

PCA Principal Components Analysis

data dependent dimensionality reduction

data set $x_1, \dots, x_n \in \mathbb{R}^d$

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

Useful for

- visualization
- interpretation
- compression
- finding intrinsic dimensionality

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 ←

Claim:

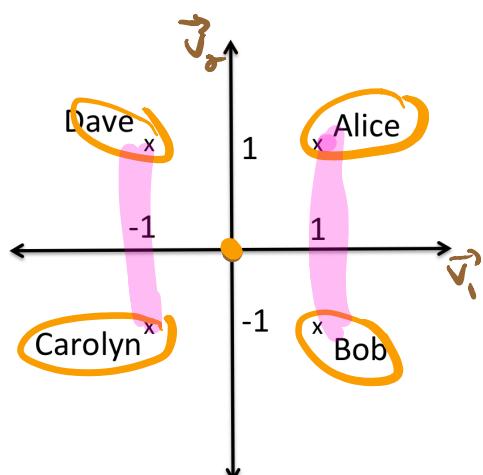
each row can be expressed \approx
as $x_i = \bar{x} + a_{i1}\vec{v}_1 + a_{i2}\vec{v}_2$

$$\begin{aligned}\bar{x} &= (5.5, 4.5, 5, 5.5) \\ + 1 \cdot \vec{v}_1 &= (3, -3, -3, 3) \quad \text{notice.} \\ + 1 \cdot \vec{v}_2 &= (1, -1, 1, -1) \quad (v_1 \cdot v_2 = 0)\end{aligned}$$

$(9.5, 0.5, 3, 7.5)$

	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

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



Claim:

$$\begin{aligned}x_A &\approx \bar{x} + 1 \cdot \vec{v}_1 + 1 \cdot \vec{v}_2 \\ x_B &\approx \bar{x} + 1 \cdot \vec{v}_1 + (-1) \cdot \vec{v}_2 \\ x_C &\approx \bar{x} + (-1) \cdot \vec{v}_1 + (-1) \cdot \vec{v}_2 \\ x_D &\approx \bar{x} + (-1) \cdot \vec{v}_1 + 1 \cdot \vec{v}_2\end{aligned}$$

$$\begin{pmatrix} (a_{ii}, a_{i2}) \\ \hline \text{Alice} & 1 & 1 \\ \text{Bob} & 1 & -1 \\ \text{Carol} & -1 & -1 \\ \text{Dave} & -1 & 1 \end{pmatrix}$$

Why bother?

① visualization

② helps interpret data

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

Goal: express n d -dimensional vectors

$\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ as linear combinations
of K vectors $\vec{v}_1, \dots, \vec{v}_K \in \mathbb{R}^d$

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

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

Preprocessing:

① Subtract mean

Henceforth we will assume

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

② Scale coordinates so that
variance is 1 in each coords

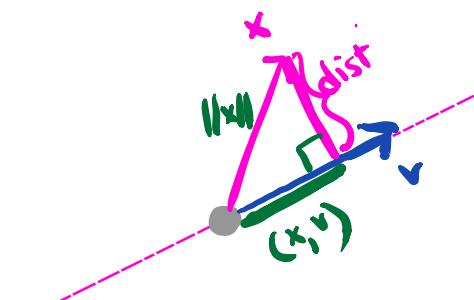
[to avoid sensitivity to units]

K=1

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

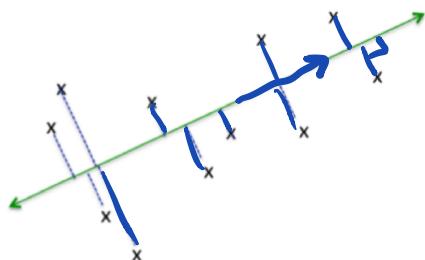
to minimize

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



$$(x_i \cdot v)^2 + \text{dist}^2 \geq \|x_i\|^2$$

$\uparrow \max \quad \uparrow \min$



objective: choose v s.t. $\|v\|=1$

$$\text{to max } \frac{1}{n} \sum_{i=1}^n (x_i \cdot v)^2$$

max avg squared projections

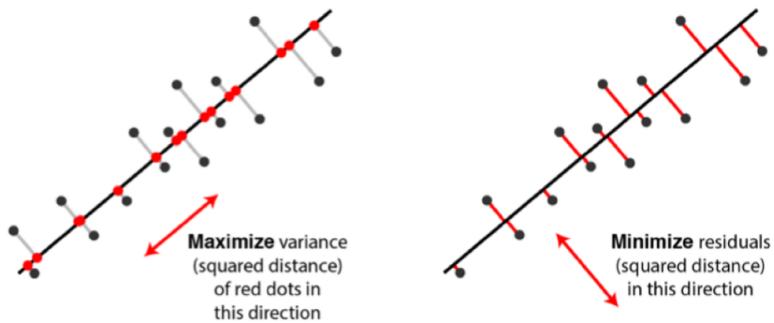
max variance of projections.

$$(x_1, v), (x_2, v) \dots (x_n, v)$$

$$\boxed{\frac{1}{n} \sum_{i=1}^n (x_i \cdot v)^2}$$

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n (x_i \cdot v) \\ &= \left(\frac{1}{n} \sum_{i=1}^n x_i \right) \cdot v = 0 \end{aligned}$$

$$\text{Var}(x) = E((X - E(X))^2)$$



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

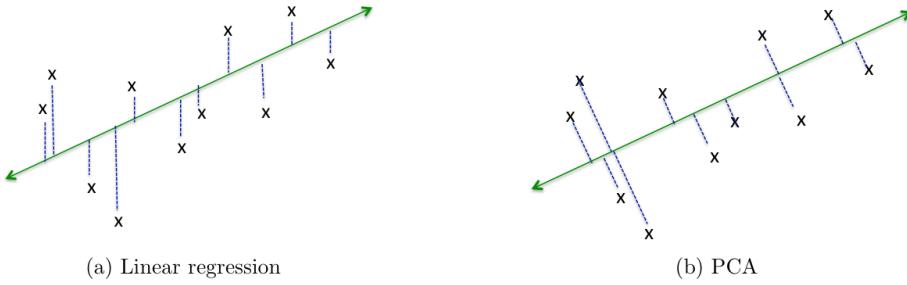


Figure 2: Linear regression minimizes the sum of squared vertical distances, while PCA minimizes the sum of squared perpendicular distances.

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

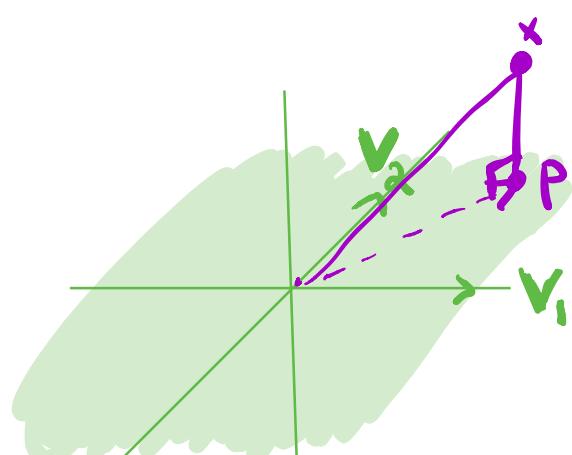
find k dim subspace
s.t.

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 $v_1, \dots, v_k \in \mathbb{R}^d$
 $\|v_i\|=1 \quad (v_i, v_j) = 0 \text{ if } i \neq j$

subspace: span of v_1, \dots, v_k



Find v_1, \dots, v_k
to max $\frac{1}{n} \sum_{i=1}^n \left[\sum_{j=1}^k (x_i, v_j)^2 \right]$

$$p = a_1 v_1 + a_2 v_2$$

$$(x, v_1)$$

$$(x - p, v_1) = 0$$

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

$$p = \underbrace{(x, v_1)}_{= a_1} v_1 + \underbrace{(x, v_2)}_{= a_2} v_2$$

Summary:

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

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

arg sum of squared projection lengths

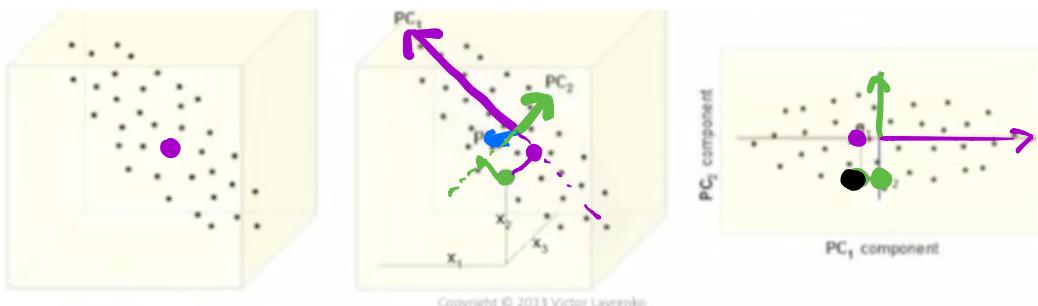
Application: Visualization

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)
⋮
" \vec{v}_k -coord" (x_i, v_k)

③ Plot points

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



Copyright © 2013 Victor Lavrenko

- * 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

$$n = 3192$$
$$d = 200,000$$

DNA positions that
tend to exhibit gene mutation.

Using only genomic data

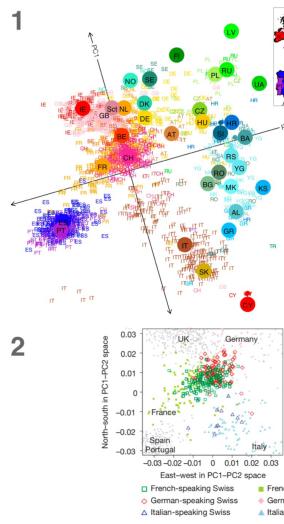


Figure 1: The genetic map of Europe using PCA, with the geographic map of Europe for reference. Figure 2: The same map, but zoomed in on Switzerland. Swiss individuals tend to cluster with countries that speak the same language. (Courtesy: John Novembre, UCLA)



courtesy: John Novembre, UCLA

How PCA works

$K=1$

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

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

$$\mathbf{X} = \begin{pmatrix} -\mathbf{x}_1- \\ -\mathbf{x}_2- \\ \vdots \\ -\mathbf{x}_n- \end{pmatrix} \quad \mathbf{Xv} = \begin{pmatrix} \mathbf{x}_1 \cdot \mathbf{v} \\ \vdots \\ \mathbf{x}_n \cdot \mathbf{v} \end{pmatrix}$$

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

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

$$\text{obj: } \frac{1}{n} (\mathbf{Xv})^T \mathbf{Xv}$$

$$= \frac{1}{n} \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v}$$

$$\mathbf{A} = \frac{1}{n} \mathbf{X}^T \mathbf{X}$$

empirical covariance matrix

$$(\mathbf{X}^T \mathbf{X})^T = \mathbf{X}^T \mathbf{X}$$

$$\frac{1}{n} \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} -x_1 & + \\ -x_2 & + \\ \vdots & \vdots \\ -x_n & + \end{pmatrix}$$

$$\mathbf{A}_{Kk} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{ik} \mathbf{x}_{ik}$$

\mathbf{X} : n docs.
d words

$x_{ij} = \# \text{ freq word } j \text{ appears in doc } i$

\mathbf{x}_k random sample of k^{th} feature.

$$x_{1k}, x_{2k}, \dots, x_{nk}$$

\mathbf{x}_e random sample of e^{th} feature.

rows Y & Z.

$$\mathbf{A}_{ke} = \operatorname{Cov}_j (\mathbf{x}_k, \mathbf{x}_e)$$

$$\frac{1}{n} \mathbf{X}^T \mathbf{X} = \begin{pmatrix} \operatorname{Var}(\mathbf{x}_k, \mathbf{x}_k) & & \\ & \operatorname{Var}(\mathbf{x}_k, \mathbf{x}_e) & \\ & & \operatorname{Cov}(\mathbf{x}_k, \mathbf{x}_e) \end{pmatrix}$$

$$\operatorname{Cov}(Y, Z) = E((Y - E(Y))(Z - E(Z)))$$

if means are 0.

$$\operatorname{Cov} \Rightarrow E(YZ)$$

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

Suppose magically A diagonal

$$\text{optimal } v = (1, 0, 0)$$

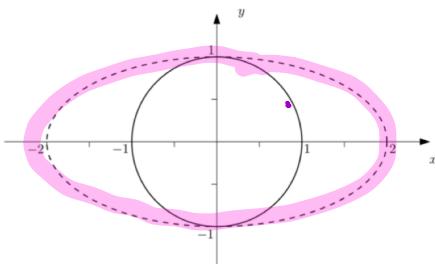


Figure 1: The point (x, y) on the unit circle is mapped to $(2x, y)$.

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

$$\max v^T A v \quad \|v\|^2 = 1$$

$$v^T A v = (v_1, v_2, v_3) \begin{pmatrix} 3 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0.5 \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix}$$

$$= 3v_1^2 + v_2^2 + 0.5v_3^2$$

$v_1^2 + v_2^2 + v_3^2 = 1$

$$v_1 = 1 \quad v_2, v_3 = 0$$

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

$$\boxed{\vec{v} = e_1}$$

Every symmetric matrix is "diagonals in disguise"

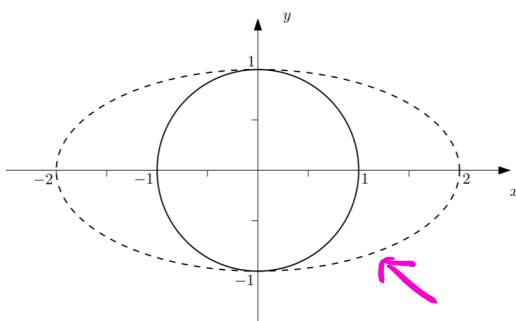


Figure 1: The point (x, y) on the unit circle is mapped to $(2x, y)$.

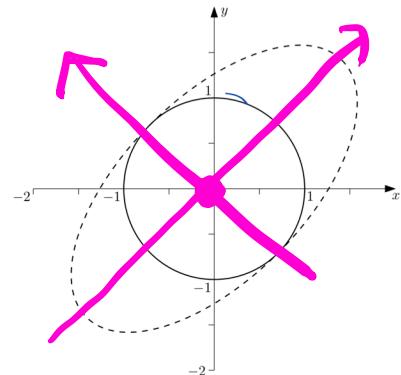


Figure 2: The same scaling as Figure 1, but now rotated 45 degrees.

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

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

$$A = Q D Q^T$$

orthogonal
matrix whose
columns are
orthonormal

$$\begin{pmatrix} \star & \star & 0 \\ 0 & \star & \star \\ \star & 0 & \star \end{pmatrix}$$

$$Q^T Q = Q Q^T = I$$

Orthogonal matrices preserve

$$\begin{aligned} \text{length } \|Q w\|^2 &= (Q w)^T Q w \\ &= w^T Q^T Q w = w^T w = \|w\|^2 \end{aligned}$$

JY

" "

$$A = Q D Q^T$$

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

also has direction of max stretch

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

"rotated axis" that gets stretched the most.

Solution to maximization problem

$$\max_{\mathbf{v} \text{ s.t. } \|\mathbf{v}\|=1} \mathbf{v}^T A \mathbf{v}$$

$$\mathbf{v}^T A \mathbf{v} = \underbrace{\mathbf{v}^T Q D Q^T \mathbf{v}}_{y^T D y} = y^T D y$$

$$D = \begin{pmatrix} \lambda_1 & & & 0 \\ & \ddots & & \\ 0 & & \ddots & \lambda_d \end{pmatrix}$$

$$\lambda_1 > \lambda_2 > \dots > \lambda_d$$

$$= y^T D y$$

↑
cpt y is e_1

$$\boxed{Q^T v = e_1} \quad \boxed{Q Q^T v = v = Q e_1}$$

↑

$$\begin{aligned} Q^T v &= e_1 \\ Q Q^T v &= Q e_1 \\ \boxed{H} \quad v &= Q e_1 \end{aligned}$$

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

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

\equiv top k eigenvectors of A

$$A = Q D Q^T$$

$$\boxed{A \vec{Q} e_i} = Q D Q^T \vec{Q} e_i$$

$$= Q D e_i$$

$$= Q \lambda_i e_i$$

$$= \lambda_i \boxed{\vec{Q} e_i}$$

$\vec{Q} e_i$ is an eigenvector w/ eigenvalue λ_i :

$$\boxed{B \vec{w} = \lambda \vec{w}}$$

...

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

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

Then they form an orthonormal basis.

$$\Rightarrow \text{If vector } \mathbf{u}: \quad \mathbf{u} = \sum_{i=1}^n c_i \mathbf{v}_i$$

$$\mathbf{A}\mathbf{u} = \mathbf{A} \sum c_i \mathbf{v}_i = \sum c_i \mathbf{A}\mathbf{v}_i$$

$$= \sum c_i \lambda_i \mathbf{v}_i$$

$$\begin{aligned} \mathbf{A}^k \mathbf{u} &= \mathbf{A}^{k-1} \mathbf{A} \mathbf{u} = \sum c_i \underbrace{\lambda_i \mathbf{A}^{k-1} \mathbf{v}_i}_i \\ &= \sum c_i \lambda_i^k \mathbf{v}_i \end{aligned}$$

$$\begin{aligned} &= c_1 \lambda_1^k \left[\mathbf{v}_1 + \left(\frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \left(\frac{\lambda_3}{\lambda_1} \right)^k \mathbf{v}_3 + \dots \right] \\ &\approx c_1 \lambda_1^k \mathbf{v}_1 \end{aligned}$$

\Rightarrow when we normalize,
get v ,

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 :

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

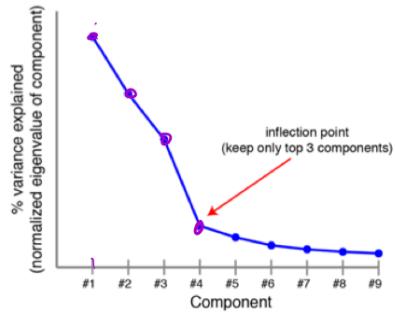
This corresponds to subtracting out the variance of the data that is already explained by the first principal component \mathbf{v}_1 .

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.

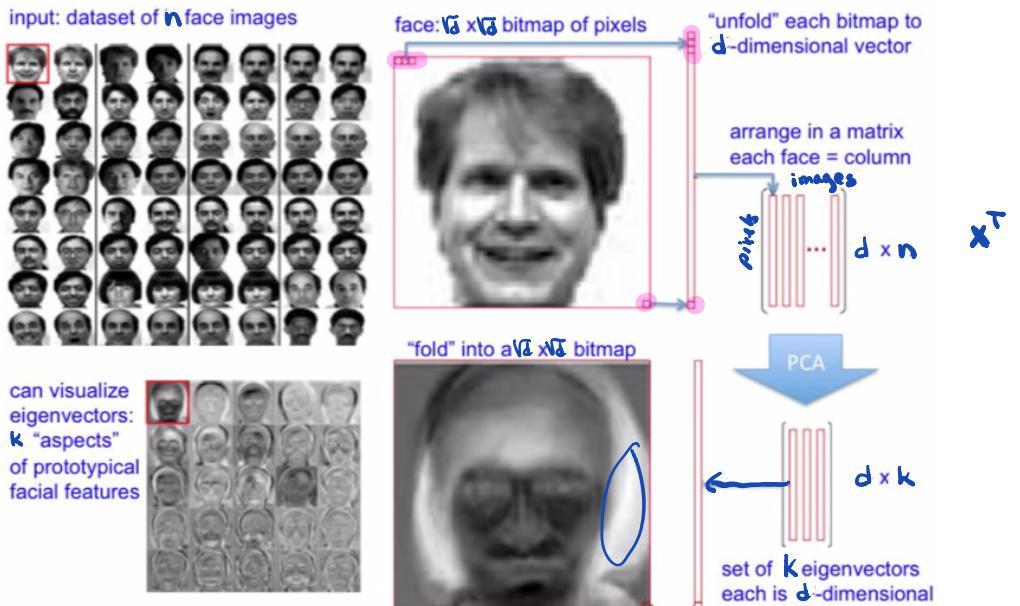
How to choose k ?

- 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



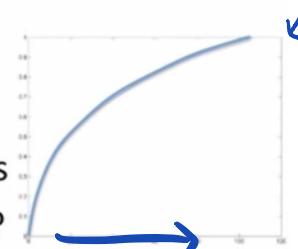
Eigenfaces Slides by Victor Lawrenko

Eigen Faces: Projection

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



- 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

