## LSH

- important technique for solving approx nearest neighbor queries

- idea: construct hash fns that are likely to map _similar_ items to same bucket

## PCA
principal components analysis

data dependent dimensionality reduction

Credit for figures:
- Roughgarden & Valiant
- John Benedetto
- Novembre et al
- Alex Williams
- Sandipan Dey
- Leskovec, Rajaraman, Ullman slides

data set $\quad x_1, \dots, x_m$
$\qquad x_i \in \mathbb{R}^n$

$n = 4$

$a_{i1} \mid a_{i2}$

| | kale | taco bell | sashimi | pop tarts |
|---|---|---|---|---|
| Alice | 10 | 1 | 2 | 7 |
| Bob | 7 | 2 | 1 | 10 |
| Carolyn | 2 | 9 | 7 | 3 |
| Dave | 3 | 6 | 10 | 2 |

Table 1: Your friends' ratings of four different foods.

$n = 4$, $x_1, x_2, x_3, x_4$

$\begin{bmatrix} 1 & 1 \\ 1 & -1 \\ -1 & 1 \\ -1 & -1 \end{bmatrix}$

$v_1 = (3, -3, -3, 3)$

$v_2 = (1, -1, 1, -1)$

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

$x_i \approx \bar{x} + a_{i1} v_1 + a_{i2} v_2$

$\bar{x} + 1 \cdot v_1 + 1 \cdot v_2$
$= (9.5, 0.5, 3, 7.5)$

| | kale | taco bell | sashimi | pop tarts |
|---|---|---|---|---|
| $x_1 - \bar{x}$ Alice | 4.5 | -3.5 | -3 | 1.5 |
| $x_2 - \bar{x}$ Bob | 1.5 | -3 | -4 | 4.5 |
| Carolyn | -3.5 | 4.5 | 2 | -2.5 |
| $x_4 - \bar{x}$ Dave | -2.5 | 1.5 | 5 | -3.5 |

$\sum (x_i - \bar{x}) = 0$

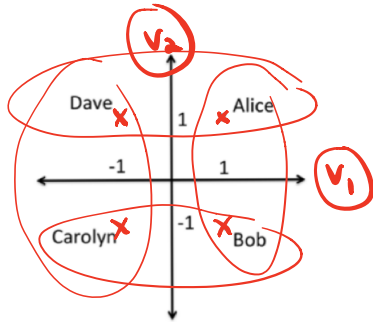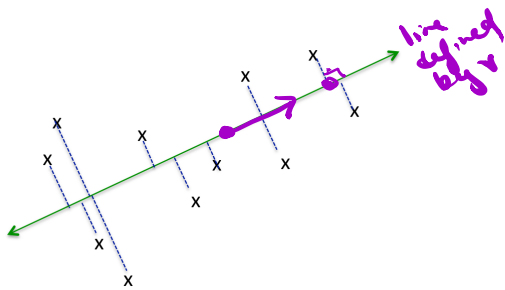| | kale | taco bell | sashimi | pop tarts |
|---|---|---|---|---|
| Alice | 1.5 | -0.5 | 0 | -2.5 |
| Bob | -2.5 | 0 | -1 | 1.5 |
| Carolyn | -0.5 | 1.5 | -1 | .5 |
| Dave | 0.5 | -1.5 | 2 | -0.5 |



Figure 1: Visualizing 4-dimensional data in the plane.

$x_1, ..., x_m$

represent them approx. $v_i$

$x_i \approx \sum_{j=1}^{k} a_{ij} v_j$

JL dim reduction    vs    PCA

- cares about preserving distances    doesn't

- coordinates had no meaning    goal to find meaning in vectors $v_1, ..., v_k$

- $\frac{c \log n}{\varepsilon^2}$ dimensions to preserve distances $1 \pm \varepsilon$    useful even if $k = 1, 2$

goal: to find "intrinsic" dimensionality

x x x x x x x x x x x x x

line defined by v

PCA: Preprocessing

1) subtract out mean

2) from a practical perspective scale coordinates.

$$x_{ij} = \frac{x_{ij}}{\sqrt{\sum_{k=1}^{m} x_{kj}^2}}$$

$k=1:$

Objective: choose $\vec{v}$     $\|\vec{v}\| = 1$
so as to minimize

minimize $\frac{1}{m} \sum_{i=1}^{m} \left[ \text{dist}(x_i \text{ to line defined by } v) \right]^2$     $\equiv$

$\sum_{i=1}^{m}:$ $\left(\langle x_i, v \rangle^2 + \text{dist}^2 = \|x_i\|^2\right)$

Claim: this is the same thing as choosing $v$ of length 1 to maximize

$\frac{1}{m} \sum_{i=1}^{m} \langle x_i, v \rangle^2$

maximize

$\|\mathbf{x}_i\|$     $\mathbf{x}_i$     $\text{dist}(\mathbf{x}_i \leftrightarrow \text{line})$

$O$     $\langle \mathbf{x}_i, \mathbf{v} \rangle$

Figure 4: The geometry of the inner product.

Maximize variance (squared distance) of red dots in this direction

Minimize residuals (squared distance) in this direction

Two equivalent views of principal component analysis.

$X = \langle x_i, v \rangle$ w/ prob $\frac{1}{m}$

$$E(X) = 0 \qquad \frac{1}{m}\sum_i (x_i, v) = \frac{1}{m}\left(\sum x_i, v\right) = 0$$
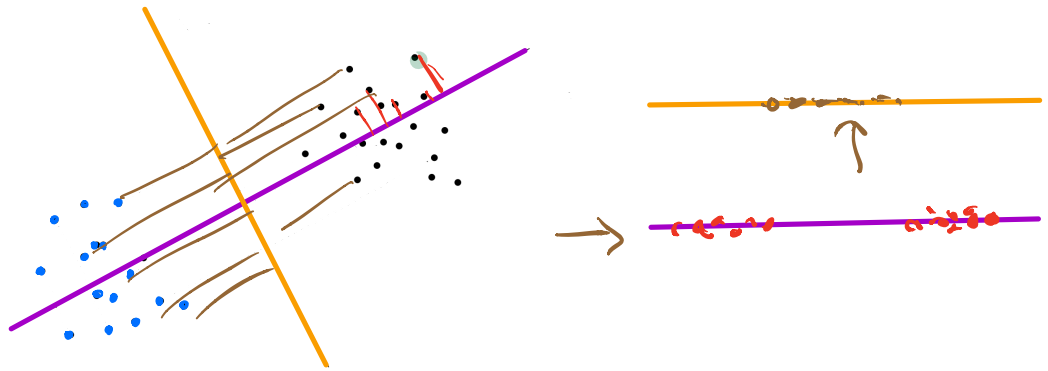
$$Var(X) = E\left((x-\mu)^2\right) = E(x^2)$$



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.
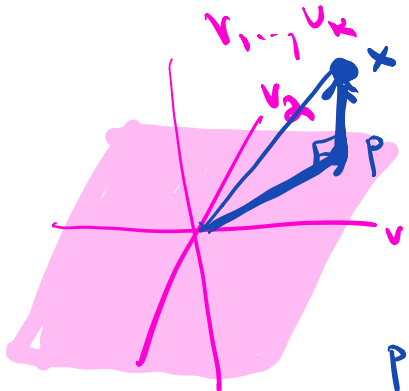
larger $k$   objective

find   $k$-dim   so as to max $\dfrac{1}{m}\displaystyle\sum_{i=1}^{m}\left(\text{length of } x_i\text{'s projection to } S\right)^2$
subspace $S$

Find set of $k$ orthonormal vectors.

$v_1, ..., v_k$    $\|v_i\|^2 = 1$    $(v_i, v_j) = 0 \;\forall\, i \neq j$

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

subspace — set of all vectors that can be written as comb

$$\sum_{j=1}^{k} c_j v_j$$

$P = a_1 v_1 + a_2 v_2$

$(x-p, v_1) = 0$

$u \cdot v$

$(u, v)$

$u^T v$

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

$\qquad = a_1 (v_1, v_1) + a_2 (v_2, v_1) \leftarrow$

$a_1 = (x, v_1)$

$\|p\|^2 = (a_1 v_1 + a_2 v_2, a_1 v_1 + a_2 v_2) = a_1^2 + a_2^2$
$\qquad\qquad = (x, v_1)^2 + (x, v_2)^2$

objective: find $v_1, \ldots, v_k$ orthonormal

to maximize $\frac{1}{m} \sum\limits_{i=1}^{m} \sum\limits_{j=1}^{k} (x_i, v_j)^2$

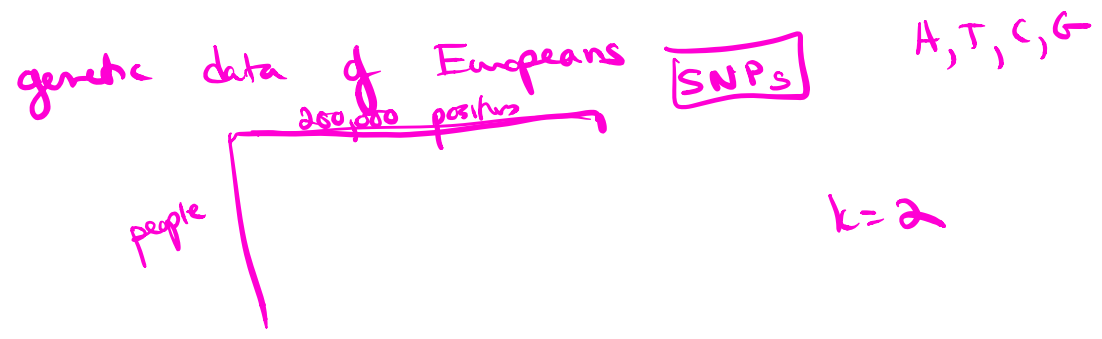avg. sum of squared projected lengths.

## Application: Visualization.

① Perform PCA $\longrightarrow$ $\underline{v_1, \ldots, v_k}$ top k "principal components"

② $\forall x_i$ define "$v_1$ coord" $\left. \begin{matrix} (x_i, v_1) & \leftarrow \\ (x_i, v_2) & \leftarrow \\ \vdots & \\ (x_i, v_k) & \leftarrow \end{matrix} \right.$

"$v_2$ coord"

⋮

"$v_k$ coord"

③ Plot your pts

$$x_i = \left( (x_i, v_1), (x_i, v_2), (x_i, v_3) \right)$$

- look for clusters
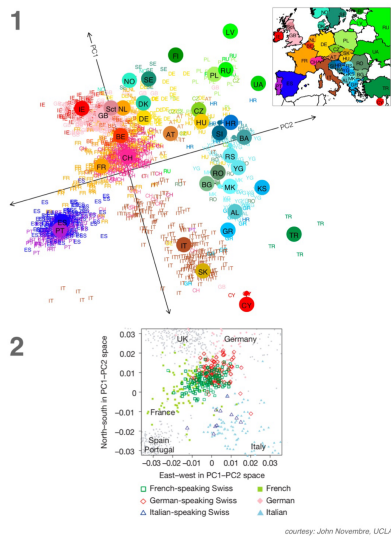- look what are pts particularly large along $v_1$

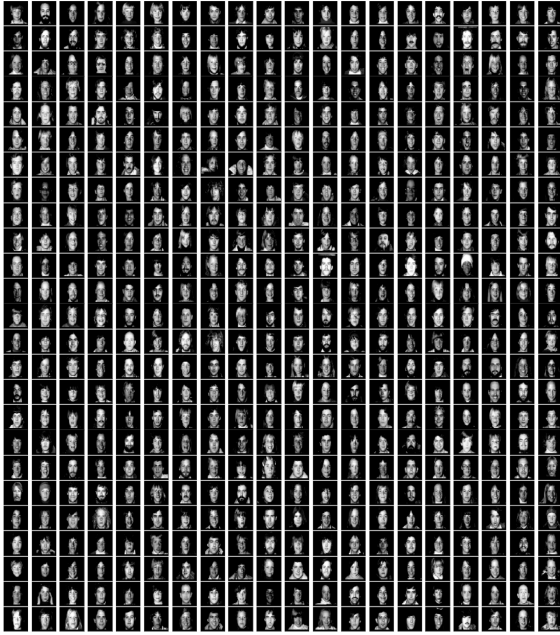genetic data of Europeans   $\boxed{SNPs}$   $A, T, C, G$

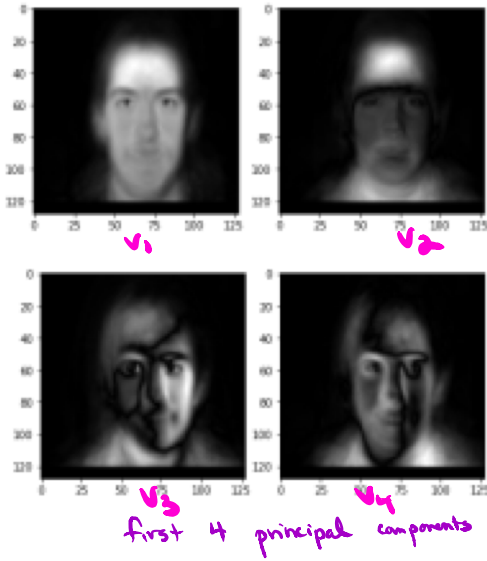250,000 positions

people

$k = 2$

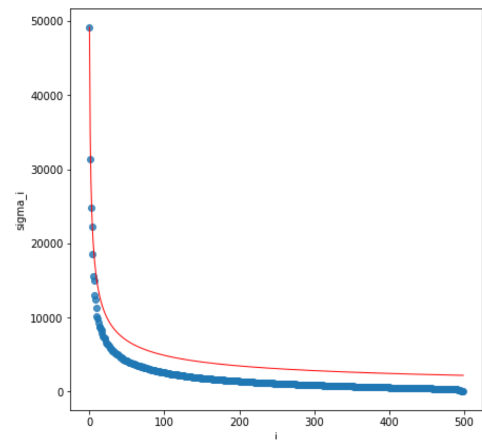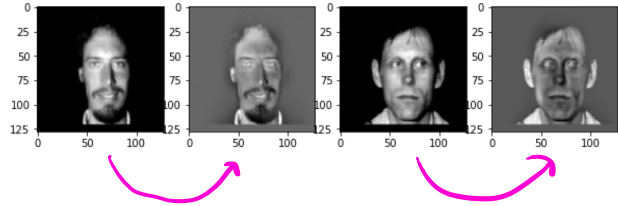Figure 1: The genetic map of Europe using PCA, with the geographic map of Europe for reference. Figure 2: The same map, but zoomed in on Switzerland. Swiss individuals tend to cluster with countries that speak the same language. *(Courtesy: John Novembre, UCLA)*

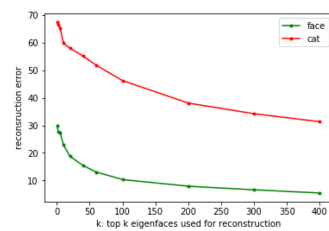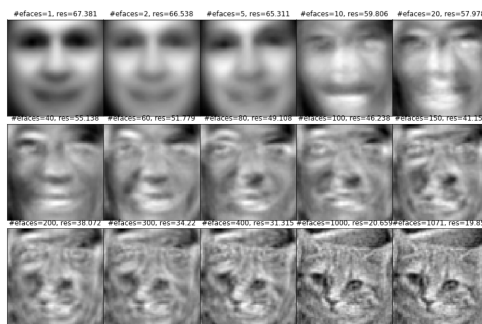Application 2: Compression

Eigenfaces.

500 faces



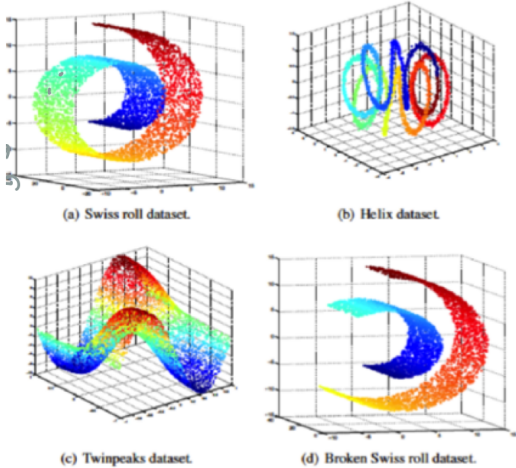2 examples of images after subtracting "mean" face.







v1

v2

v3

v4

first 4 principal components

**Original Face Image**





**eigenfaces space**



#efaces=1, res=29.769   #efaces=2, res=27.586   #efaces=5, res=27.347   #efaces=10, res=23.01   #efaces=20, res=18.755
#efaces=40, res=15.416   #efaces=60, res=13.066   #efaces=80, res=11.821   #efaces=100, res=10.342   #efaces=150, res=8.813
#efaces=200, res=7.924   #efaces=300, res=6.626   #efaces=400, res=5.454   #efaces=1000, res=1.963   #efaces=1071, res=1.617



#efaces=1, res=57.804   #efaces=2, res=57.611   #efaces=5, res=54.054   #efaces=10, res=52.01   #efaces=20, res=45.897
#efaces=40, res=35.868   #efaces=60, res=29.624   #efaces=80, res=24.103   #efaces=100, res=20.317   #efaces=150, res=16.154
#efaces=200, res=13.257   #efaces=300, res=9.581   #efaces=400, res=6.908   #efaces=1000, res=0.924   #efaces=1071, res=0.653



#efaces=1, res=67.381   #efaces=2, res=66.538   #efaces=5, res=65.311   #efaces=10, res=59.806   #efaces=20, res=57.978
#efaces=40, res=55.138   #efaces=60, res=51.779   #efaces=80, res=49.108   #efaces=100, res=46.238   #efaces=150, res=41.153
#efaces=200, res=38.072   #efaces=300, res=34.22   #efaces=400, res=31.315   #efaces=1000, res=20.659   #efaces=1071, res=19.855

(a) Swiss roll dataset.

(b) Helix dataset.

(c) Twinpeaks dataset.

(d) Broken Swiss roll dataset.

How PCA works

$k=1$     find $v$ st. $\|v\|=1$

to maximize $\quad \frac{1}{m} \boxed{\sum_{i=1}^{m} (x_i, v)^2}$

$$X = \begin{pmatrix} -x_1- \\ \vdots \\ -x_m- \end{pmatrix} \Leftarrow \qquad Xv = \begin{pmatrix} (x_1, v) \\ \vdots \\ (x_m, v) \end{pmatrix} \Leftarrow$$

$$(Xv)^T \qquad \boxed{\left( (x_1,v) \cdots (x_m,v) \right) \cdot Xv}$$

$$\underline{(Xv)^T \; Xv}$$

$(Xv)^T = v^T X^T$

$\left( (x_1, v) \quad \cdots \quad (x_m, v) \right) \begin{pmatrix} (x_1, v) \\ \vdots \\ (x_m, v) \end{pmatrix}$ 

$Xv$

$\underbrace{v^T X^T X}_{A} \, v = \sum_{i=1}^{m} (x_i, v)^2$

$\boxed{A = X^T X} \longrightarrow$ correlation covariance

find $v$ to

max $\quad v^T A v$

$\boxed{A_{k\ell} = \sum_{i=1}^{m} x_{ik} \, x_{i\ell}}$

$k \begin{pmatrix} | & & | \\ x_1 & \cdots & x_m \\ | & & | \end{pmatrix} \begin{pmatrix} -x_1- \\ \vdots \\ -x_m- \end{pmatrix}$ $\ell^{th}$

$A$ is symmetric.

Suppose rows of $X$ are documents, cols words

max $\underline{\underline{v^T A v}}$

$\|v\| = 1$

Suppose $A = \begin{pmatrix} 2 & & 0 \\ & 1 & \\ 0 & & \frac{1}{2} \end{pmatrix}$

$v = (v_1, v_2, v_3)$

$v^T A v$

$\rightarrow \dfrac{2 \cdot \boxed{v_1^2} + 1 \cdot \boxed{v_2^2} + \frac{1}{2} \cdot \boxed{v_3^2}}{v_1^2 + v_2^2 + v_3^2 = 1}$

$v = (1, 0, 0)$

$v_1 = 1$

Every symmetric

$\begin{pmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_n \end{pmatrix}$

$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
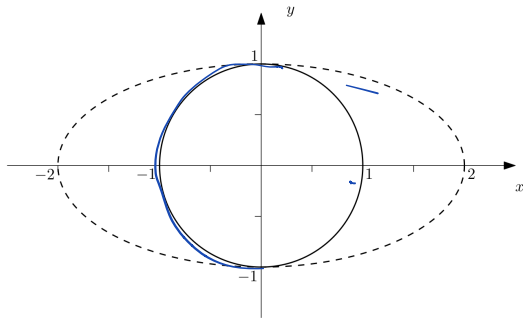
$\dfrac{\sum \lambda_i v_i^2}{\sum v_i^2 = 1} =$

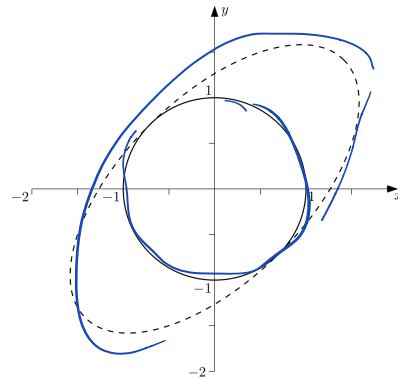Figure 1: The point $(x, y)$ on the unit circle is mapped to $(2x, y)$.



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

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

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

$$A = Q \ D \ Q^T$$

$n \times n$

diagonal

$$Q = \begin{pmatrix} q_1 & q_2 & \cdots & q_n \end{pmatrix} \quad e_1 \qquad Q^T = \begin{pmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{pmatrix}$$

orthogonal $\equiv$ cols are orthonormal

$\|q_i\|^2 = 1 \qquad (q_i, q_j) = 0 \quad \forall i \neq j$

special

$$(Q v)^T Q v$$

$$= v^T \underbrace{Q^T Q}_{I} v = v^T v$$

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

orthogonal matrices preserve length.

$$\max \quad v^T A v = v^T Q \underbrace{D Q^T v}$$

$$y = Q^T v \quad \Leftarrow \quad \sum y_j^2 = 1$$

$$\text{if } \|v\| = 1 \Rightarrow \|y\| = 1$$

$$y^T D y = \sum_{j=1}^{n} \lambda_j \boxed{y_j^2}$$

$$y_1 = 1 \quad y_j = 0 \quad j > 1$$

$$y = e_1 \qquad e_1 = Q^T v \qquad \lambda_1 > \lambda_2 > \cdots$$

$$Q e_1 = \underbrace{Q Q^T}_{I} v$$

$$v = Q e_1 = q_1 \quad \left(\text{first col of } Q\right)$$

---

$M . z$   $z$ is eigenvector w/ eigenvalue $\lambda$

$\qquad$ if $\qquad \underline{Mz = \lambda z}$

Observatn $\quad \boxed{Q e_1}$ is eigenvector of matrix $A$ corresponding to eigenvalue $\lambda_1$

$$\longrightarrow Q e_j \qquad '' \qquad '' \qquad \lambda_j$$

$$\underline{A Q e_j} = Q D Q^T Q e_j = Q D e_j$$
$$= Q \lambda_j e_j$$
$$= \lambda_j \underline{Q e_j}_{j^{th} col of Q}$$

$$v^T A v$$

soln for $k=1$ : largest eigenvector of $\hat{A} = X^T X$

$$A = Q D Q^T$$

to do PCA  w/ $k=1$

find unit vector $v$ that maximizes $v^T \underbrace{X^T X}_{A} v$

principal eigenvector of $A$

$$A = Q D Q^T$$

set of all eigenvectors

$v$ $\boxed{Q e_1 = q_1}$ $\lambda_1$

$\quad \vdots$

$Q e_n = q_n$ $\quad \lambda_n$

diag of $D$.
all nonnegative.

$u_0$ random vector.   eigenvectors $q_1, \dots, q_n$ are a basis

$$\boxed{u_0 = \sum_{j=1}^{n} c_j q_j}$$

$$A u_0 = \sum_j c_j A q_j = \sum_j c_j \lambda_j q_j =$$

$$A \sum_j c_j \lambda_j q_j = \sum_j c_j \lambda_j A q_j = \sum_j c_j \lambda_j^2 q_j$$

$$A \sum_j c_j \lambda_j^2 q_j = \sum_j c_j \lambda_j^3 q_j$$

$$A^k u_0 = \sum_j c_j \lambda_j^k q_j$$

$$= \lambda_1^k \left[ \underline{c_1 q_1} + \underbrace{c_2 \frac{\lambda_2^k}{\lambda_1^k} q_2} + c_3 \frac{\lambda_3^k}{\lambda_1^k} q_3 + \dots \right]$$

$\dots \leq \lambda_3 \leq \lambda_2 < \lambda_1$

$$= \boxed{c_1 \lambda_1^k q_1}$$

rescale $\Rightarrow q_1$

**Algorithm 1**
POWER ITERATION

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

- Select random unit vector $\mathbf{u}_0$

- For $i = 1, 2, \ldots$, 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 $\mathbf{u}_i / \|\mathbf{u}_i\|$.

*vector not changing.*

$O\left( \dfrac{\log n}{\log \left| \frac{\lambda_1}{\lambda_2} \right|} \right)$ iterations.

$A^2 u_0$

$A \quad A^2 A^4 A^8$

$u_1 = A u_0$
$u_2 = A u_1$
$u_3 = A u_2$
$\vdots$
$u_i = A^i u_0$

$\lambda_1 > \lambda_2$

Pagerank.



0.0005   webpage   $\frac{1}{2k}$   0.3.

$2k$   $2^k$   webpage   $0$

$\pi$

$j$   $\frac{1}{k} \ \frac{1}{k} \ \frac{1}{k}$

$i$   $1$

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

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

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

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.

*Singular value decomposition.*

$$O(n^2 m)$$
$$O(m^2 n)$$



**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.