

lecture follows notes on webpage.

PCA Principal Components Analysis

data dependent dimensionality reduction

we have a dataset x_1, x_2, \dots, x_n in \mathbb{R}^d

Credit for figures:
Roughgarden & Valiant
John Benedetto
Novembre et al
Alex Williams
Sandipan Dey
Victor Lawrence

Useful for

- visualization
- interpretation
- compression

find intrinsic dimensionality of the data.

Example: $n=4$ data pts
 $d=4$ dimensions

	kale	taco bell	sashimi	pop tarts
x_A	10	1	2	7
x_B	7	2	1	10
x_C	2	9	7	3
x_D	3	6	10	2

\bar{x} 5.5 4.5 5 5.5

x_i = row of table, i.e. one person's preferences.

Claim:

each row can be expressed \approx
as $\vec{x}_i = \bar{\vec{x}} + a_{ii} \vec{v}_1 + a_{12} \vec{v}_2$

reduce person's preferences from a 4D system
(kale, taco bell, sashimi, pop tarts) to a 2D system (v_1, v_2)

$$\bar{\vec{x}} = (5.5, 4.5, 5, 5.5)$$

$$\vec{v}_1 = (3, -3, -3, 3)$$

$$\vec{v}_2 = (1, -1, 1, -1)$$

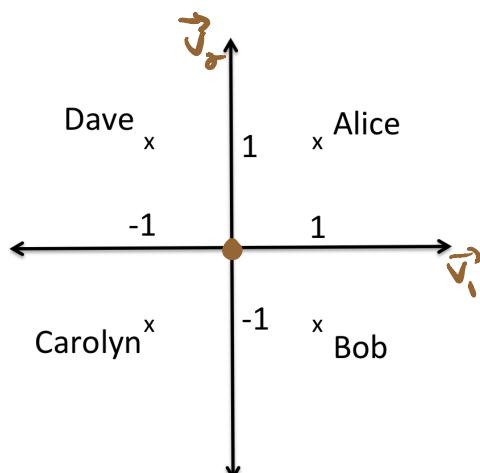
3's on kale and pop tarts (the two are correlated) and -3's on taco bell / sashimi (two are correlated). Captures some structure in data / correlation

notice.
 $(v_1 \cdot v_2 = 0)$

orthogonal / uncorrelated

	kale	taco bell	sashimi	pop tarts
\vec{x}_A	10	1	2	7
\vec{x}_B	7	2	1	10
\vec{x}_C	2	9	7	3
\vec{x}_D	3	6	10	2

$$\bar{\vec{x}} \quad 5.5 \quad 4.5 \quad 5 \quad 5.5$$



$$\begin{aligned}\vec{x}_A &\approx \bar{\vec{x}} + 1 \cdot \vec{v}_1 + 1 \cdot \vec{v}_2 \\ \vec{x}_B &\approx \bar{\vec{x}} + 1 \cdot \vec{v}_1 + (-1) \cdot \vec{v}_2 \\ \vec{x}_C &\approx \bar{\vec{x}} + (-1) \cdot \vec{v}_1 + (-1) \cdot \vec{v}_2 \\ \vec{x}_D &\approx \bar{\vec{x}} + (-1) \cdot \vec{v}_1 + 1 \cdot \vec{v}_2\end{aligned}$$

average columnwise

$$(a_{ii}, a_{12})$$

Alice	1	1
Bob	1	-1
Carol	-1	-1
Dave	-1	1

origin is $\bar{\vec{x}}$

Why bother?

① visualization

② helps interpret data

each pt has " \vec{v}_1 " coordinate
" \vec{v}_2 " coordinate

develops the best coordinate system!

Goal: express n d -dimensional vectors

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ as linear combinations
of k vectors $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$
with $k \ll n$ ideally

$$\text{so } \vec{x}_i \approx \sum_{j=1}^k a_{ij} \vec{v}_j$$

(assuming $\sum_{i=1}^n \vec{x}_i = 0$)

Preprocessing:

① Subtract mean

② Scale coordinates so that
variance is 1 in each coords

[to avoid sensitivity to units]

SVD chooses coordinates to maximize variance. If our units create different magnitudes in our data this can naturally generate more/less variance, but that does not mean the data is more/less correlated in that direction.

Framing PCA as an optimization problem.

K=1

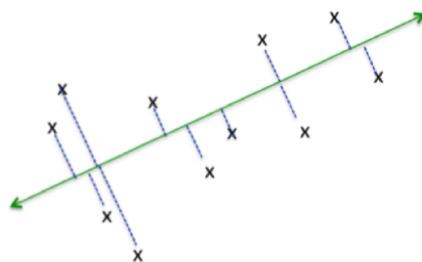
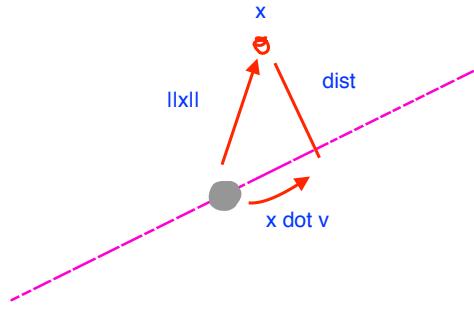
choose $\vec{v} \in \mathbb{R}^d$ with $\|\vec{v}\|=1$

unit vector!

to minimize

least squared error between the line specified by vector v (thru origin?) and data

$$\frac{1}{n} \sum_{i=1}^n (\text{dist}(x_i \leftrightarrow \text{line defined by } \vec{v}))^2$$

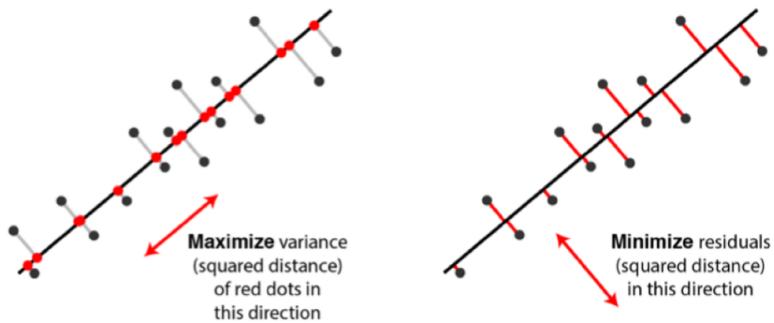


$$(x \cdot v)^2 + \text{dist}^2 = \|x\|^2 = \text{constant!}$$

Thus minimizing dist is equivalent to maximizing the average $(x \cdot v)$ = squared projections

MAXIMIZING VARIANCE OF PROJECTED POINTS

Choose v with $\|v\|=1$ that maximizes $\frac{1}{n} \sum_i (x_i \cdot v)$



Two equivalent views of principal component analysis.

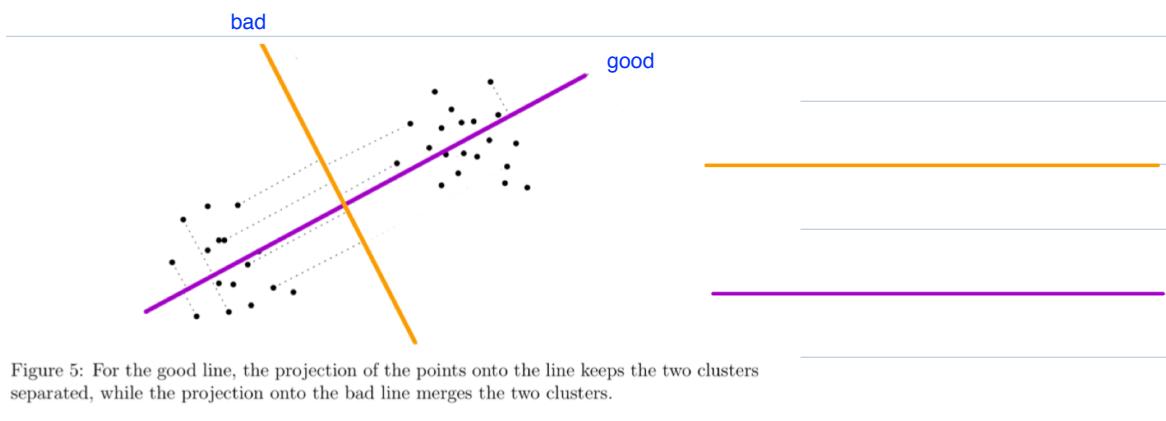


Figure 5: For the good line, the projection of the points onto the line keeps the two clusters separated, while the projection onto the bad line merges the two clusters.

Comparison to linear regression

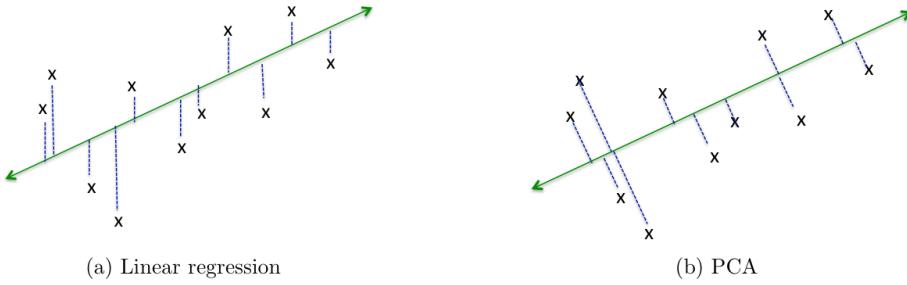


Figure 2: Linear regression minimizes the sum of squared vertical distances, while PCA minimizes the sum of squared perpendicular distances.

Goal: express n d -dimensional vectors $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ as linear combinations of k vectors $\vec{v}_1, \dots, \vec{v}_k \in \mathbb{R}^d$

so $\vec{x}_i \approx \sum_{j=1}^k a_{ij} \vec{v}_j$

Larger k

find k dim subspace
s.t.

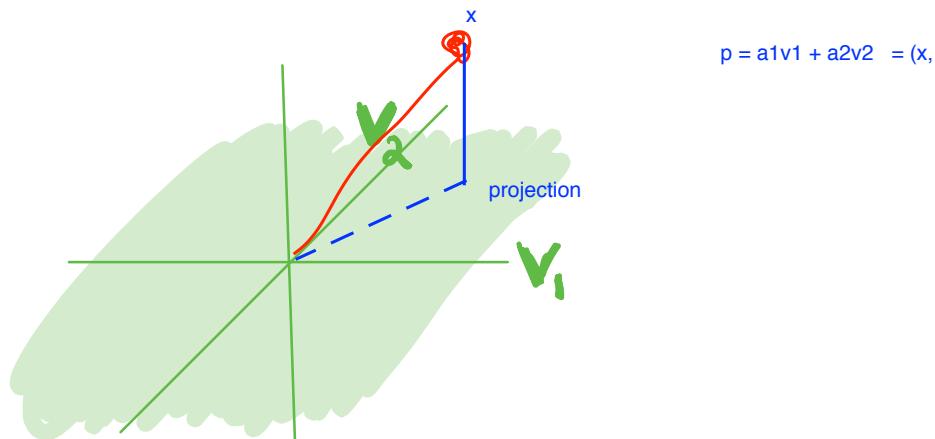
Squared distances = variance of
minimized projection
maximized.

$$\max_{\text{k dim subspaces } S} \frac{1}{n} \sum_{i=1}^n (\text{length of } x_i \text{'s projection})^2$$

define subspace using k orthonormal vectors v_1, \dots, v_k in R^d
 $\|v_i\| = 1$ and $(v_i, v_j) = 0$ for $i \neq j$

Subspace = $\text{span}\{v_1, \dots, v_k\}$

$$\max \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k (x_i \cdot v_j)^2$$



$$p = a_1 v_1 + a_2 v_2 = (x,$$

Summary:

Problem: compute orthonormal vectors v_1, \dots, v_k
to maximize

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k (x_i, v_j)^2$$

|| projection onto $\text{span}\{v_1, \dots, v_k\}$ ||^2

arg sum of squared projection lengths

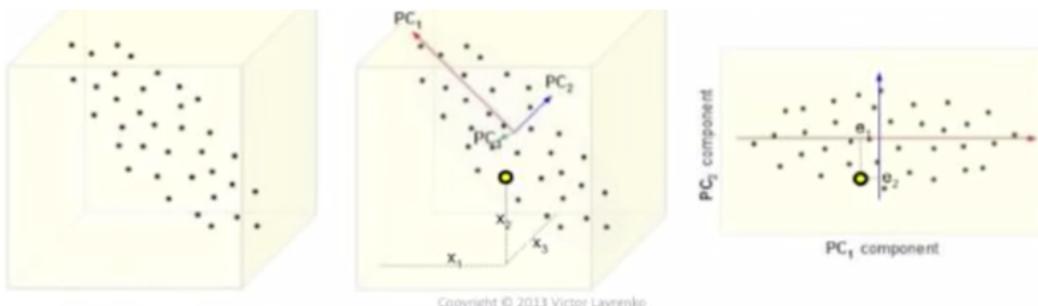
Application: Visualization

① Perform PCA $\rightarrow \vec{v}_1, \dots, \vec{v}_k$ top k
"principal components"

② $\forall \vec{x}_i$ define
 " \vec{v}_1 -coord" (x_i, v_1)
 " \vec{v}_2 -coord" (x_i, v_2)
 " \vdots "
 " \vec{v}_k -coord" (x_i, v_k)

③ Plot points

$$x_i \rightarrow ((x_i, v_1), (x_i, v_2), \dots, (x_i, v_k))$$



- * look for clusters
- * look for pts particularly large along \vec{v}_i

Cool example

Genetic data of Europeans

[Novembre et al]

3192 people

200,000 SNPs

DNA positions that
tend to exhibit gene mutation.

Using only genomic data

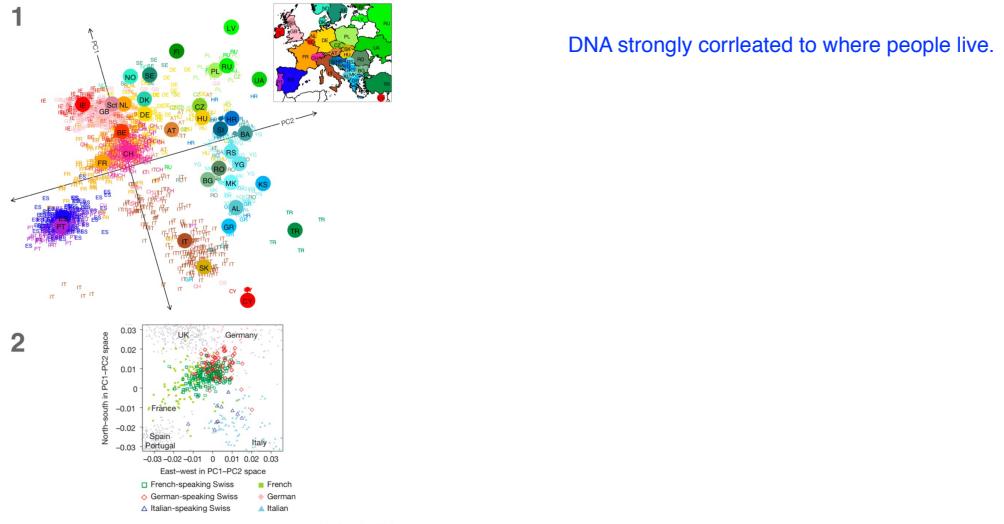


Figure 1: The genetic map of Europe using PCA, with the geographic map of Europe for reference. Figure 2: The same map, but zoomed in on Switzerland. Swiss individuals tend to cluster with countries that speak the same language. (Courtesy: John Novembre, UCLA)

How PCA works

K=1

finding the first principal component

$$\underset{\mathbf{v}}{\operatorname{arg \max}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i, \mathbf{v})^2$$

$\|\mathbf{v}\|=1$

$$\mathbf{x}_i \in \mathbb{R}^d$$
$$\mathbf{v} \in \mathbb{R}^d$$

$$\mathbf{X} = \begin{pmatrix} -\mathbf{x}_1- \\ \vdots \\ -\mathbf{x}_n- \end{pmatrix}$$

$n \times d$

means are zero as we subtracted out expected value.

finding the first principal component.

Problem: find \vec{v} with $\|v\|=1$
to maximize $v^T A v$
where $A = X^T X$

Suppose magically A diagonal

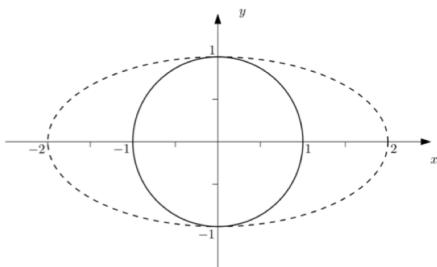


Figure 1: The point (x, y) on the unit circle is mapped to $(2x, y)$.

$$A = \begin{pmatrix} 3 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0.5 \end{pmatrix}$$

Let $v = (v_1, v_2, v_3)$

$$v^T A v = 3v_1^2 + v_2^2 + 0.5v_3^2 \text{ such that } v_1^2 + v_2^2 + v_3^2 = 1$$

Obviously to maximize the first quantity, just choose the coordinate that has the largest coefficient to one! i.e. set $v_1 = 1$ and $v_2, v_3 = 0$

$$\begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$$

Every symmetric matrix is "diagonals in disguise"

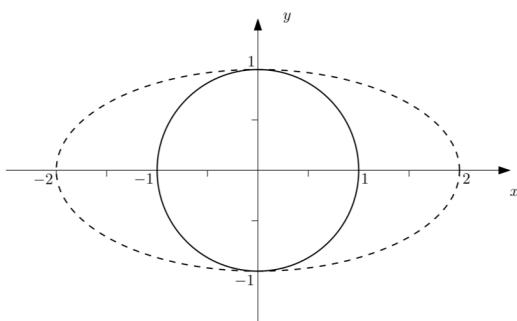


Figure 1: The point (x, y) on the unit circle is mapped to $(2x, y)$.

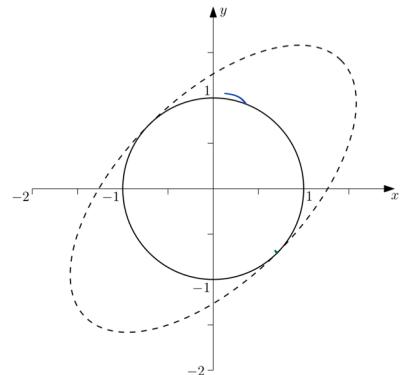


Figure 2: The same scaling as Figure 1, but now rotated 45 degrees.

$$\begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\begin{pmatrix} \frac{3}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{3}{2} \end{pmatrix} = \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{rotate back } 45^\circ} \cdot \underbrace{\begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}}_{\text{stretch}} \cdot \underbrace{\begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}}_{\text{rotate clockwise } 45^\circ}.$$

For any symmetric matrix A we can write A as

$$A = Q D Q^T$$

Text

↑ ↑ ↑
 orthogonal diagonal $\begin{pmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_n & \\ & & & 0 \end{pmatrix}$

(orthonormal columns; eigenvectors corresponding to eigenvalues in D)

ORTHOGONAL MATRICES preserve length; i.e. if Q is symmetric then for any x : $\|x\| = \|Qx\|$

with $\lambda_1 > \lambda_2 > \dots > 0$ as matrix is positive semi-definite ($X^T X$ is always positive semi-definite)

$$A = Q D Q^T$$

$$Q^T = \begin{pmatrix} z_1^T & \dots & z_n^T \end{pmatrix}$$

also has direction of max stretch

"rotated axis" that gets stretched the most.

$$Q = \begin{pmatrix} z_1 & \dots & z_n \end{pmatrix}$$

eigenvector corresponding to the largest eigenvalue (since $X^T X$ is positive semi-definite, we don't have to worry about negative eigenvectors)

Solution to maximization problem

$$\begin{aligned} \max_{\mathbf{v}} \quad & \mathbf{v}^T A \mathbf{v} \\ \text{s.t.} \quad & \|\mathbf{v}\| = 1 \end{aligned}$$

Every matrix of form $A = X^T X$ can
be written as $Q^T D Q$ where
all $\lambda_i \geq 0$

$$Q^T = \begin{pmatrix} z_1 & \dots & z_n \end{pmatrix}$$

$$Q = \begin{pmatrix} z_1 & \dots & z_n \end{pmatrix}$$

$$D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$$

Conclusion: for $k=1$, soln is first column of Q

for general k , soln
is first k cols of Q

\equiv top k principal components

Computing principal components:

- ① SVD. [cubic]
- ② power iteration

Algorithm 1
POWER ITERATION

Given matrix $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$:

- Select random unit vector \mathbf{u}_0
- For $i = 1, 2, \dots$, set $\mathbf{u}_i = \mathbf{A}^i \mathbf{u}_0$. If $\mathbf{u}_i / \|\mathbf{u}_i\| \approx \mathbf{u}_{i-1} / \|\mathbf{u}_{i-1}\|$, then return $\mathbf{u}_i / \|\mathbf{u}_i\|$.

To see why this works:

Suppose v^1, v^2, \dots, v^n are eigenvectors of \mathbf{A}

Then they form an orthonormal basis.

$$\Rightarrow \text{If vector } u: u = \sum_{i=1}^n c_i v_i$$

$$Au = A \sum c_i v_i = \sum c_i Av_i$$

$$= \sum c_i \lambda_i v_i$$

$$A^k u = A^{k-1} Au = \underbrace{\sum c_i \lambda_i A^{k-1} v_i}_{c_i \lambda_i A^{k-2} \lambda_i v_i \dots} \\ = \sum c_i \lambda_i^k v_i$$

$$= c_1 \lambda_1^k \left[v_1 + \left(\frac{\lambda_2}{\lambda_1} \right)^k v_2 + \left(\frac{\lambda_3}{\lambda_1} \right)^k v_3 + \dots \right] \\ \approx c_1 \lambda_1^k v_1$$

\Rightarrow when we normalize,

get v_i .

this depended on $\lambda_2 < \lambda_1$
& then can show that $O\left(\frac{\log n}{\log(\lambda_1/\lambda_2)}\right)$
iterations suffice to get
very close.

1. Find the top component, \mathbf{v}_1 , using power iteration.

2. Project the data matrix orthogonally to \mathbf{v}_1 :

$$\begin{bmatrix} \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \vdots & & \vdots \\ \cdots & \mathbf{x}_m & \cdots \end{bmatrix} \mapsto \begin{bmatrix} \cdots & (\mathbf{x}_1 - \langle \mathbf{x}_1, \mathbf{v}_1 \rangle \mathbf{v}_1) & \cdots \\ \cdots & (\mathbf{x}_2 - \langle \mathbf{x}_2, \mathbf{v}_1 \rangle \mathbf{v}_1) & \cdots \\ \vdots & & \vdots \\ \cdots & (\mathbf{x}_m - \langle \mathbf{x}_m, \mathbf{v}_1 \rangle \mathbf{v}_1) & \cdots \end{bmatrix}.$$

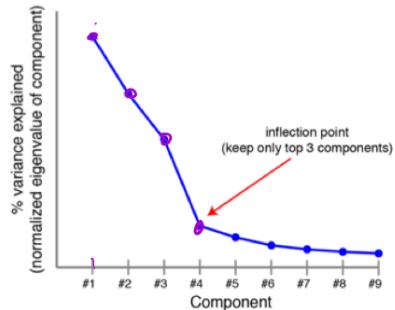
This corresponds to subtracting out the variance of the data that is already explained by the first principal component \mathbf{v}_1 .

3. Recurse by finding the top $k-1$ principal components of the new data matrix.

"Greedy alg" correctness: k -dim subspace that maximizes norms of projections contains $k-1$ -dimensional space that maximizes norms of projections.

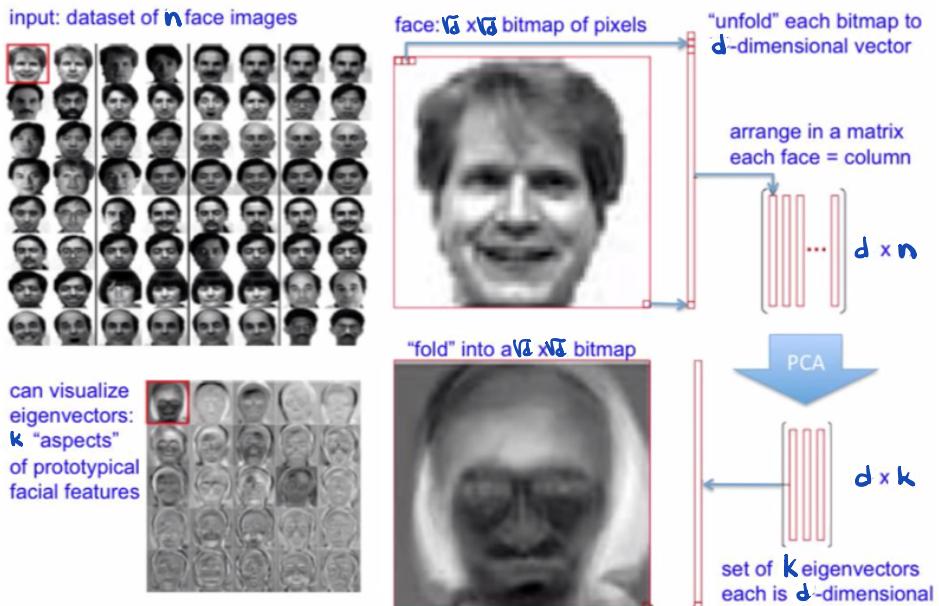
How to choose k ?

- For data visualization: a few
- compression: Look at eigenvalues.
As soon as small enough; happy.



Scree plot. Principal components are ranked by the amount of variance they capture in the original dataset, a scree plot can provide some sense of how many components are needed.

PCA example: Eigen Faces



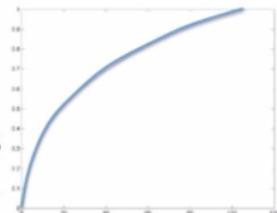
Eigenfaces
Slides by
Victor
Lawrenko

Eigen Faces: Projection

$$= \text{mean} + 0.9 * \text{eigenvector}_1 - 0.2 * \text{eigenvector}_2 + 0.4 * \text{eigenvector}_3 + \dots$$



- Project new face to space of eigen-faces
- Represent vector as a linear combination of principal components
- How many do we need?



(Eigen) Face Recognition

- Face similarity
 - in the reduced space
 - insensitive to lighting expression, orientation
- Projecting new “faces”
 - everything is a face



new face

projected to eigenfaces

Copyright © 2013 Victor Lavrenko

