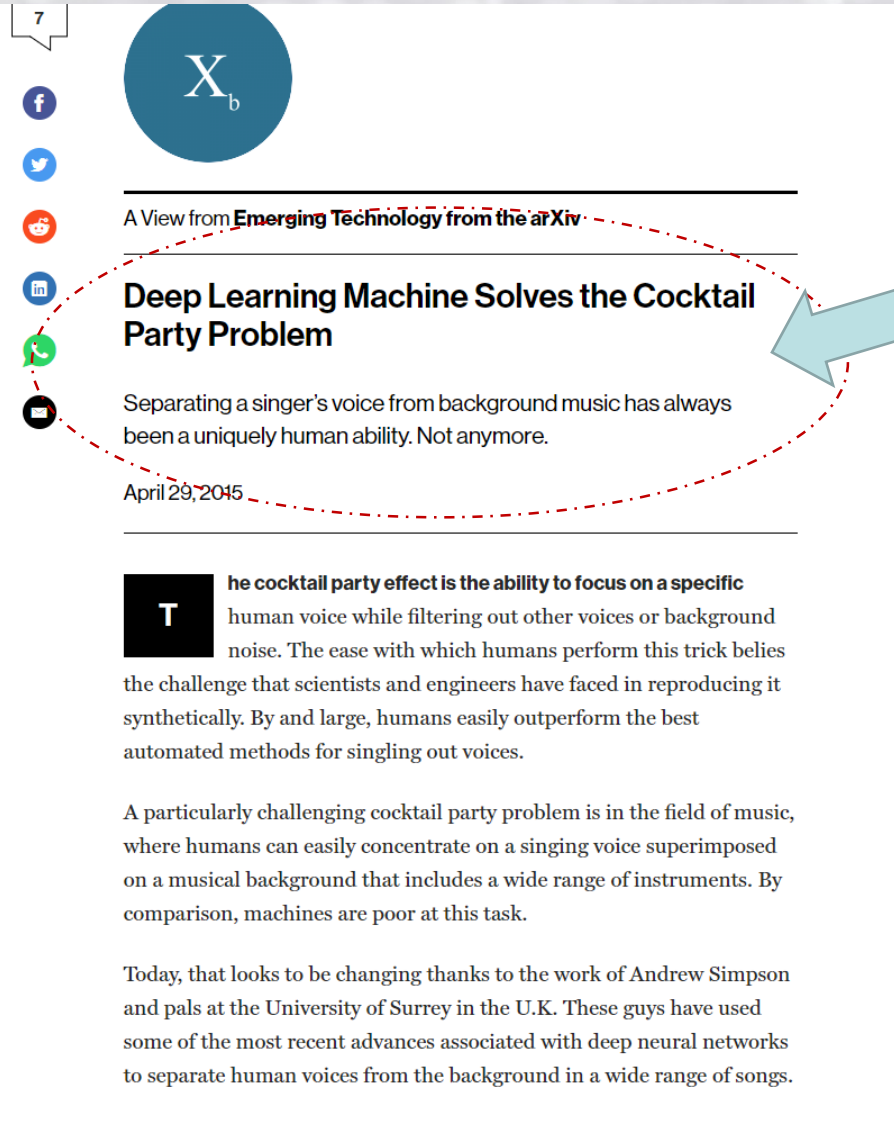
(*) Here, using the SVD-based, low-rank approximation to $\mathbf{A}$, the fidelity of the image is very strong – even after discarding roughly 85% of the image data!

# SVD: Audio Denoising & the "cocktail party problem"

- The **cocktail party effect** is the phenomenon of the brain's ability to focus one's auditory attention (an effect of selective attention in the brain) on a particular stimulus while filtering out a range of other stimuli.

- Listeners have the ability to both <u>segregate different stimuli into different streams</u>, and subsequently decide which streams are most pertinent to them. Thus, it has been proposed that one's sensory memory subconsciously siphons through all stimuli, and when an important word or phrase with high meaning appears, it stands out to the listener.

- This effect is what allows most people to "tune into" a single voice and "tune out" all others.



Sources          Mixtures          Separated Sources

# SVD: Audio Denoising & the "cocktail party problem"

A View from **Emerging Technology from the arXiv**

## Deep Learning Machine Solves the Cocktail Party Problem

Separating a singer's voice from background music has always been a uniquely human ability. Not anymore.

April 29, 2015

**T**he cocktail party effect is the ability to focus on a specific human voice while filtering out other voices or background noise. The ease with which humans perform this trick belies the challenge that scientists and engineers have faced in reproducing it synthetically. By and large, humans easily outperform the best automated methods for singling out voices.

A particularly challenging cocktail party problem is in the field of music, where humans can easily concentrate on a singing voice superimposed on a musical background that includes a wide range of instruments. By comparison, machines are poor at this task.

Today, that looks to be changing thanks to the work of Andrew Simpson and pals at the University of Surrey in the U.K. These guys have used some of the most recent advances associated with deep neural networks to separate human voices from the background in a wide range of songs.

(*) Recent research with DNNs claims to have solved the problem.

(*) Another approach: use SVD to separate signal from noise!

# PCA & SVD

Q: How do PCA and SVD relate?

- In fact, we can use SVD to perform PCA.

- Given a data matrix **X**, perform the SVD decomposition:

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

# PCA & SVD

Q: How do PCA and SVD relate?

- In fact, we can use SVD to perform PCA.

- Given a data matrix **X**, perform the SVD decomposition:

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

- Next, we compute the covariance **C** of **X** (where the data is assumed centered):

$$C = XX^T = U\Sigma V^T \left(U\Sigma V^T\right)^T = U\Sigma V^T V\Sigma^T U^T$$

# PCA & SVD

- Given a data matrix **X**, perform the SVD decomposition:

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

- Next, we compute the covariance **C** of **X** (where the data is assumed centered):

$$C = XX^T = U\Sigma V^T \left(U\Sigma V^T\right)^T = U\Sigma V^T V \Sigma^T U^T$$

$$= U\Sigma \left(V^T V\right) \Sigma U^T$$

Why?

# PCA & SVD

- Given a data matrix **X**, perform the SVD decomposition:

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

- Next, we compute the covariance **C** of **X** (where the data is assumed centered):

$$C = XX^T = U\Sigma V^T \left(U\Sigma V^T\right)^T = U\Sigma V^T V\Sigma^T U^T$$

$$= U\Sigma \left(V^T V\right)\Sigma U^T = U\Sigma \left(V^T V\right)\Sigma U^T$$

Why?

# PCA & SVD

- Given a data matrix **X**, perform the SVD decomposition:

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

- Next, we compute the covariance **C** of **X** (where the data is assumed centered):

$$C = XX^T = U\Sigma V^T \left(U\Sigma V^T\right)^T = U\Sigma V^T V\Sigma^T U^T$$
$$= U\Sigma\left(V^T V\right)\Sigma U^T = U\Sigma\left(V^T V\right)\Sigma U^T = U\Sigma^2 U^T$$

Why?

# PCA & SVD

- Given a data matrix **X**, perform the SVD decomposition:

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

- Next, we compute the covariance **C** of **X** (where the data is assumed centered):

$$C = XX^T = U\Sigma V^T \left(U\Sigma V^T\right)^T = U\Sigma V^T V\Sigma^T U^T$$

$$= U\Sigma \left(V^T V\right)\Sigma U^T = U\Sigma \left(V^T V\right)\Sigma U^T = U\Sigma^2 U^T$$

(*) Recall that with the PCA *eigendecomposition* algorithm we generate the of C.

(*) **This is precisely what we've done here, where: D=$\Sigma^2$, which confirms that the singular values of X are indeed equivalent to the square root of the eigenvalues of XX$^T$.**

## PCA

- Here is the **PCA algorithm**:

(1) Write $N$ data points $x_i=(x_{1i},x_{2i},\ldots,x_{Mi})$ as row vectors.

(2) Put these vectors into the data matrix **X** (of size $N$ x $M$).

(3) Center the data by subtracting off the mean of each column, place into matrix **B**.

(4) Computer the covariance matrix: $\mathbf{C} = \dfrac{1}{N}\mathbf{B}^T\mathbf{B}$

(5) Computer the *eigenvalues* and *eigenvectors* of **C**, so: $\mathbf{C} = \mathbf{VDV}^T$
 where **D** is the diagonal matrix of eigenvalues; V is the matrix of corresponding eigenvectors.

(6) <u>Sort of the columns of D</u> into order of decreasing eigenvalues, and apply the same order to the columns of V.

(7) Reject those with eigenvalues less than some given threshold, leaving $L$ dimensions in the data.

# Fin