

Features?

Images



Videos



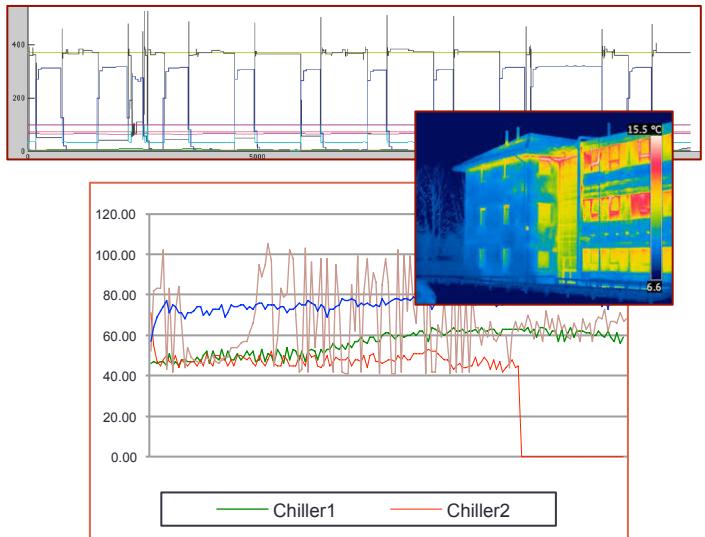
Social network



Books

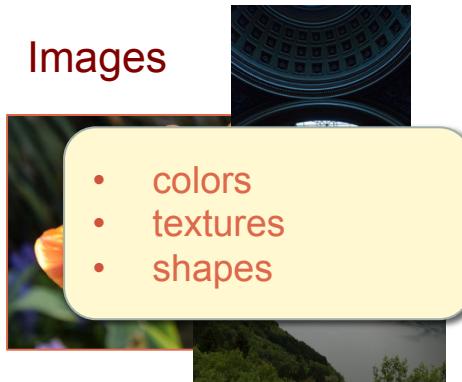


Sensor readings



Features?

Images



Videos



- actors
 - ratings
 - directors

Social network

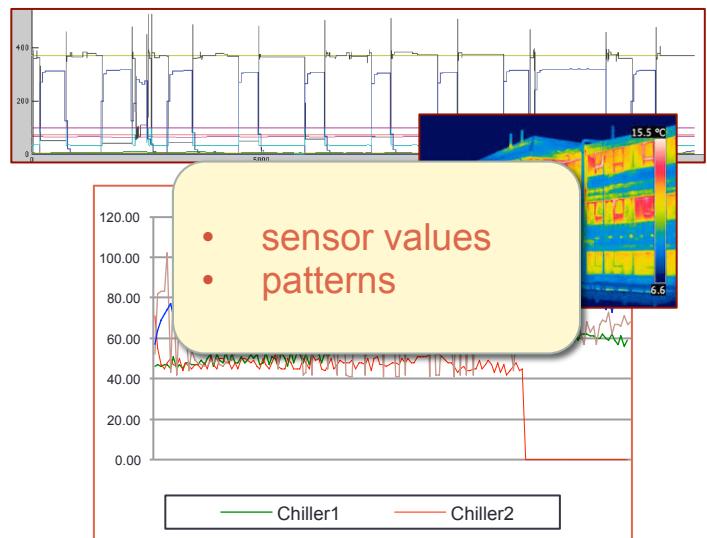


Books



- words
 - authors
 - publishers

Sensor readings

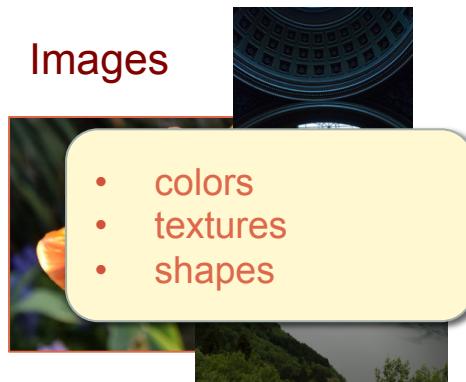


Data representation?

- Sets
 - Does a feature instance **exists** in a given object?
- Vectors
 - How many **times** a feature instance exists in the given object?
 - What is the **strength** of a feature instance is in the given object?
- Strings
 - What is the **order** of feature instances in the given object?
- Trees
 - What is the **hierarchy** of the feature instance in the given object?
- Graphs
 - What are the **pairwise relationships** among feature instances in the given object?

Data representation?

Images



Videos

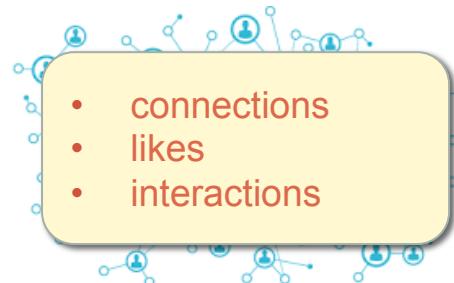


- actors
- ratings
- directors



count

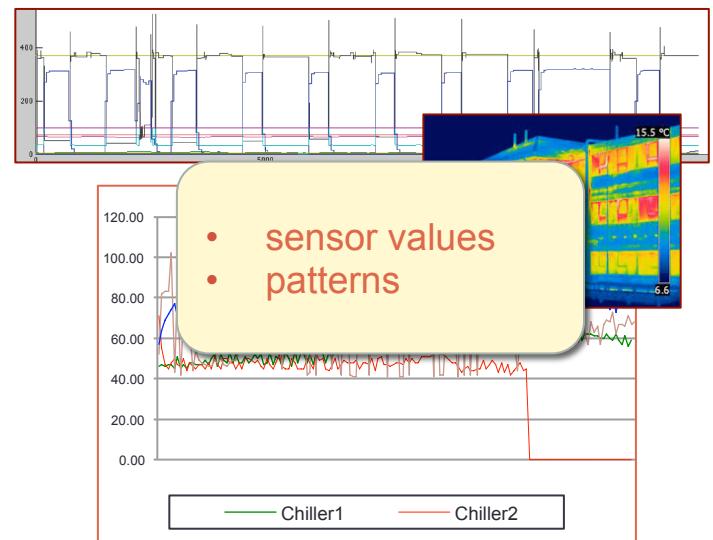
Social network



Books



Sensor readings



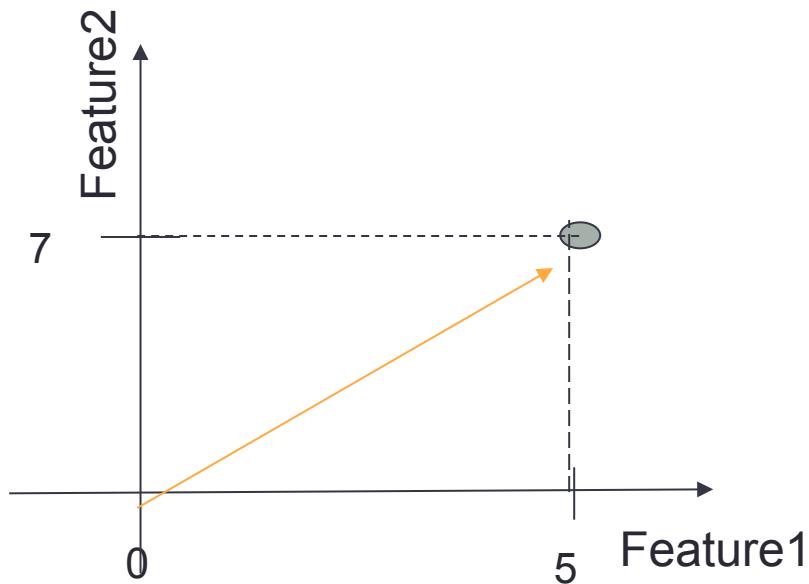
How do we represent counts?

- Feature1 occurs 5 times



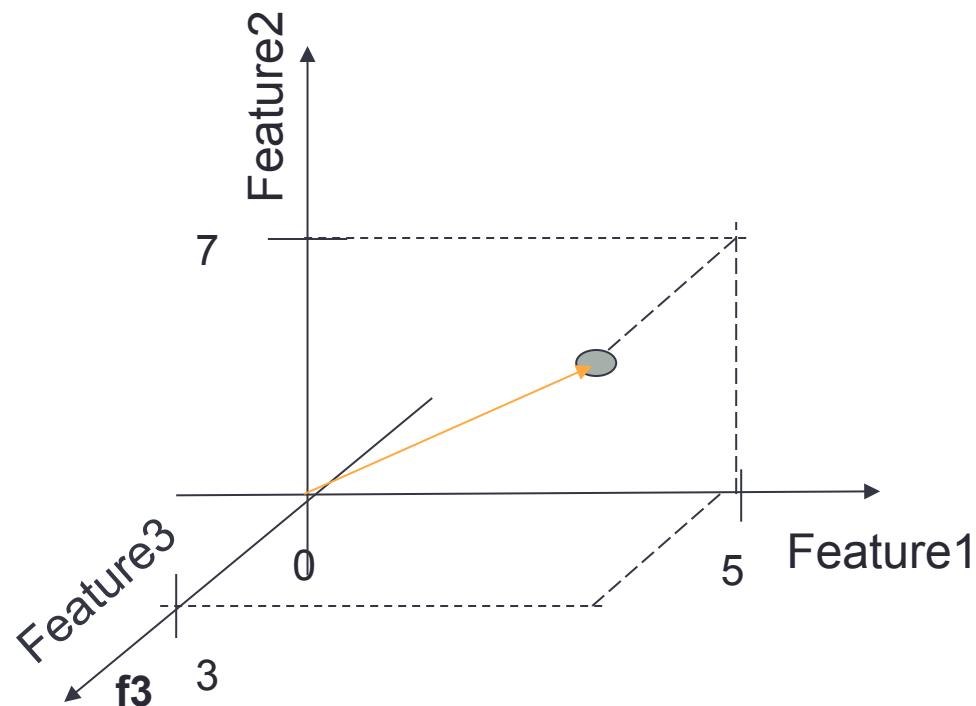
How do we represent counts?

- Feature1 occurs 5 times
- Feature2 occurs 7 times

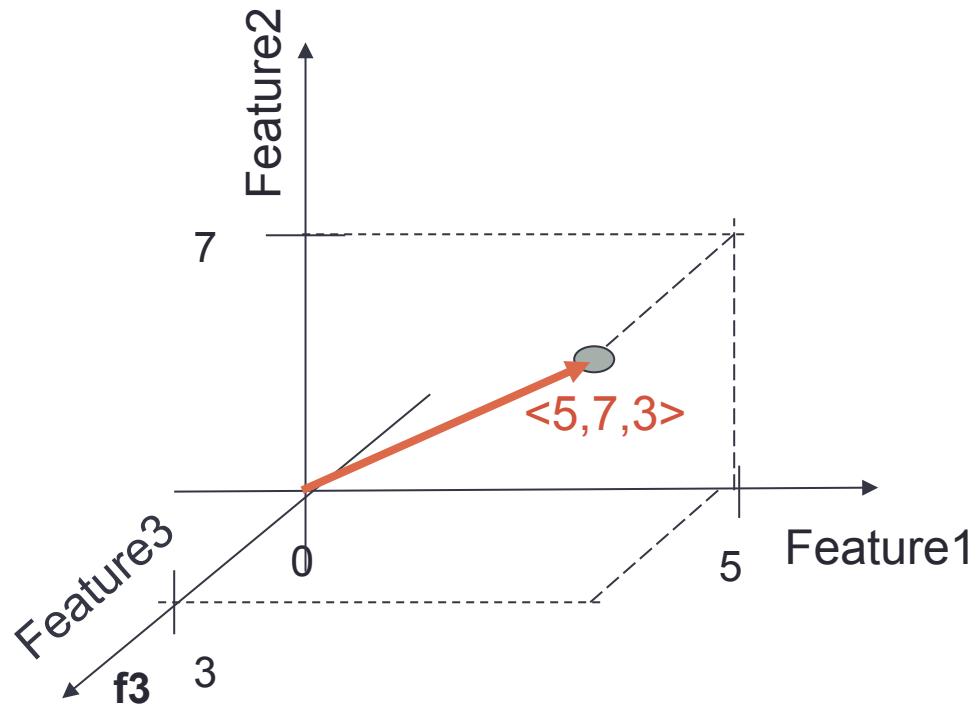


How do we represent counts?

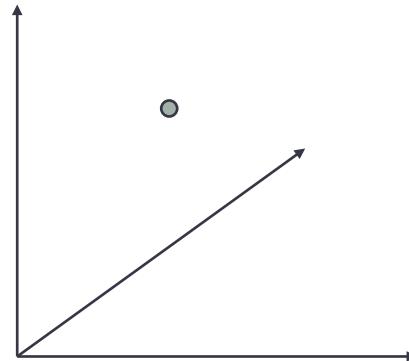
- Feature1 occurs 5 times
- Feature2 occurs 7 times
- Feature3 occurs 3 times



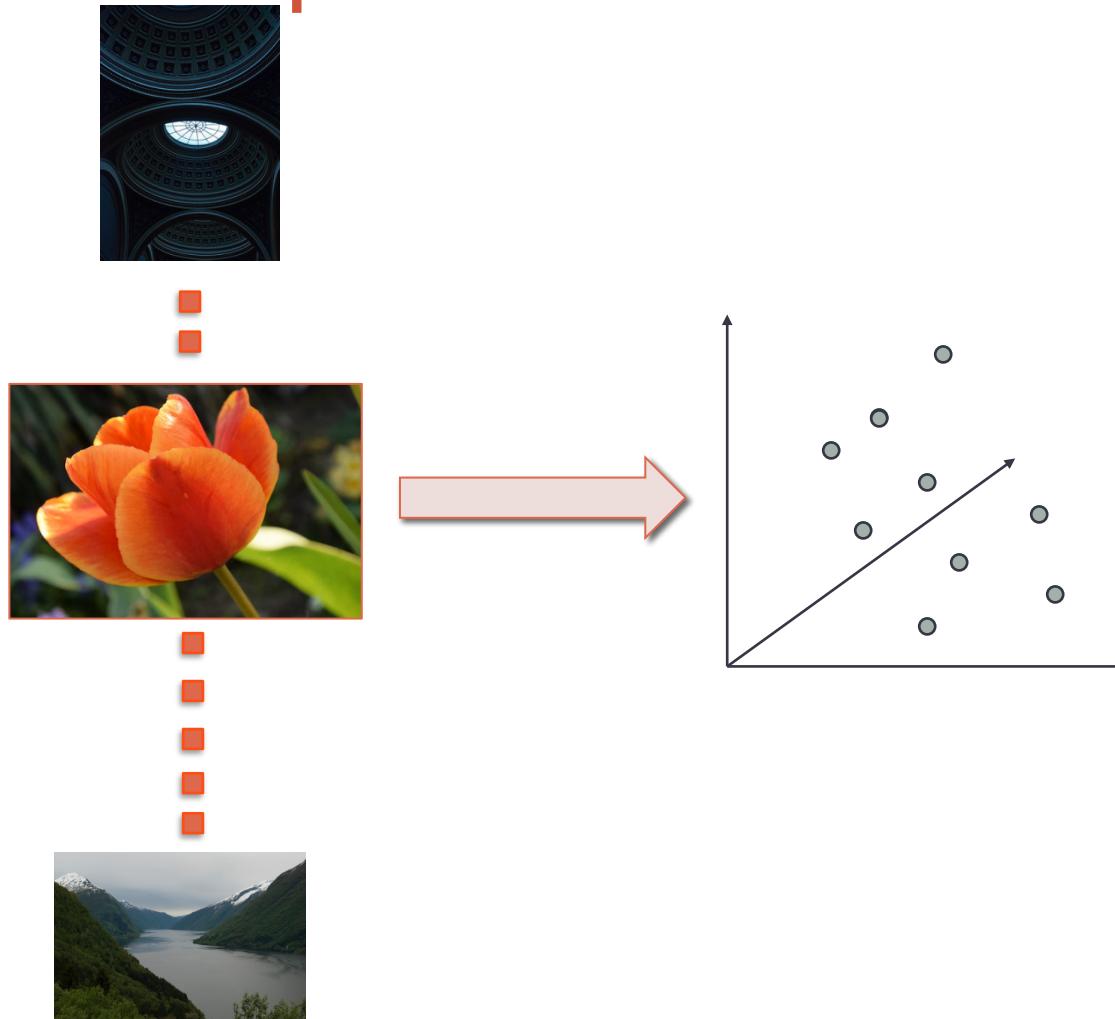
Vector representation



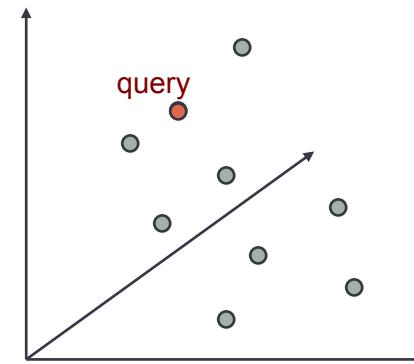
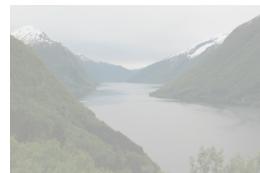
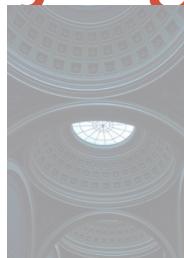
Vector representation



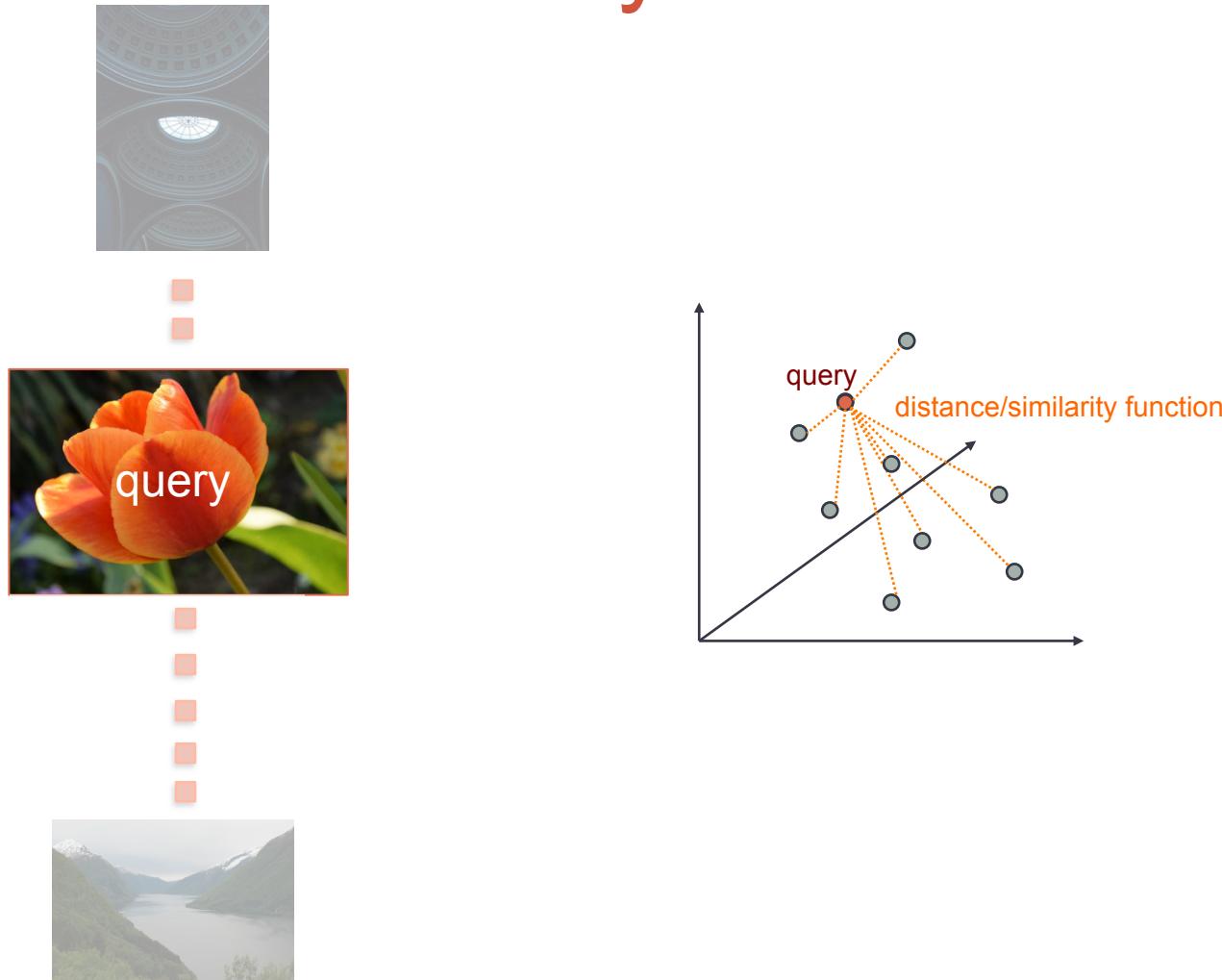
Vector representation



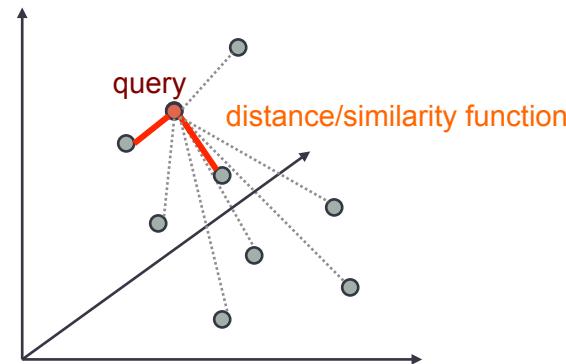
Querying



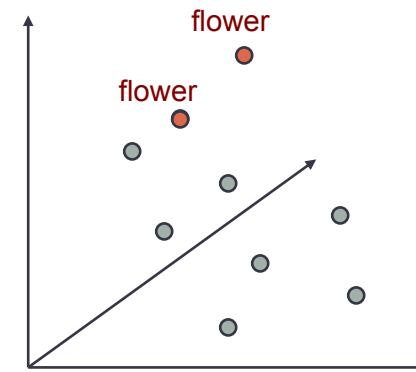
Distance/similarity measurement



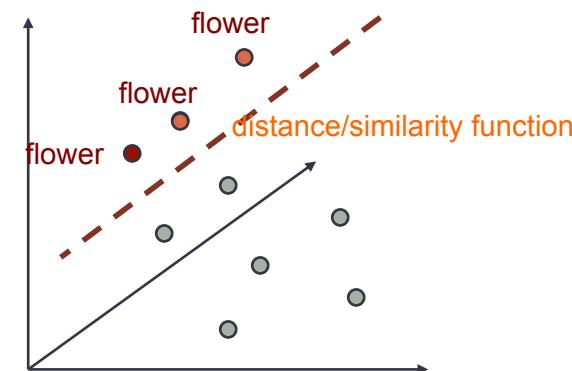
Top-k retrieval...



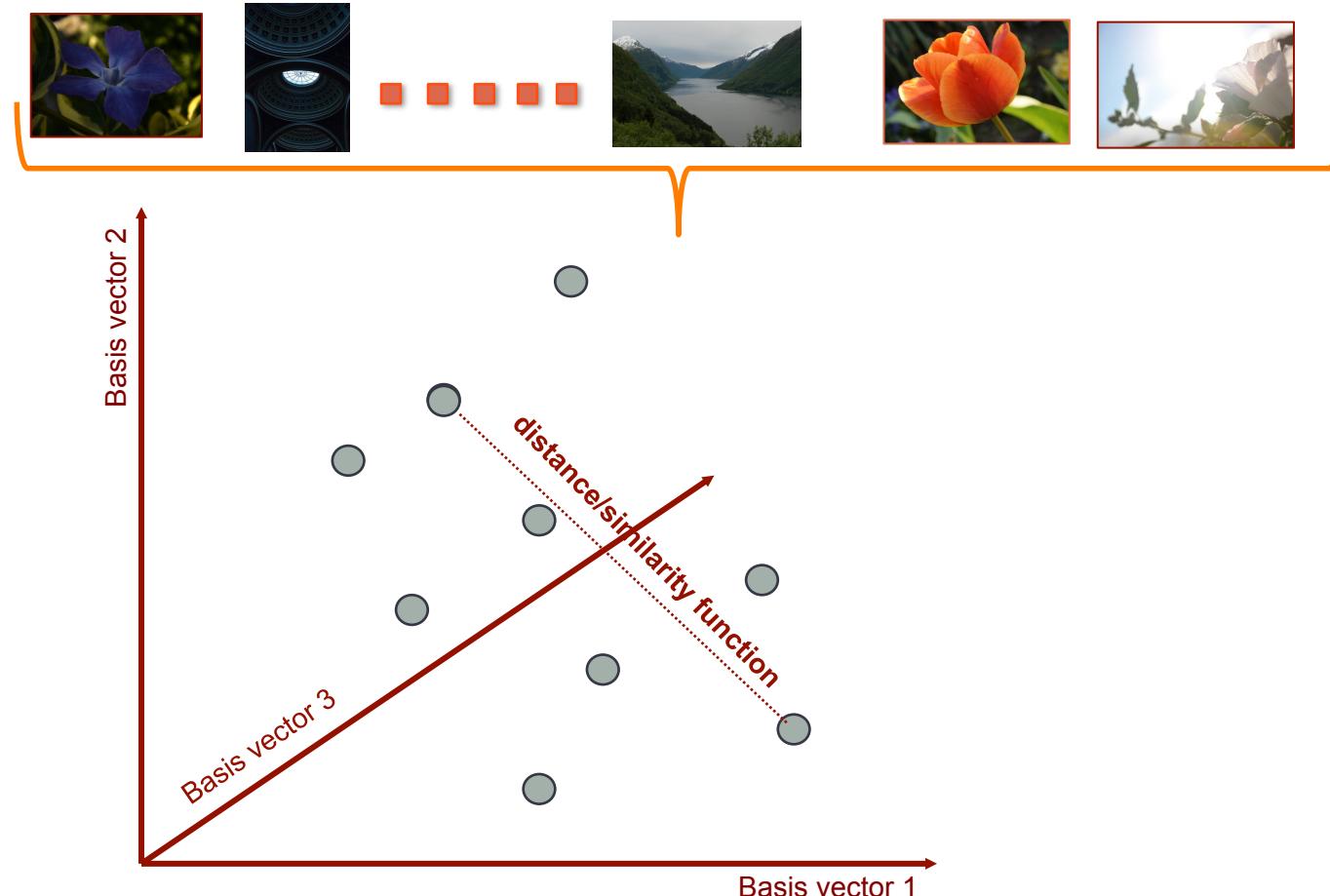
Classification



Classification



Vector spaces



Vector Space

Definition: **(Vector space):** The set \mathbb{S} is a vector space iff for all $\vec{v}_i, \vec{v}_j, \vec{v}_k \in \mathbb{S}$ and for all $c, d \in \mathbb{R}$, the following axioms hold:

- $\vec{v}_i + \vec{v}_j = \vec{v}_j + \vec{v}_i$
- $(\vec{v}_i + \vec{v}_j) + \vec{v}_k = \vec{v}_j + (\vec{v}_i + \vec{v}_k)$
- $\vec{v}_i + \vec{0} = \vec{v}_i$ (for some $\vec{0} \in \mathbb{S}$)
- $\vec{v}_i + (-\vec{v}_i) = \vec{0}$ (for some $-\vec{v}_i \in \mathbb{S}$)
- $(c + d)\vec{v}_i = (c\vec{v}_i) + (d\vec{v}_i)$
- $c(\vec{v}_i + \vec{v}_j) = c\vec{v}_i + c\vec{v}_j$
- $(cd)\vec{v}_i = c(d\vec{v}_i)$
- $1.\vec{v}_i = \vec{v}_i$

The elements of \mathbb{S} are called vectors.

Basis of a Vector Space

Definition (Linear independence and basis): Let $V = \{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n\}$ be a set of vectors in a vector space \mathbb{S} . The vectors in V are said to be linearly independent if

$$\left(\sum_{i=1}^n c_i \vec{v}_i = \vec{0} \right) \iff c_1 = c_2 = \dots = c_n = 0.$$

The linearly independent set V is said to be a basis for \mathbb{S} if for every vector, $\vec{u} \in \mathbb{S}$, there exist constants c_1 through c_n such that

$$\vec{u} = \sum_{i=1}^n c_i \vec{v}_i.$$

- The basis functions should be
 - non-redundant
 - complete

Inner product

Definition (Inner product and orthogonality): *The inner product on a vector space \mathbb{S} is a function $\mathbb{S} \times \mathbb{S} \rightarrow \mathbb{R}$ such that*

- $\vec{u} \cdot \vec{v} = \vec{v} \cdot \vec{u}$,
- $(c_1\vec{u} + c_2\vec{v}) \cdot \vec{w} = c_1(\vec{u} \cdot \vec{w}) + c_2(\vec{v} \cdot \vec{w})$, and
- $\forall_{\vec{v} \neq \vec{0}} \vec{v} \cdot \vec{v} > 0$.

The vectors \vec{u} and \vec{v} are said to be orthogonal if

$$\vec{u} \cdot \vec{v} = 0.$$

- A set of vectors whose projections onto each other are of 0 length can be used as basis vectors (if they are also complete)

A set of mutually orthogonal vectors are also linearly independent

Norm

Definition (Norms and orthonormal basis): A norm (*commonly denoted as $\|\cdot\|$*) is a function that measures the length of vectors. A vector, \vec{v} , is said to be normalized if $\|\vec{v}\| = 1$. A basis, $V = \{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n\}$, of the vector space \mathbb{S} is said to be orthonormal if

$$\forall_{\vec{v}_i, \vec{v}_j} \quad \vec{v}_i \cdot \vec{v}_j = \delta_{i,j},$$

such that if $i = j$, $\delta_{i,j} = 1$ and 0 otherwise.²

P-norms

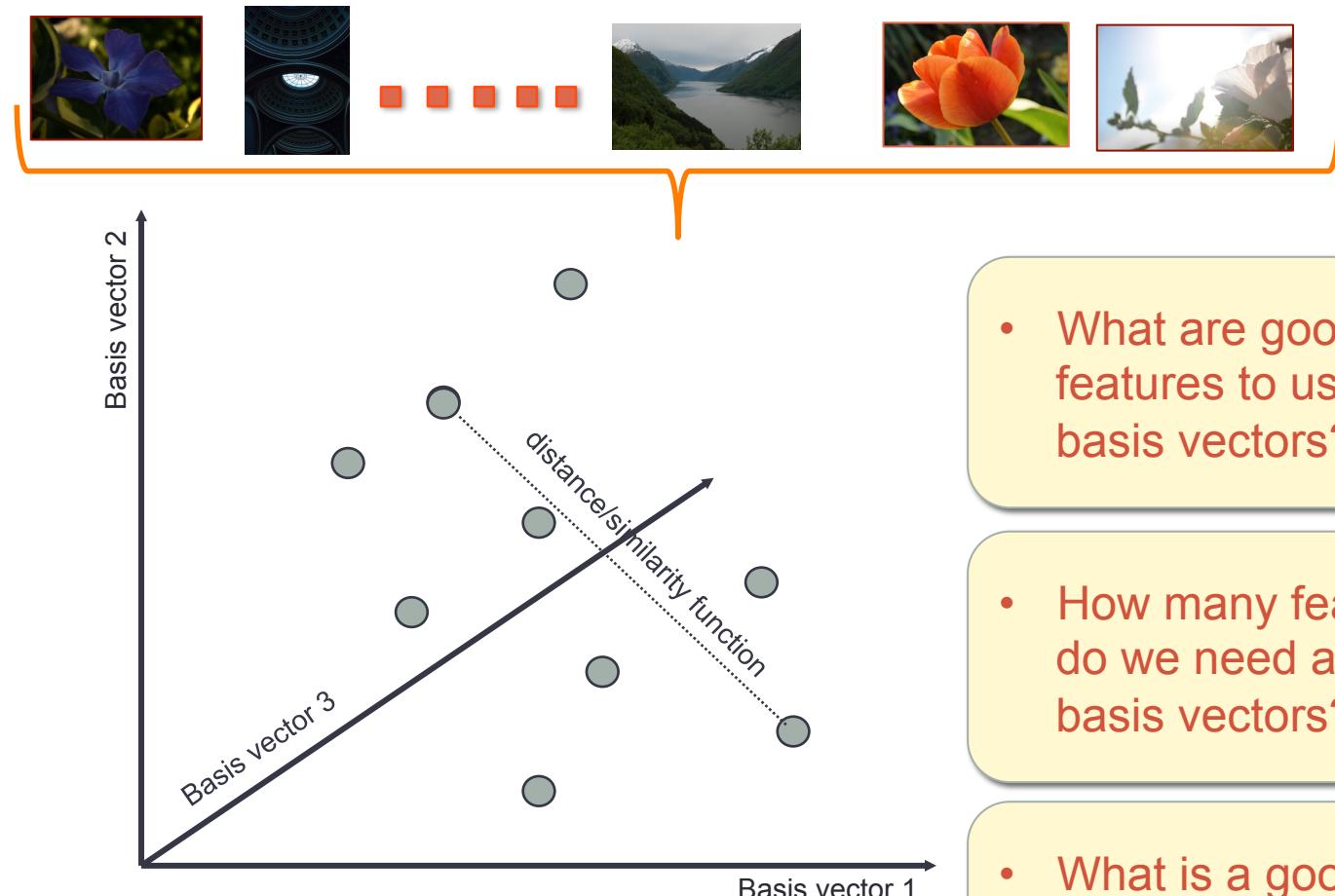
The most commonly used family of norms are the p -norms. Given a vector $\vec{v} = \langle w_1, \dots, w_n \rangle$, the p -norm is defined as

$$\|\vec{v}\|_p = \left(\sum_{i=1}^n |w_i|^p \right)^{\frac{1}{p}}.$$

At the limit, as p goes to infinity, this gives the max-norm

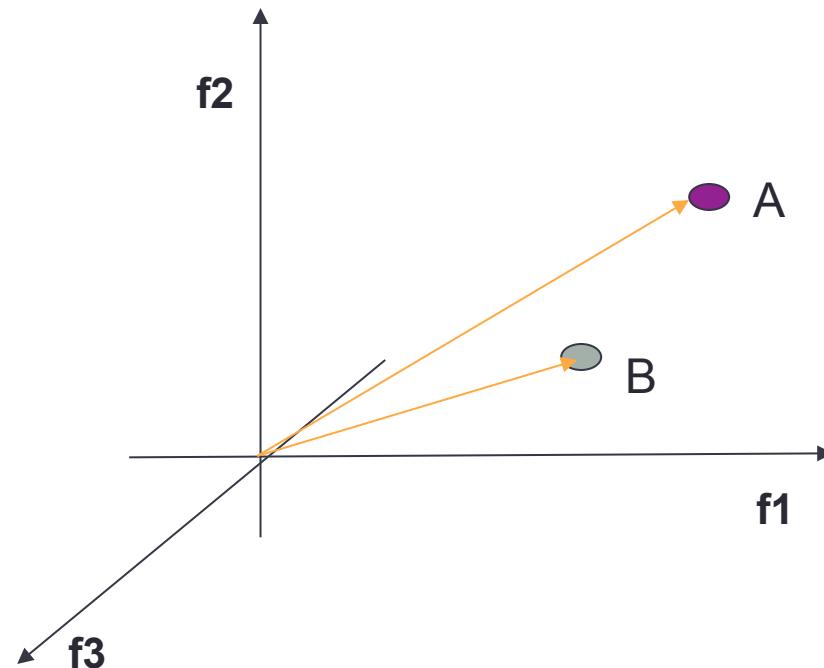
$$\|\vec{v}\|_\infty = \max_{i=1\dots n} \{|w_i|\}.$$

Vector spaces



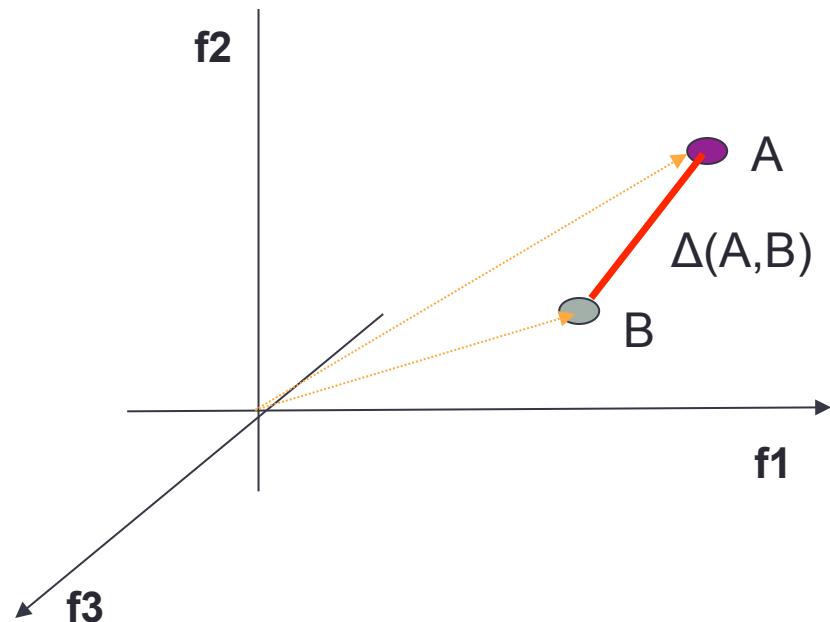
Distance between two images???

- Given A and B vectors, can we measure how different they are?

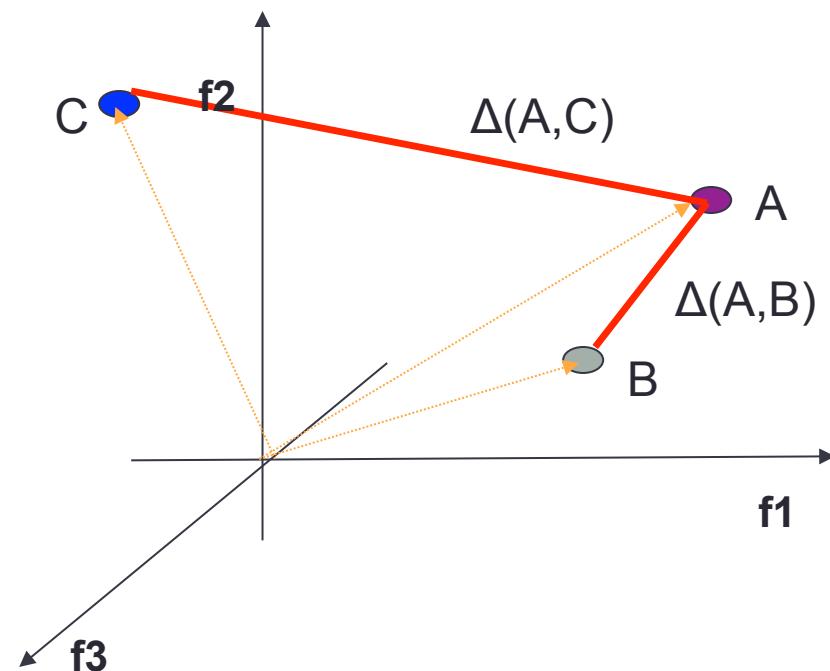


Euclidean distance

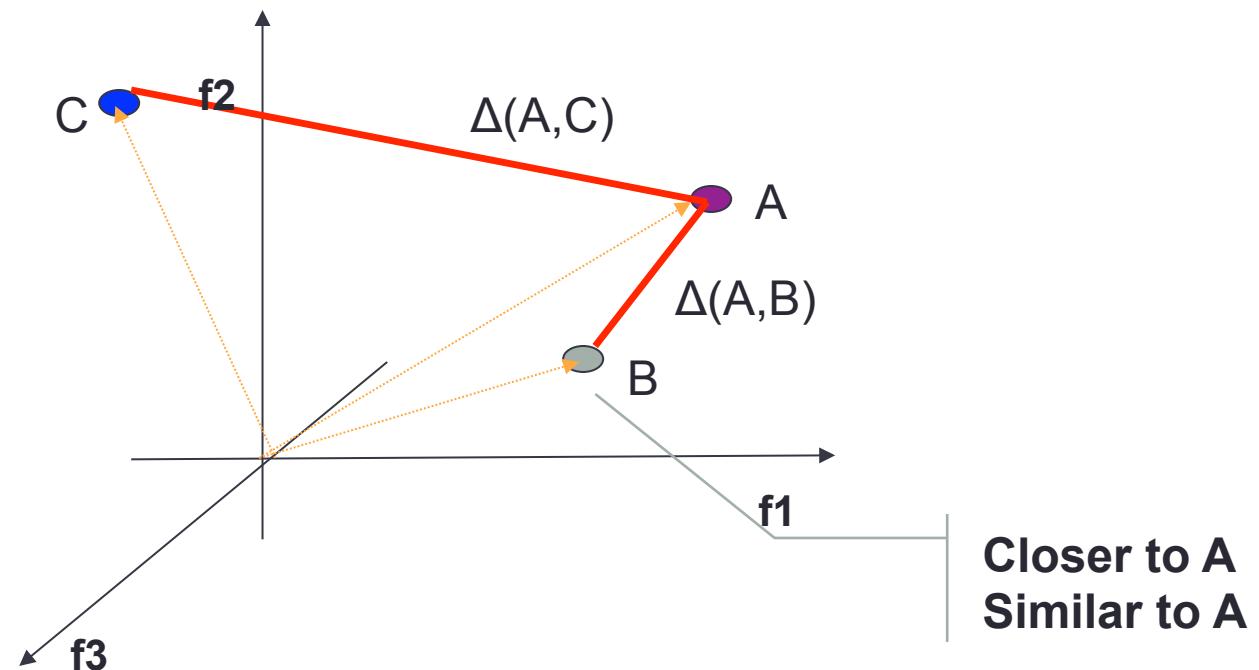
$$\Delta(A, B) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + (a_3 - b_3)^2}$$



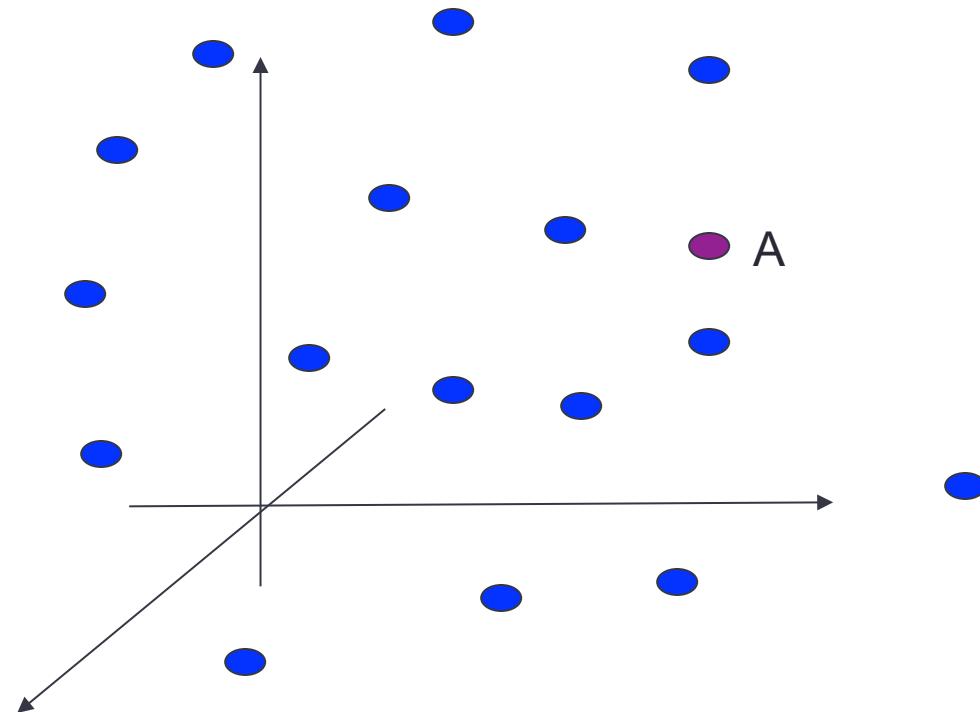
Which image is more similar to A?



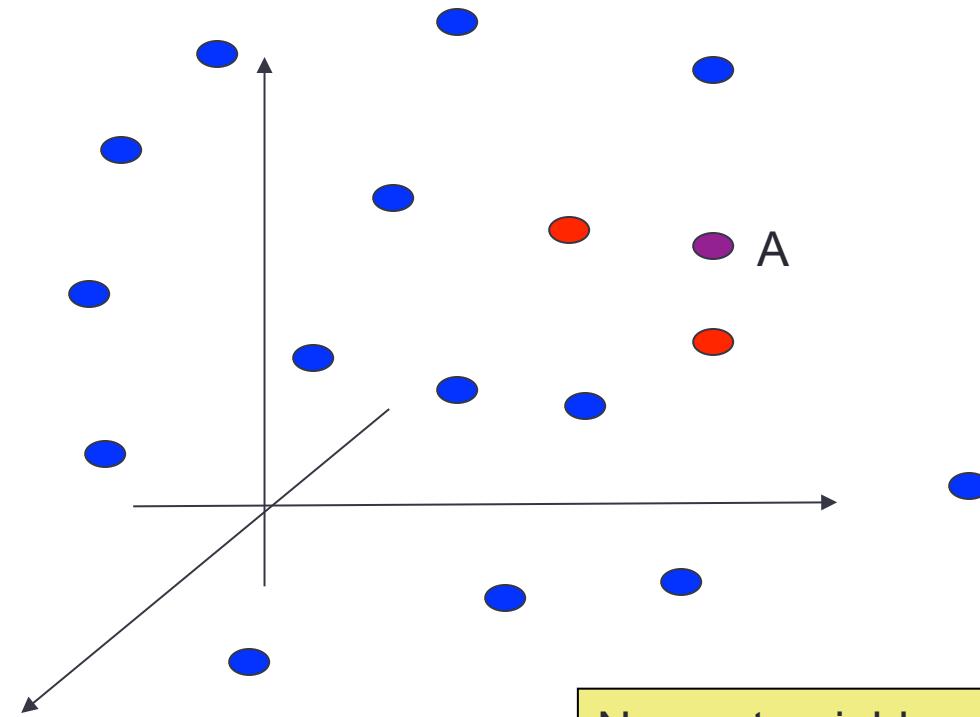
Which image is more similar to A?



“Find 2 most similar images to A”

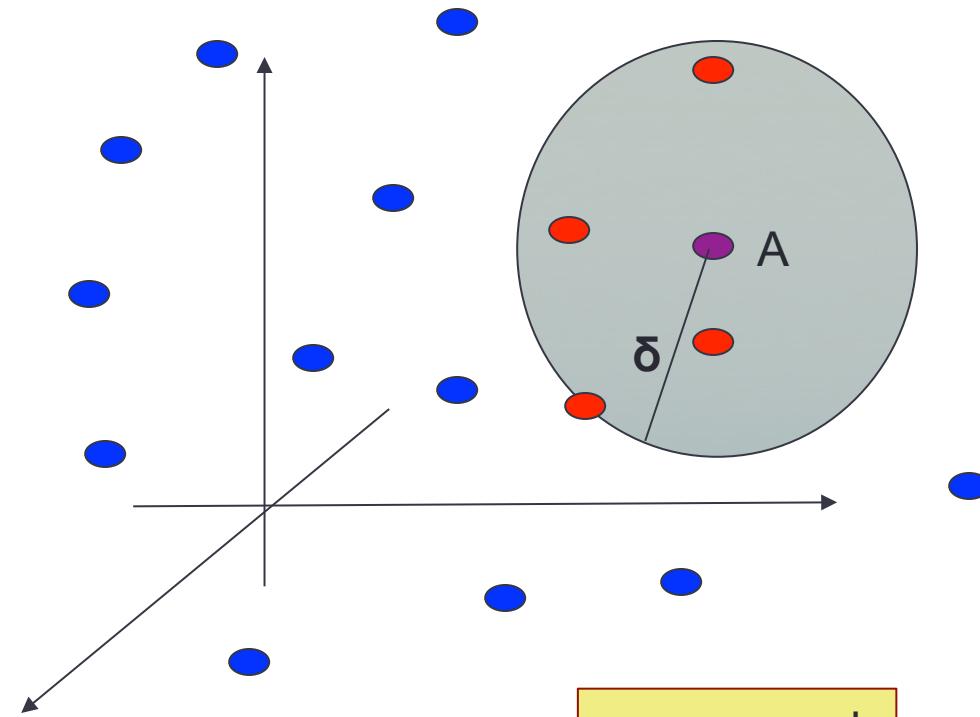


“Find 2 most similar images to A”



Nearest-neighbor search

“Find images at most δ different from A”



range search

What is a good measure then??

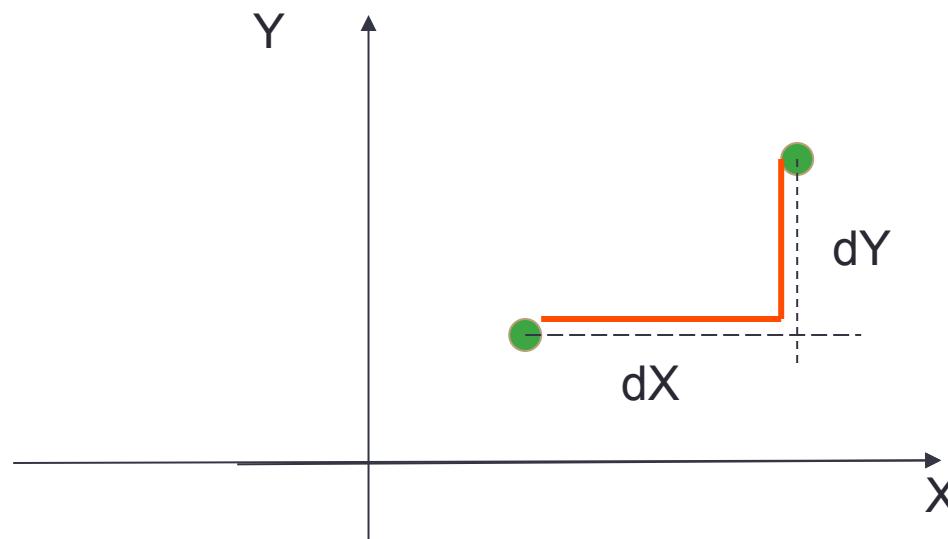
- Application dependent...
- ...but, distances in a metric space help indexing!

Metric model: axioms

- Any function d expressing a metric distance must satisfy the following axioms:
 - self-minimality: $d(s,s) = 0$
 - minimality $d(s_1,s_2) \geq d(s_1,s_1)$
 - symmetry $d(s_1, s_2) = d(s_2, s_1)$
 - triangular inequality $d(s_1,s_2) + d(s_2,s_3) \geq d(s_1,s_3)$

Metric distances (Minkowski metrics)

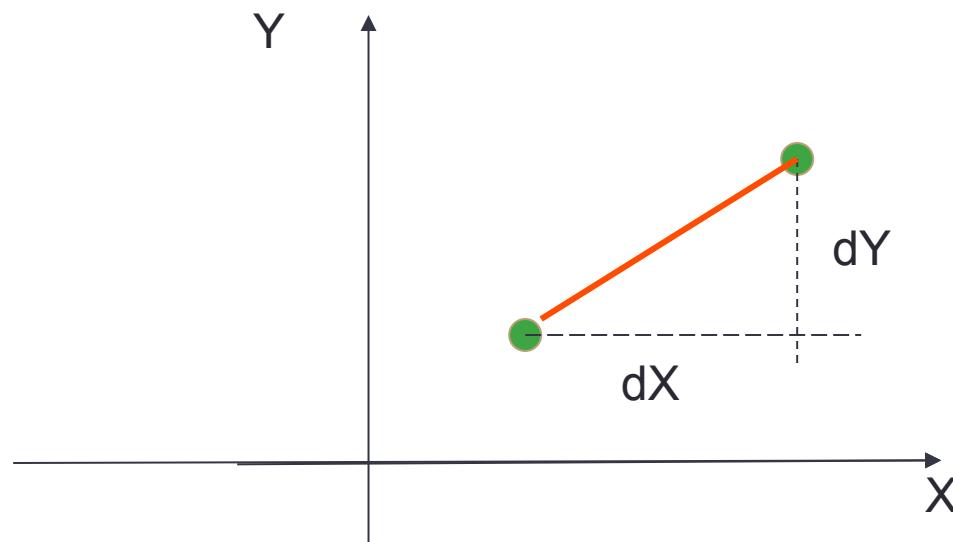
- L1-metric: $d = (dX+dY)$



Also called Manhattan Distance

Metric distances (Minkowski metrics)

- L2-metric: $d = (dX^2+dY^2)^{1/2}$



Also called Euclidean Distance

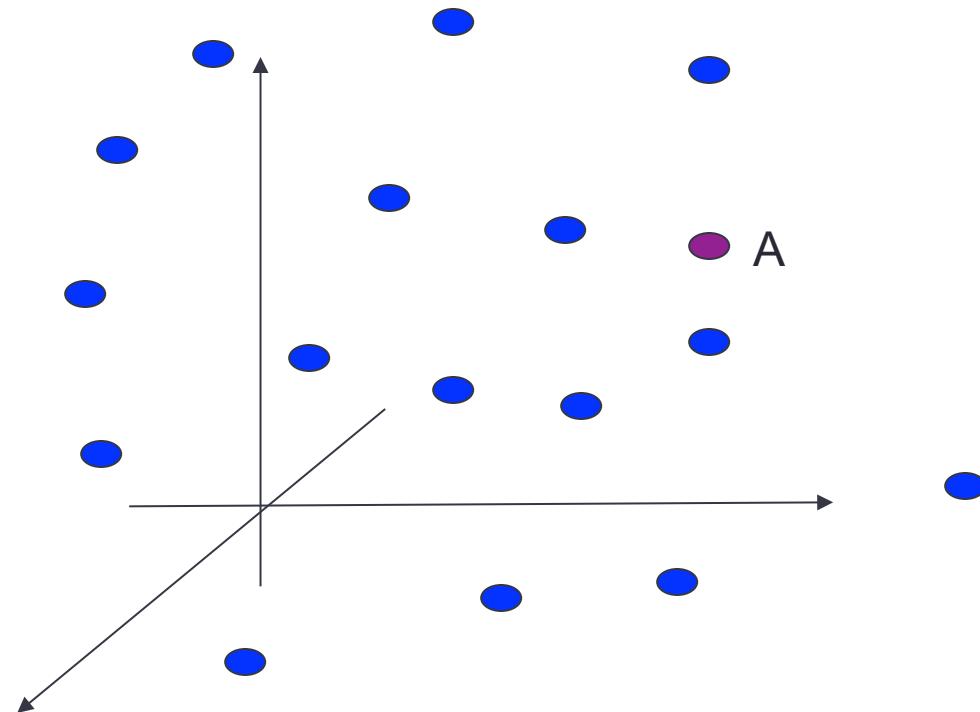
Metric distances (Minkowski metrics)

- L3-metric; $d = (dX^3+dY^3)^{1/3}$
-
-
- L(infinity): $d = \max\{X,Y\}$

...metric model

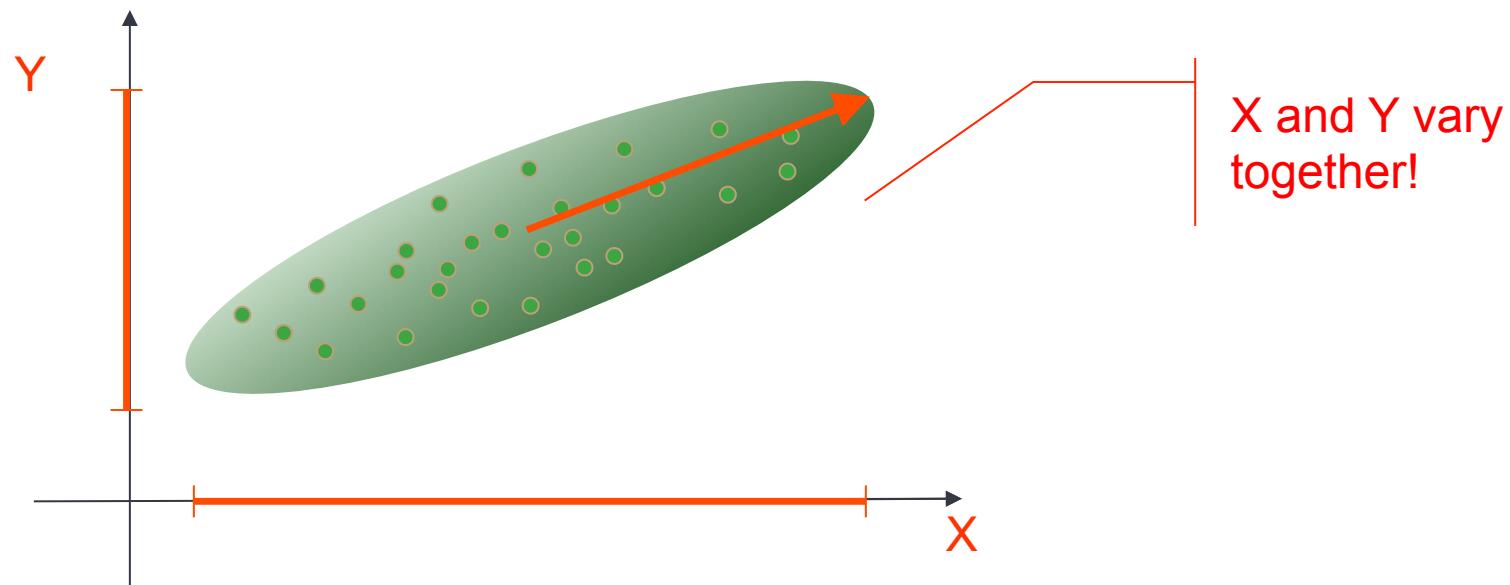
- Well suited for certain kinds of similarity evaluation, such as color based comparisons
- Consistent with widely used approaches from computer vision and pattern recognition communities
 - results suggest that the L1 metric may better capture human notions of image similarity.
- Makes it relatively easy to index data, modeled as vectors of properties, in terms of classical multi-dimensional indexing techniques.

Are there other similarity/distance measures?



Mahalonobis Distance

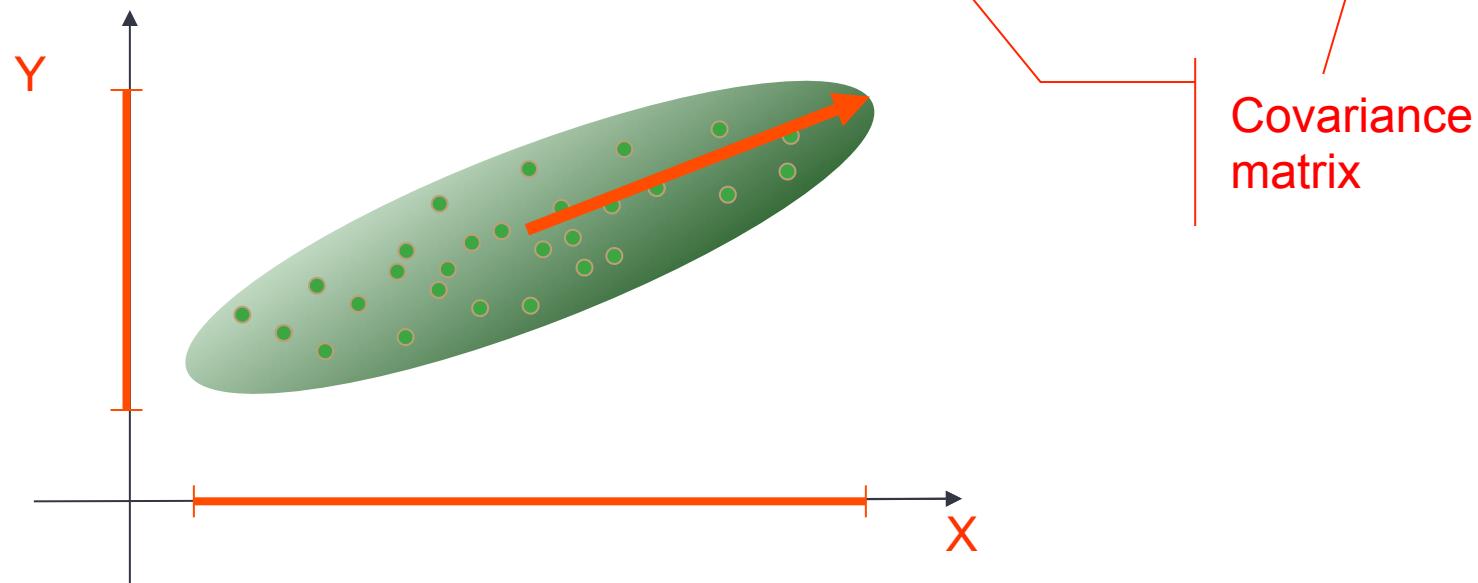
$$\Delta_{Euc}(\vec{a}, \vec{b}) = \sqrt{(\vec{a} - \vec{b})^T (\vec{a} - \vec{b})}$$



$$S[x, y] = Cov(x, y) = E((x - \mu_x)(y - \mu_y))$$

Mahalonobis Distance

$$\Delta_{Mah}(\vec{a}, \vec{b}) = \sqrt{(\vec{a} - \vec{b})^T S^{-1} (\vec{a} - \vec{b})}$$

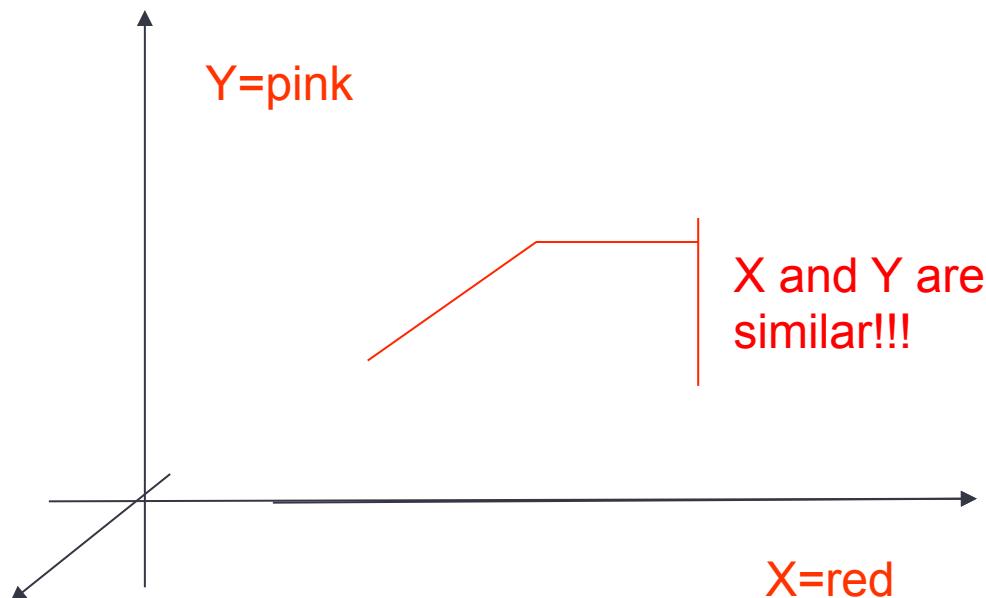


Covariance matrix

$$S[i, j] = Cov(i, j) = E((i - \mu_i)(j - \mu_j))$$

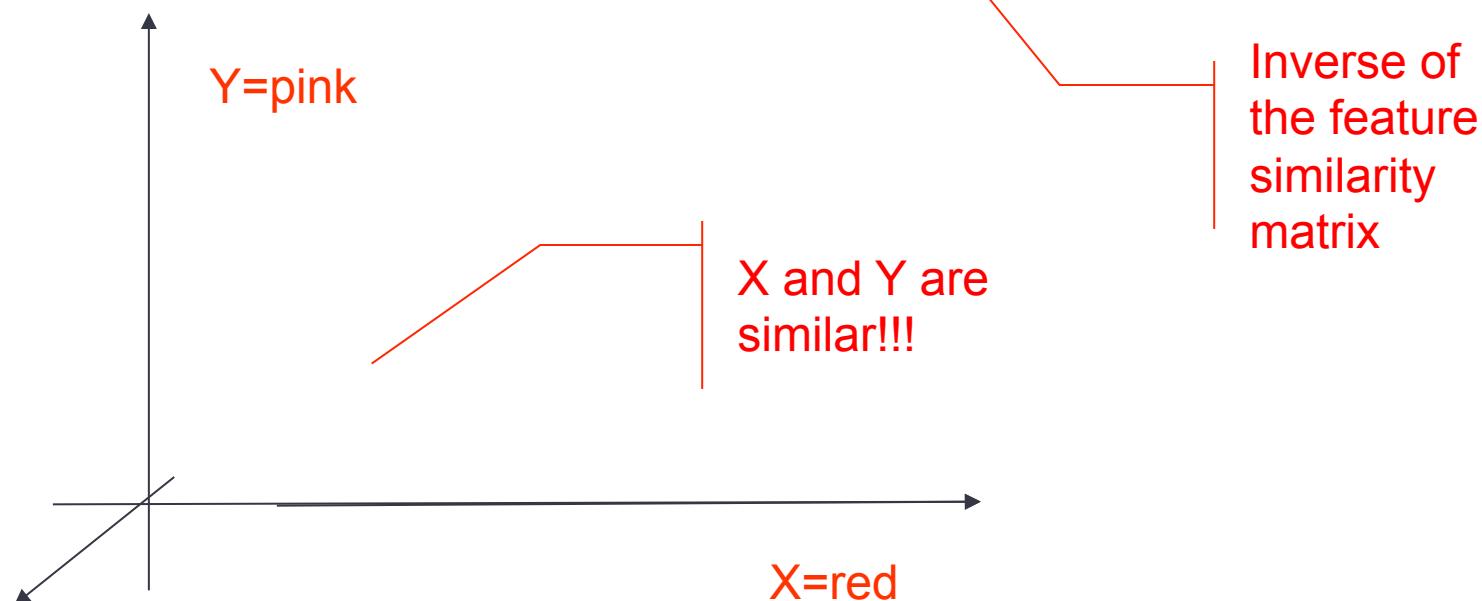
Quadratic Distance

$$\Delta_{Euc}(\vec{a}, \vec{b}) = \sqrt{(\vec{a} - \vec{b})^T (\vec{a} - \vec{b})}$$

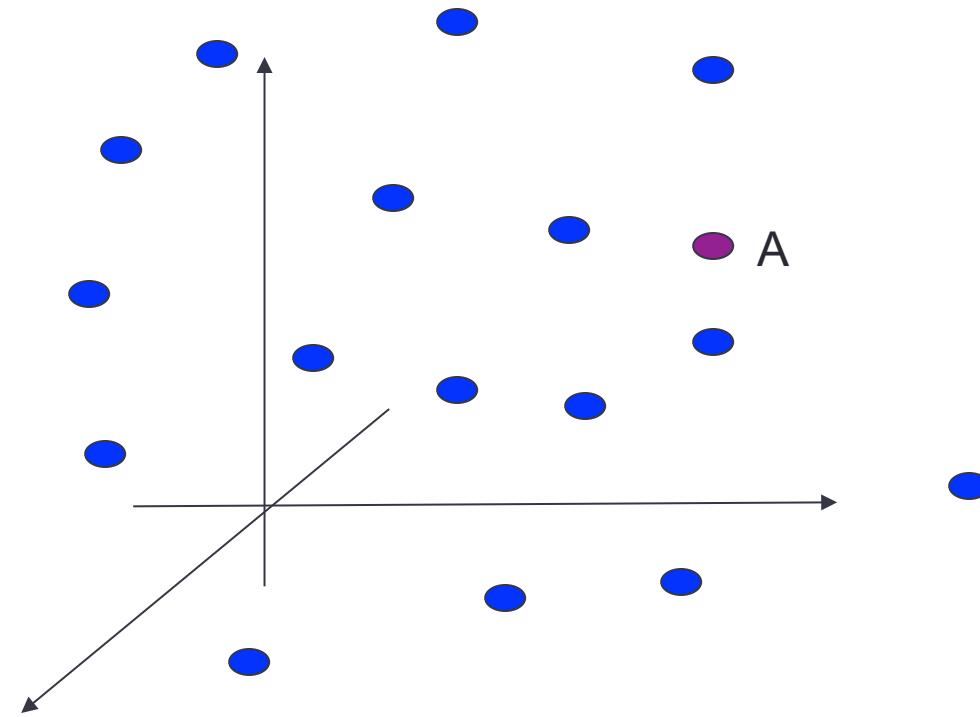


Quadratic Distance

$$\Delta_{quad}(\vec{a}, \vec{b}) = \sqrt{(\vec{a} - \vec{b})^T A^{-1} (\vec{a} - \vec{b})}$$



Are there other similarity measures?



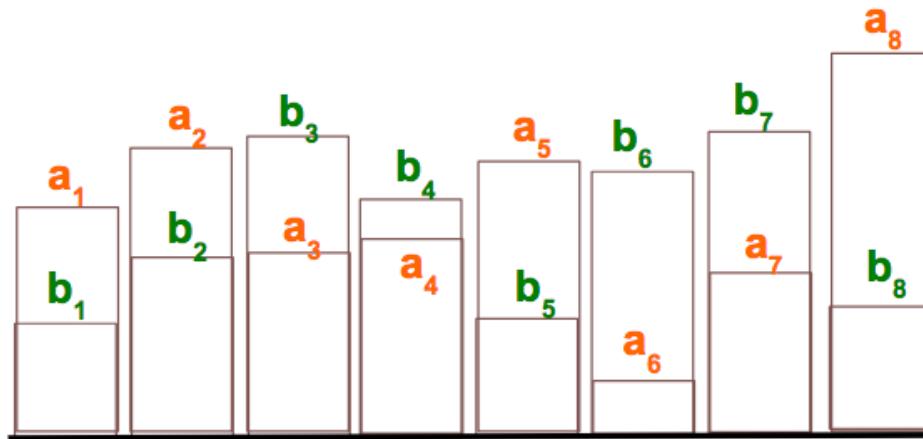
Intersection similarity

$$\text{sim}_{\text{int}}(\vec{a}, \vec{b}) = \frac{\sum_{i=1..n} \min(a_i.b_i)}{\sum_{i=1..n} \max(a_i.b_i)}$$

Intersection Similarity

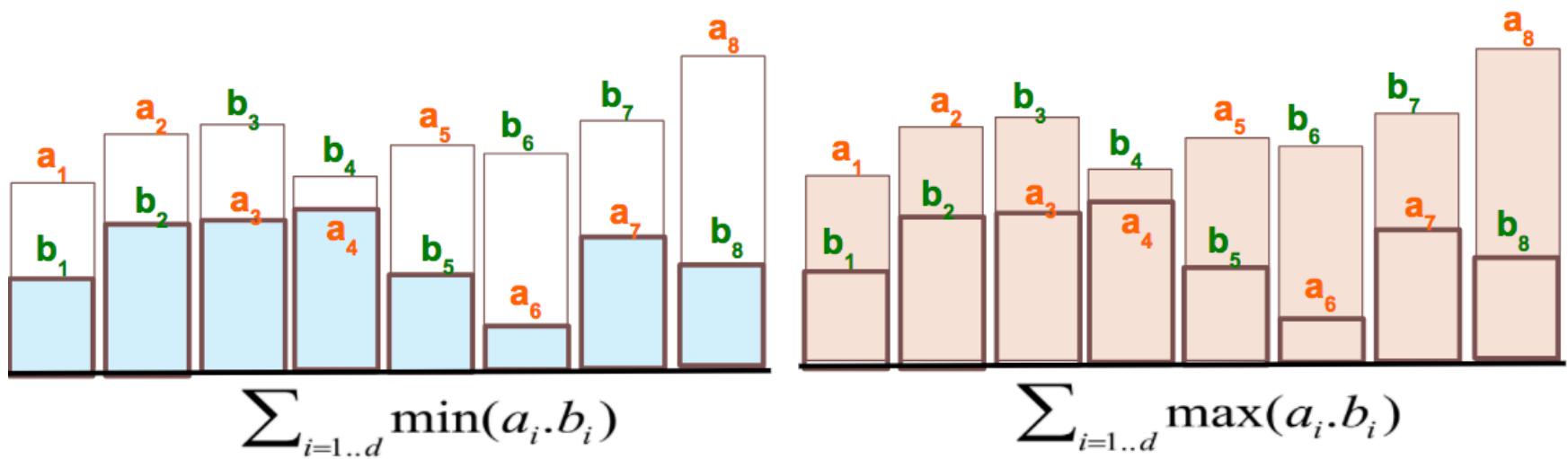
| Consider two vectors

$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$

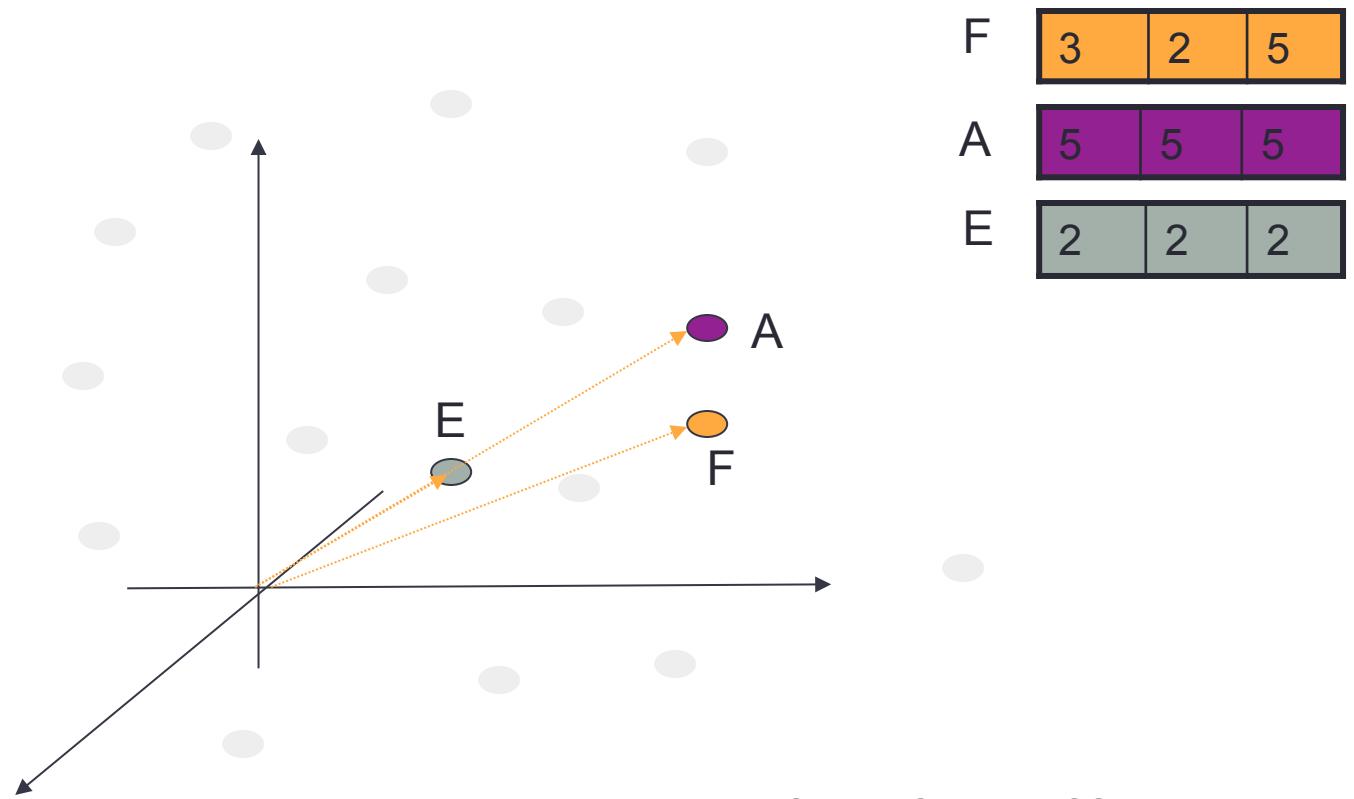


Intersection Similarity

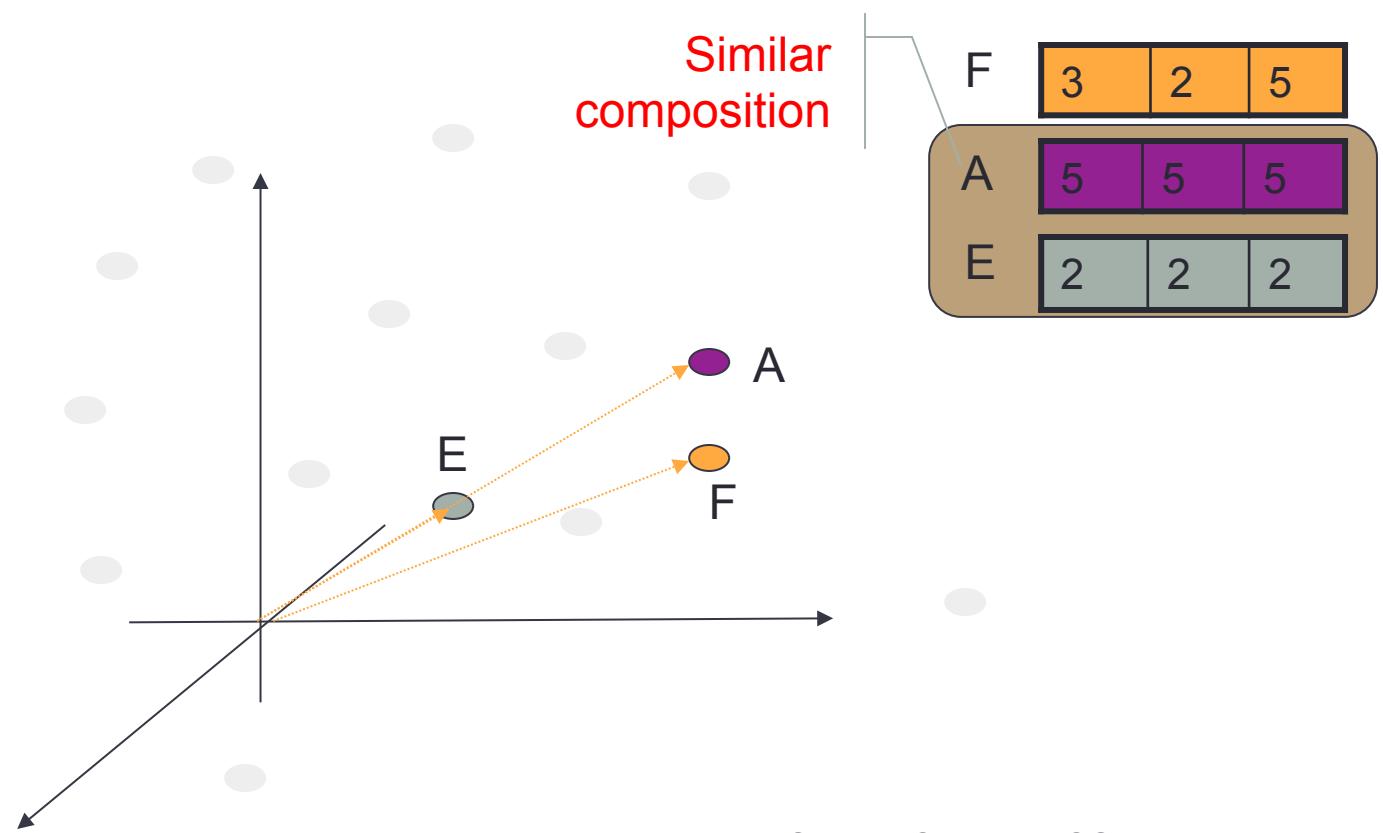
$$\text{sim}_{\text{int}}(\vec{a}, \vec{b}) = \frac{\sum_{i=1..d} \min(a_i.b_i)}{\sum_{i=1..d} \max(a_i.b_i)}$$



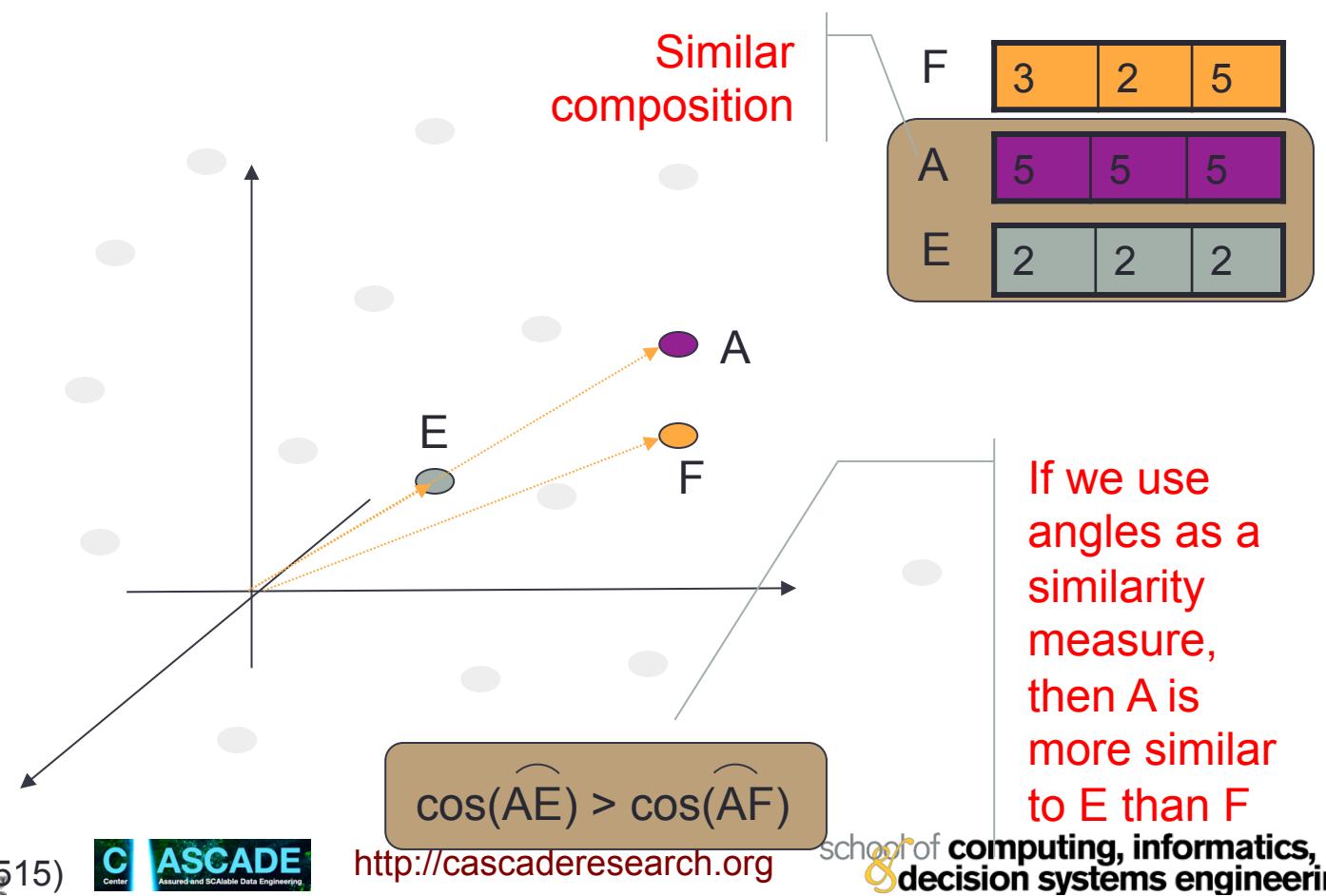
Let's try angles...



Let's try angles...



Let's try angles...



Angle-based measures

- Given

- Dot product $\vec{a} = \langle a_1, a_2, \dots, a_n \rangle$ $\vec{b} = \langle b_1, b_2, \dots, b_n \rangle$

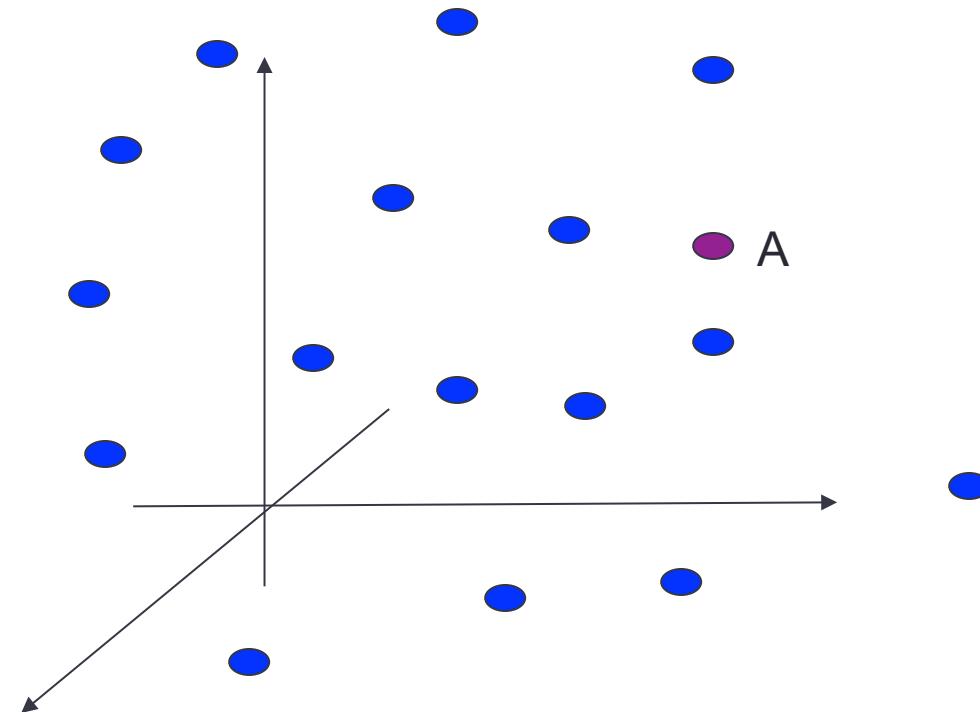
$$sim_{dot}(\vec{a}, \vec{b}) = \vec{a} \cdot \vec{b} = \sum_{i=1}^n a_i b_i$$

- Cosine similarity

$$sim_{cos}(\vec{a}, \vec{b}) = \cos(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|}$$

Cosine and dot product are the same if the vectors are unit length

Are there other similarity measures?



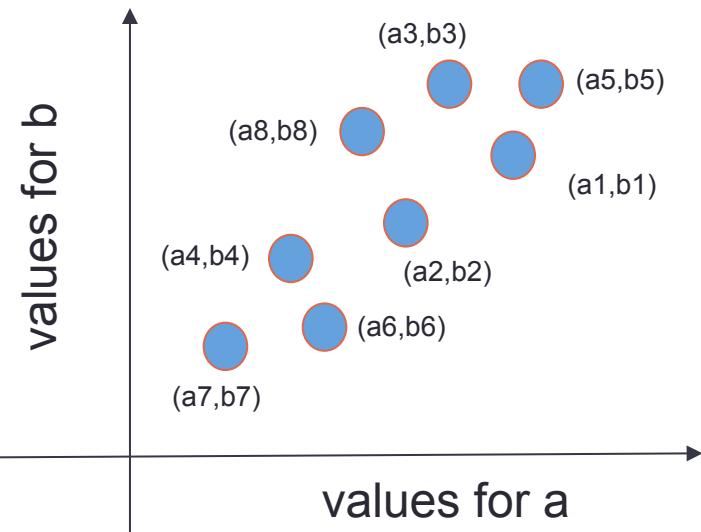
Other Commonly used Similarity / Distance Measures

- Pearson's correlation (similarity measure)
 - linear correlation (the strength of linear association) among the corresponding components of two vectors
- KL-Divergence (distance measure)
 - how one vector (interpreted as a probability distribution) diverges from the other
- Earth-movers distance (distance measure)
 - how one vector (interpreted as a probability distribution) diverges from the other
- Signal-to-noise ratio (similarity measure)
 - how one signal (represented as a vector) approximates another signal

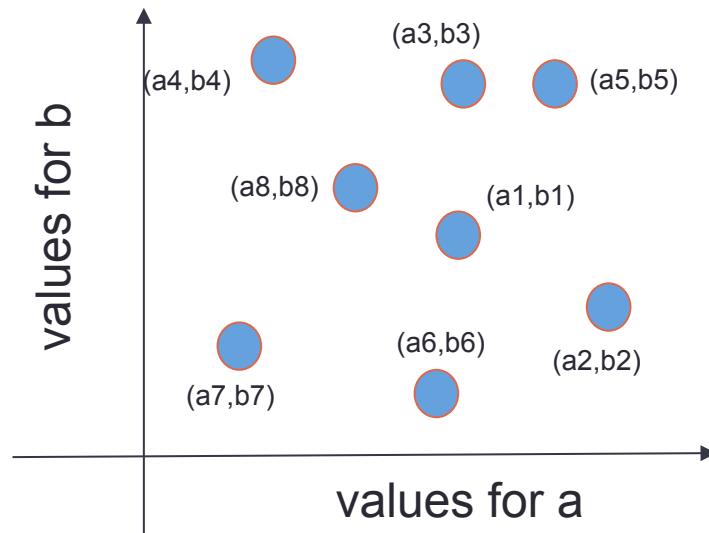
Pearson Correlation

| Consider two object vectors

$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$



Two positively correlated objects



Two un-correlated objects

KL-distance

(Kullback-Leibler divergence)

$$\Delta_{KL}(\vec{a}, \vec{b}) = \sum_{i=1..n} a_i \log \frac{a_i}{b_i}$$

$$\sum_{i=1}^n a_i = \sum_{i=1}^n b_i = 1$$

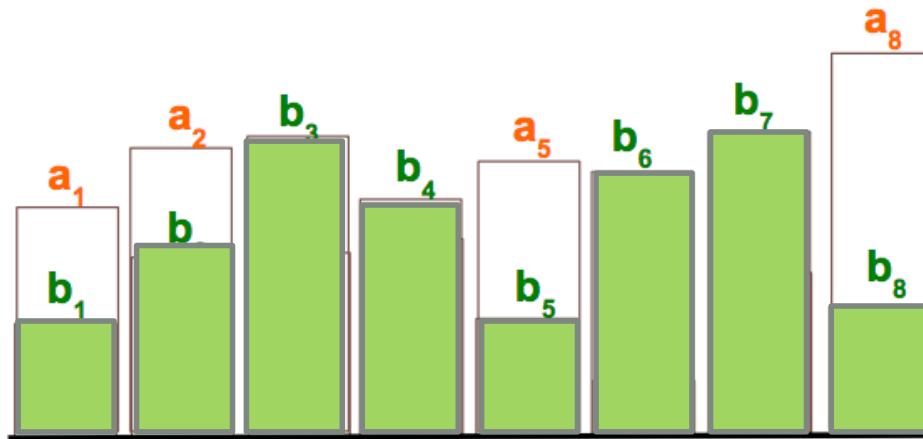
Relative entropy between
two probability
distributions

Earth Mover's Distance (EMD)

(Wasserstein metric)

| Consider two probability distributions

$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$



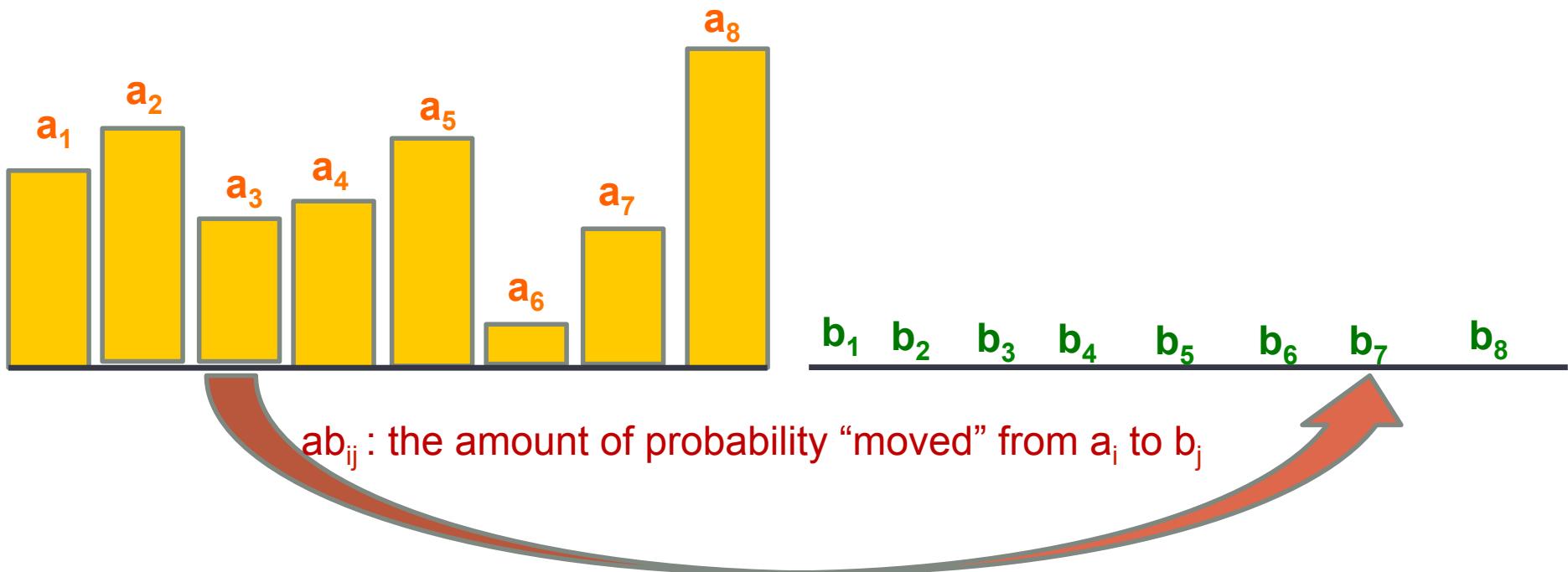
How can we convert vector a to vector b by redistributing the weights of vector a?

Earth Mover's Distance (EMD)

(Wasserstein metric)

| Consider two probability distributions

$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$

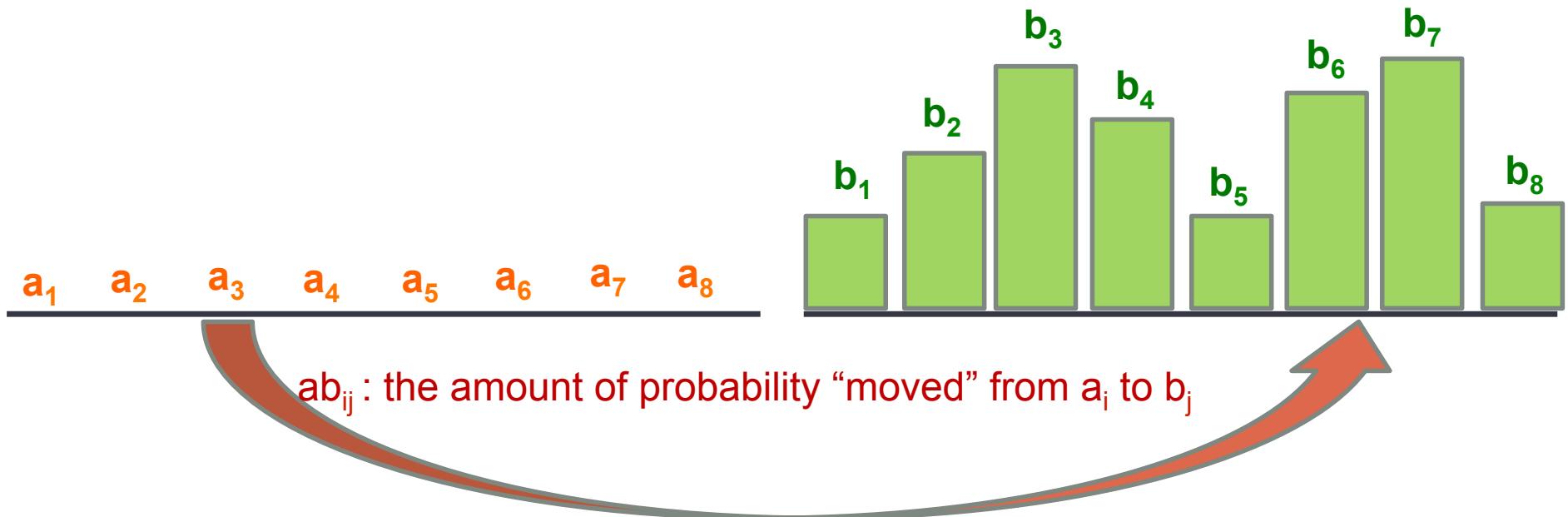


Earth Mover's Distance (EMD)

(Wasserstein metric)

| Consider two probability distributions

$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$

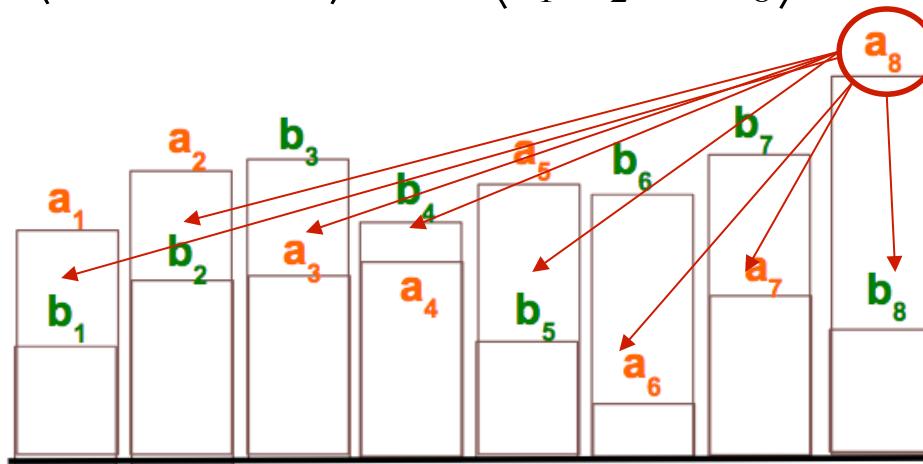


Earth Mover's Distance (EMD)

(Wasserstein metric)

| Consider two probability distributions

$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$



ab_{ij} : the amount of “probability” transferred from a_i to b_j

Earth Mover's Distance (EMD)

(Wasserstein metric)

| Consider two probability distributions

$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$

| Minimize $\text{cost}_{ij} = \sum_{i=1..n} \sum_{j=1..n} c_{ij} \times ab_{ij}$

| Such that

$$ab_{ij} \geq 0$$

$$\sum_{j=1..n} ab_{ij} = a_i$$

$$\sum_{i=1..n} ab_{ij} = b_j$$

$$\sum_{i=1..n} \sum_{j=1..n} ab_{ij} = 1$$

Earth Mover's Distance (EMD)

(Wasserstein metric)

| Consider two probability distributions

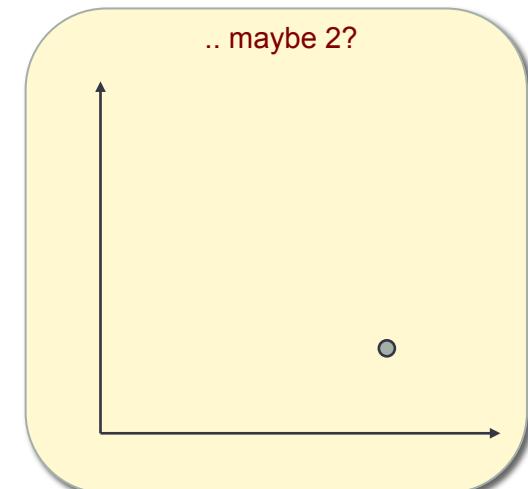
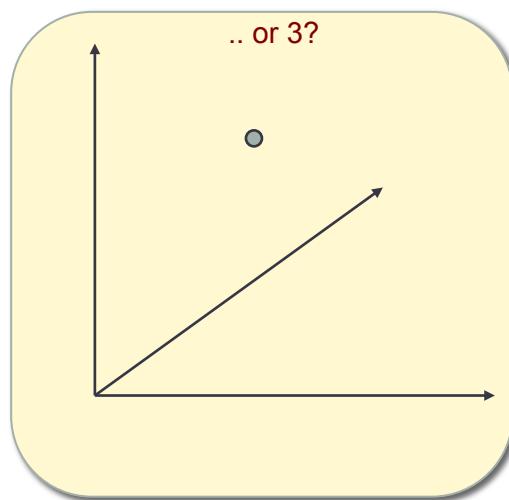
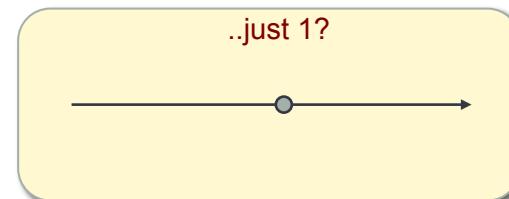
$$\vec{a} = \langle a_1, a_2, \dots, a_8 \rangle \quad \vec{b} = \langle b_1, b_2, \dots, b_8 \rangle$$

$$\Delta_{EMD}(\vec{a}, \vec{b}) = \frac{\sum_{i=1..n} \sum_{j=1..n} \text{cost}_{ij}}{\sum_{i=1..n} \sum_{j=1..n} ab_{ij}} = \sum_{i=1..n} \sum_{j=1..n} \text{cost}_{ij}$$

Signal-to-noise ratio (SNR)

$$sim_{snr}(\vec{a}, \vec{b}) = 20 \log_{10} \frac{\sqrt{\sum_{i=1..n} a_i^2}}{\sqrt{\sum_{i=1..n} (a_i - b_i)^2}}$$

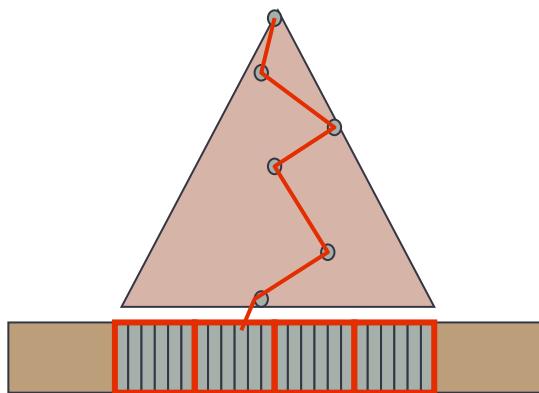
How many features do we need?



- ..or, many many more (100s, 1000s)???

We want as few features as possible

- **Dimensionality curse:** The more dimensions we have, the less efficient and effective search and analysis become
- **Efficiency:** Search data structures are not very efficient at high dimensions
- **Effectiveness:** The more dimensions we have, the more data we need to prevent overfitting

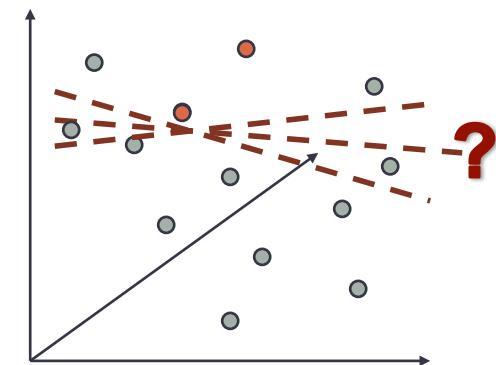
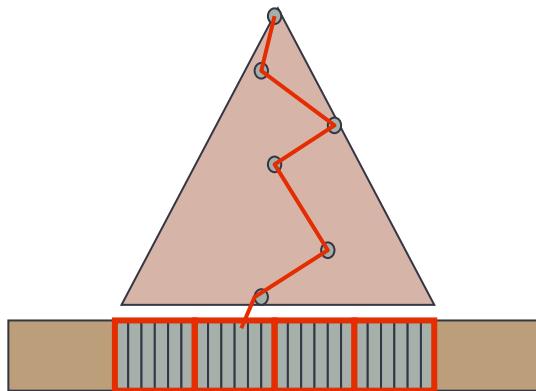


We want as few features as possible

- **Dimensionality curse:** The more dimensions we have, the less efficient and effective search and analysis become

- **Efficiency:** Search data structures are not very efficient at high dimensions

- **Effectiveness:** The more dimensions we have, the more data we need to prevent overfitting

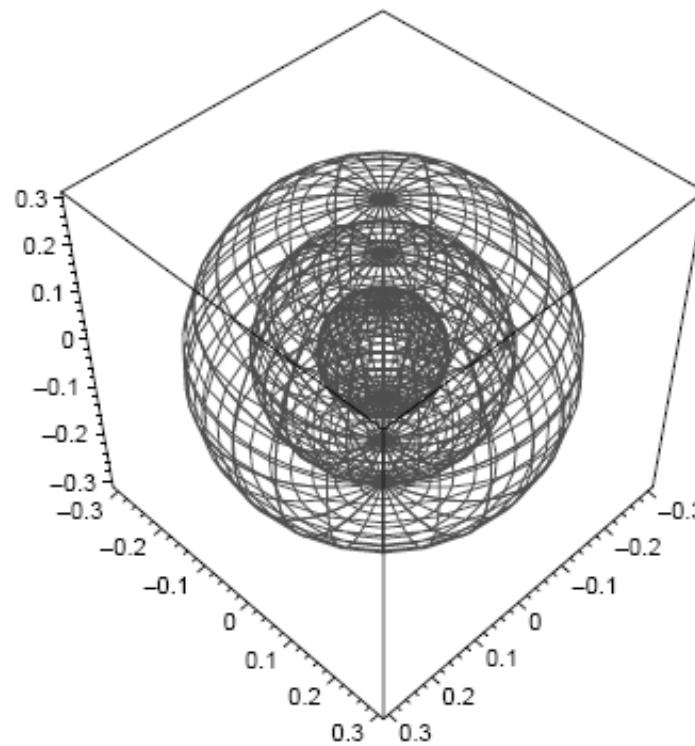


Feature selection

- more features mean more **storage space**, more feature extraction **time**, and higher cost of **index management**.
- some of the index structures require **exponential storage space** in terms of the features that are used for indexing
- more features means pairwise object **similarity/distance computations are more expensive**.
- searches in multi-dimensional vector spaces suffer from an inherent problem, called “**dimensionality curse**”

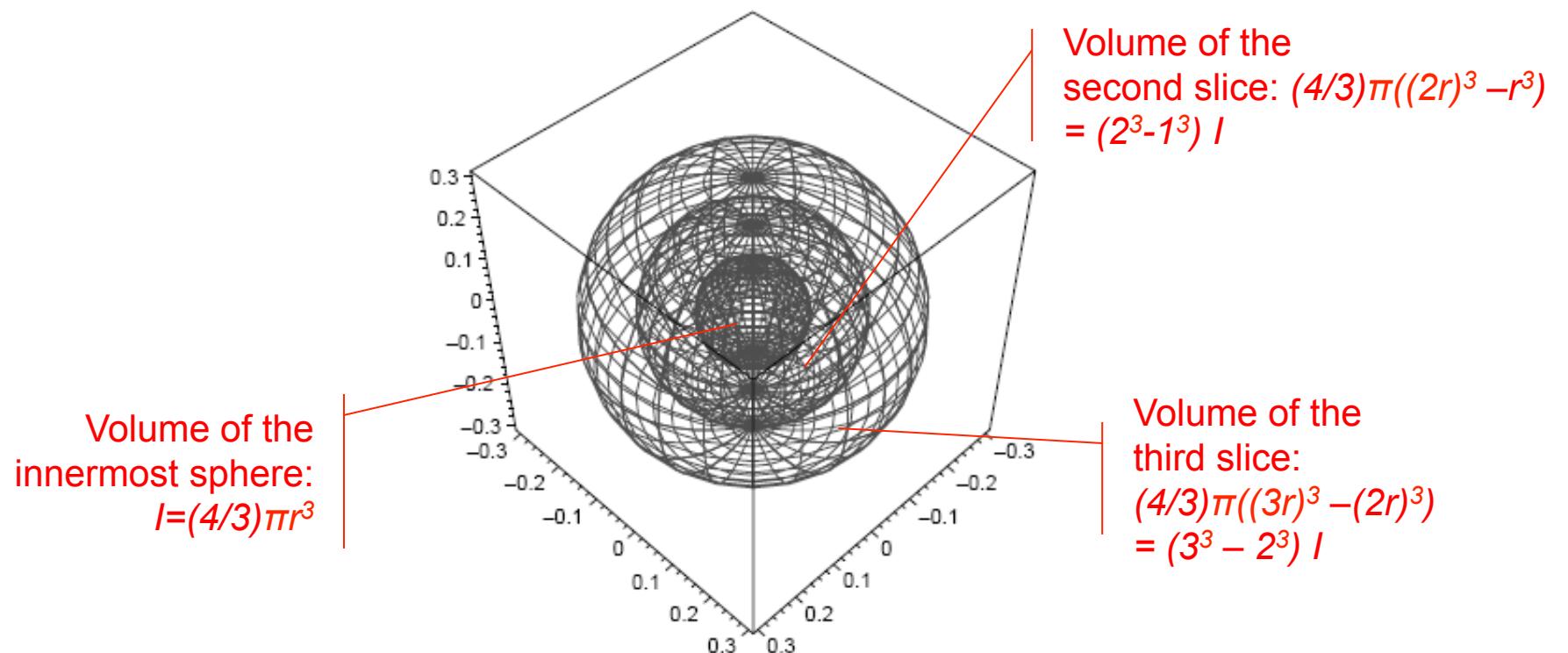
Dimensionality Curse

- Consider a query point and three alternative ranges:



Dimensionality Curse

- Consider a query point and three alternative ranges:



Dimensionality Curse

- In n-dimensional space, if the number of points in the inner most sphere is I , then
 - number of points in the second slice is $O(2^{n-1} I)$
 - number of points in the second slice is $O(3^{n-1} I)$
 - number of points in the second slice is $O(4^{n-1} I)$
- This means that most of the points lie in the outermost slice!!!!

Power Law and Dimensionality Curse

- **Power law:** most real world data is such that given d -dimensional space, the number of pairs of elements within a given distance, r , follows the formula

$$\text{pairs}(r) = c \times r^d$$

Dimensionality curse

- If

$$\lim_{n \rightarrow \infty} \frac{\text{var}(\Delta_n(\vec{q}, \vec{o}))}{\text{expected}(\Delta_n(\vec{q}, \vec{o}))} = 0$$

Distance between closest and furthest points converge

..then, nearest neighbor search does not make any sense.

What do we need to consider when selecting features?

- Application semantics
- Perception power
- Object description power
- Query description power
- Query workload
- Discrimination power

Perception-based feature selection

- E.g., lossy compression
 - Image may contain details that human eye can not recognize
 - Use a color model that reflects user perception
 - Compress those colors that are not easy to perceive
- Color table:
 1. reduce the number of colors to 256 (1 byte per pixel)
 2. Cluster similar colors into a single bucket and assign a single color to the bucket
 3. the set of buckets is called **color table**

Good feature..

- A good feature is **significant** and enables us to **differentiate** objects from others as much as possible
- A good feature corresponds to users' perception as much as possible
 - Relevance feedback!!!!

What does “significant” mean

- Information theoretic sense:
 - An event is more significant if it carries more information

What does “significant” mean

- Information theoretic sense:
 - An event is more significant if it carries more information
 - An event which has high occurrence rate carries less information
 - Solar eclipse is more interesting than sunset

High frequency ----- less information
Low frequency ----- high information

Entropy

- Total (or expected) information content
 - Uncertainty!!

$$H(X) \equiv \sum_{x \in \mathcal{A}_X} P(x) \log \frac{1}{P(x)}$$

Information content of
the event

Entropy (example)

- Total information content (uncertainty)

$$H(X) \equiv \sum_{x \in \mathcal{A}_X} P(x) \log \frac{1}{P(x)}$$

$p(A)$	$p(B)$	$-log_2 p(A)$	$-log_2 p(B)$	$\frac{p(A)}{log_2 p(A)}$	$\frac{p(B)}{log_2 p(B)}$	$H(E)$
0.05	0.95	1.3	0.022	0.216	0.07	0.29
0.5	0.5	1	1	0.5	0.5	1
0.95	0.05	0.022	1.3	0.07	0.216	0.29

Entropy (example)

- Total information content (uncertainty)

$$H(X) \equiv \sum_{x \in \mathcal{A}_X} P(x) \log \frac{1}{P(x)}$$

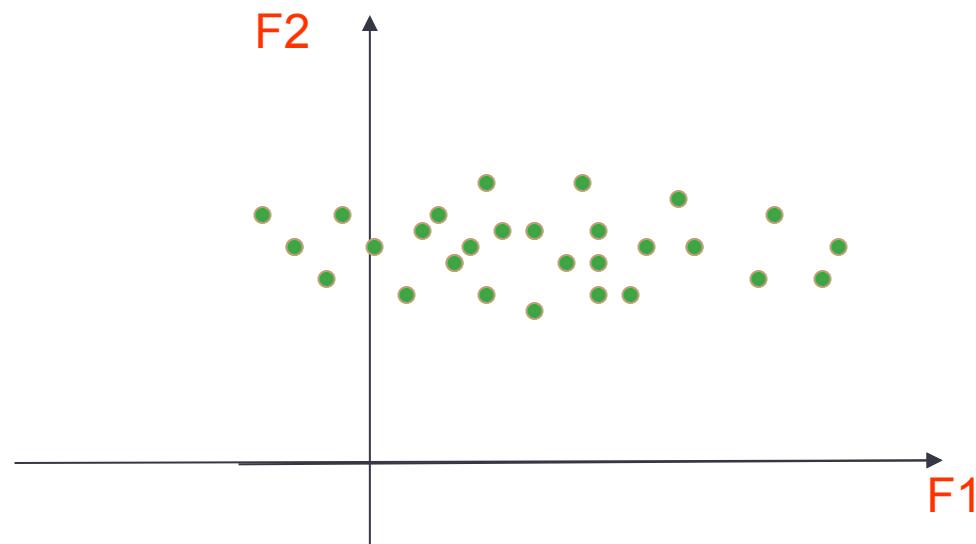
$P(a) = 0.5, P(b) = 0.5 \longrightarrow H = 1$ more uncertain
more information

$P(a) = 0.95, P(b) = 0.05 \longrightarrow H = 0.29$ less uncertain
less information

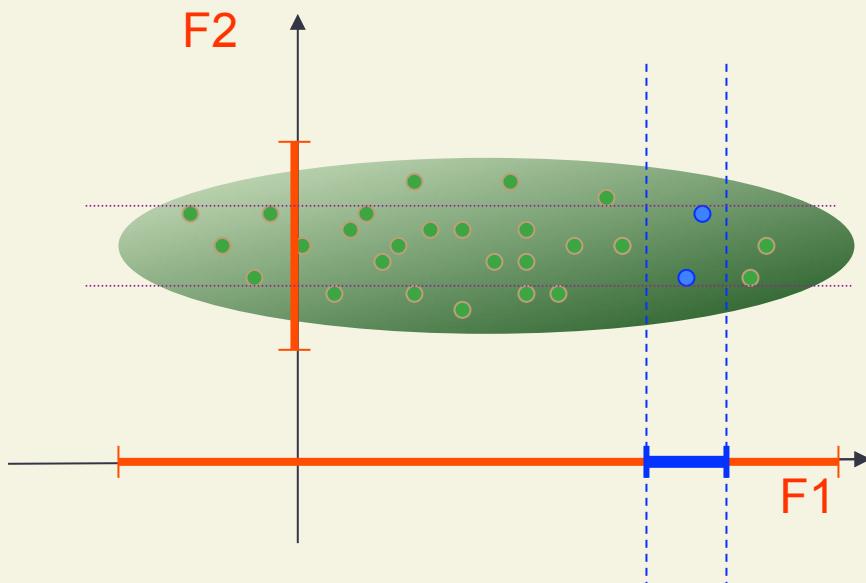
Conditional Entropy

- X takes values from $\{x_1, \dots, x_n\}$
 - $\sum_i p(X=x_i) = 1$
- Y takes values from $\{y_1, \dots, y_m\}$
 - $\sum_i p(Y=y_i) = 1$
- Conditional entropy of X given $Y=y_j$ is:
 - $H(X | Y=y_j) = \sum_i p(X=x_i | Y=y_j) * \lg 1/p(X=x_i | Y=y_j)$
- Conditional entropy of X given Y is:
 - $H(X | Y) = \sum_j p(Y=y_j) \sum_i p(X=x_i | Y=y_j) * \lg 1/p(X=x_i | Y=y_j)$

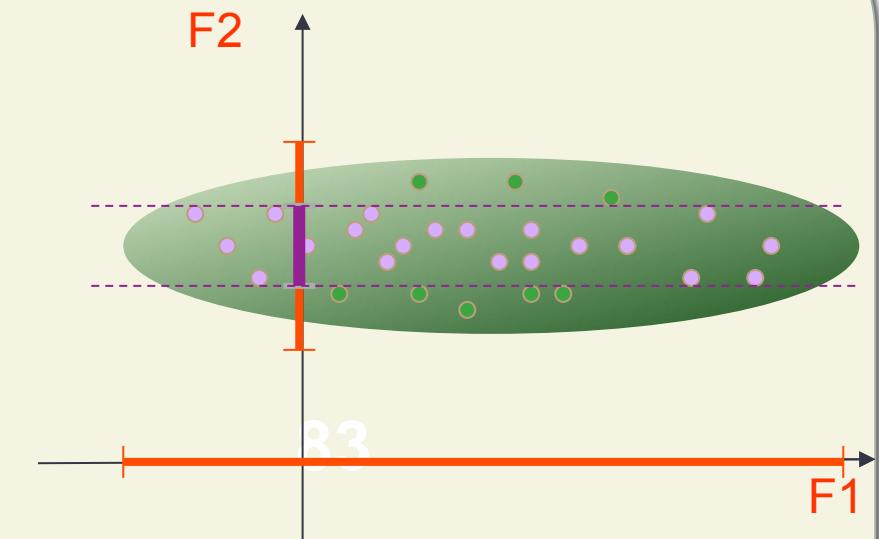
Feature selection...



Which feature is better?



F_1 : good discrimination



F_2 : poor discrimination

Text (as a collection of keywords).....

- Each document is represented as a multi-set of keywords
 - Content words (terms)
 - Non-content words (stop words)

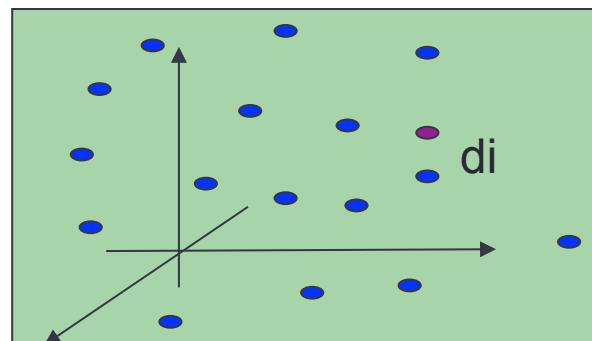
Vector representation

- Given a set of keywords, each document is represented as a vector:

$$d_i = \langle w_{i1}, w_{i2}, w_{i3}, \dots, w_{in} \rangle$$

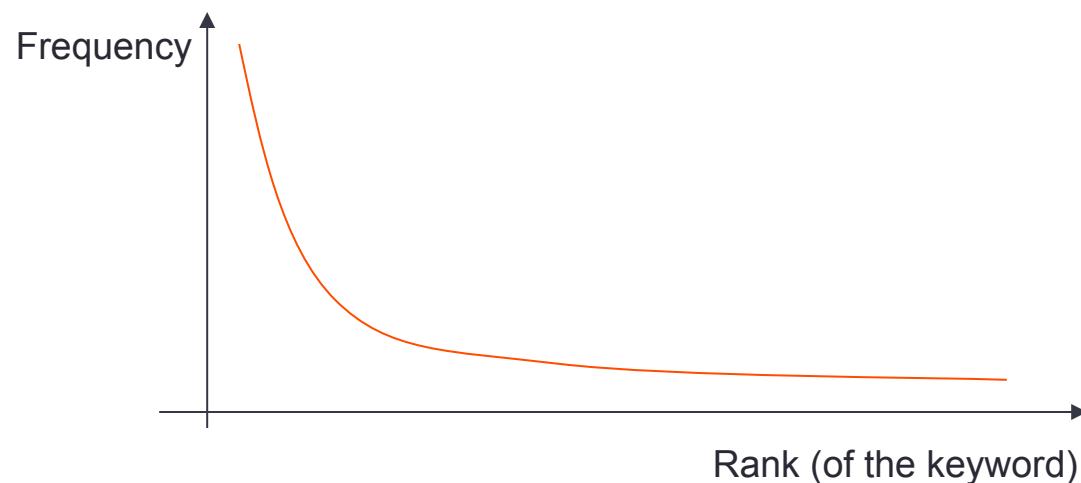
where

- $w_{ij} = 0$, if the keyword does not occur in d_i
- $w_{ij} > 0$, if the keyword occurs in d_i



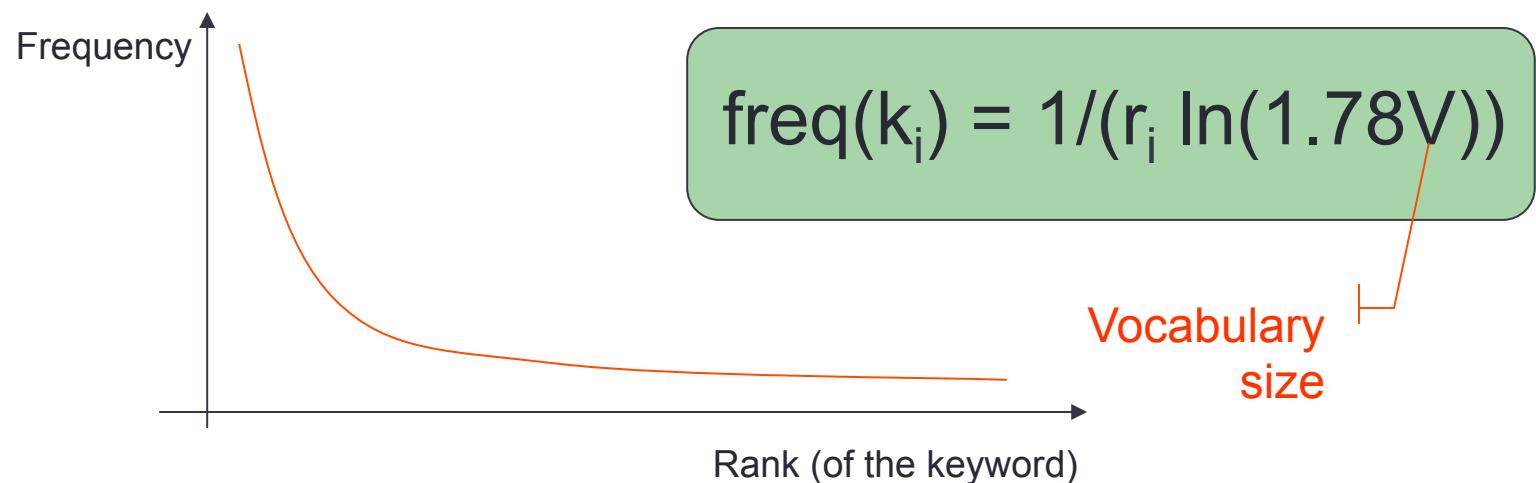
Zipfian Distribution

- The frequency of the k^{th} most frequent word in a collection is $(1/k)^{\Theta}$ times the most frequent word.



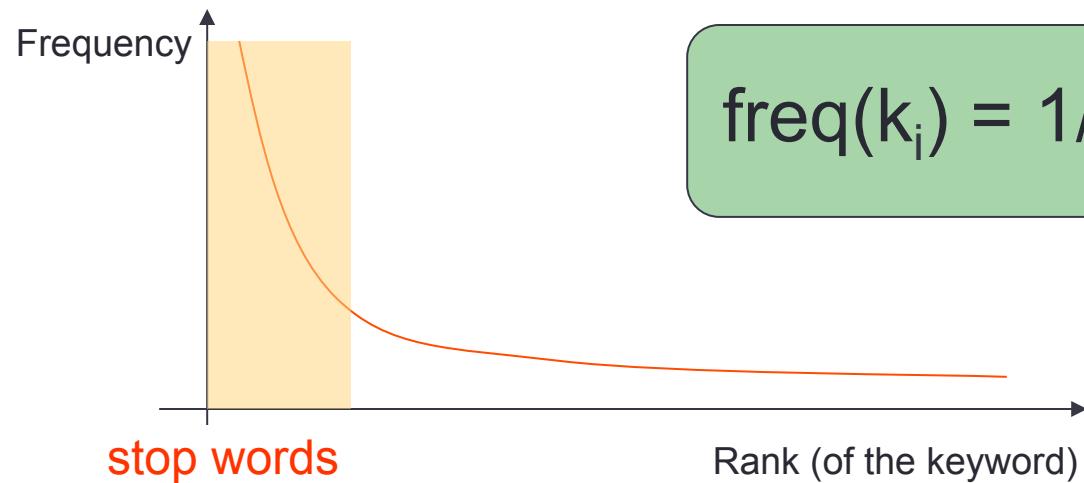
Zipfian Distribution

- The frequency of the k^{th} most frequent word in a collection is $(1/k)^{\ominus}$ times the most frequent word.



Zipfian Distribution

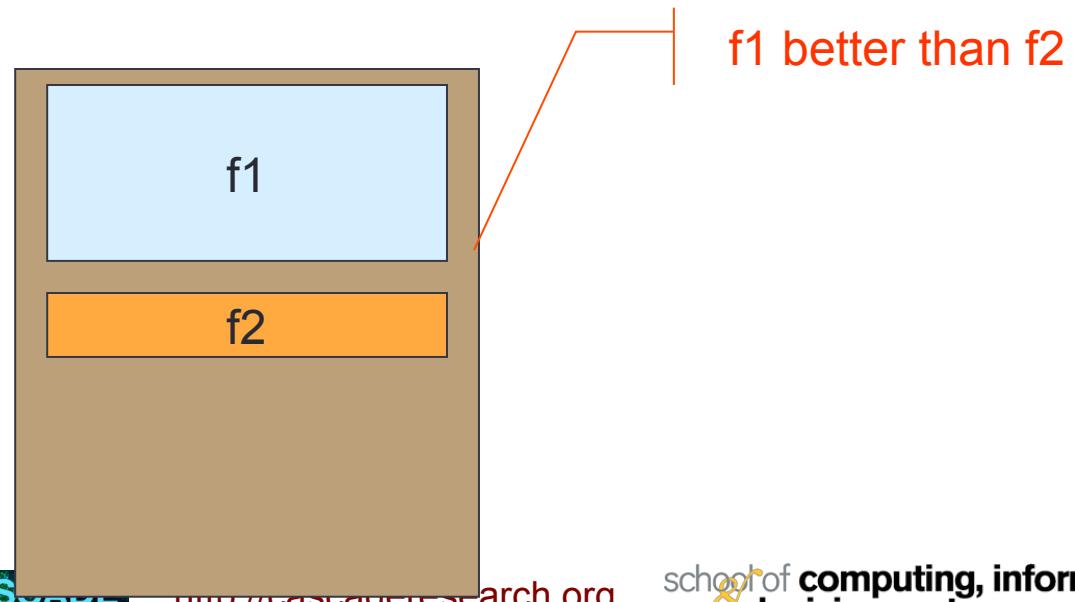
- The frequency of the k^{th} most frequent word in a collection is $(1/k)^{\Theta}$ times the most frequent word.



$$\text{freq}(k_i) = 1/(r_i \ln(1.78V))$$

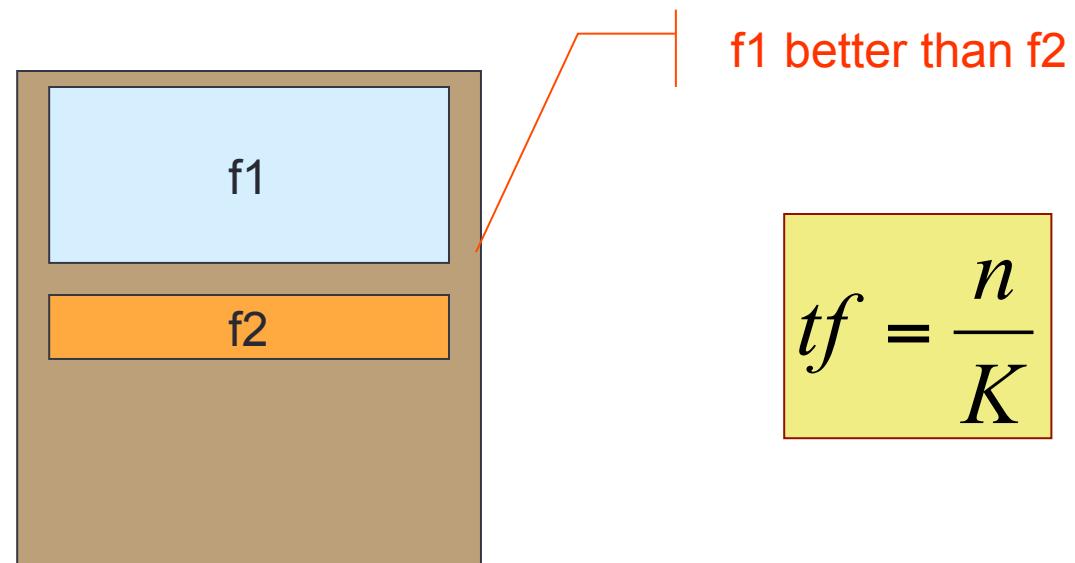
What are the weights????

- They need to capture how
 - good the term (feature) is in describing the content of the object



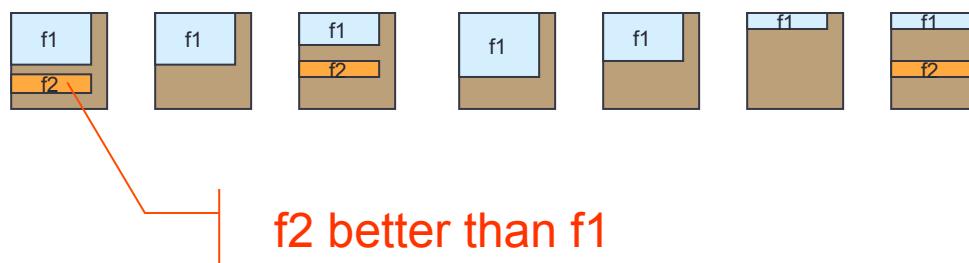
What are the weights????

- They need to capture how
 - good the term (feature) is in describing the content of the object



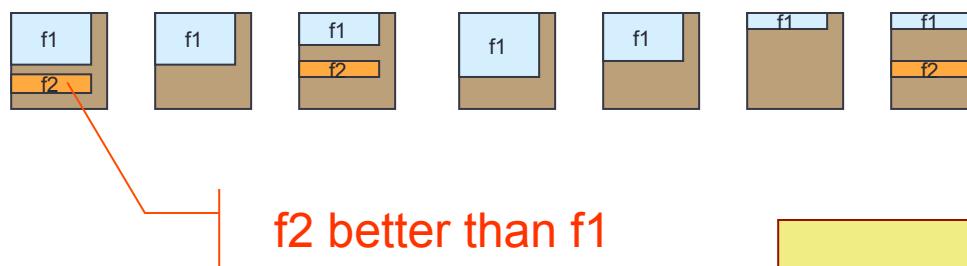
What are the weights????

- They need to capture how
 - differentiating the term (feature) is..



What are the weights????

- They need to capture how
 - differentiating the term (feature) is..



$$idf = \log\left(\frac{N}{m}\right)$$

What are the weights????

- They need to capture how
 - good the term (feature) is in describing the content of the object
 - differentiating the term (feature) is..

$$tfidf = \frac{n}{K} \log\left(\frac{N}{m}\right)$$

What are the weights????

- They need to capture how
 - good the term (feature) is in describing the content of the object
 - differentiating the term (feature) is..

$$norm_tfidf = \frac{n}{K} \frac{\log(\frac{N}{m})}{\max idf}$$

Idf of the keyword

Experiment results suggest that

- Poor terms have high document frequency
- Good terms have low document frequency
 - Problem: may not be queried often enough to be useful
- Best terms have medium document frequency

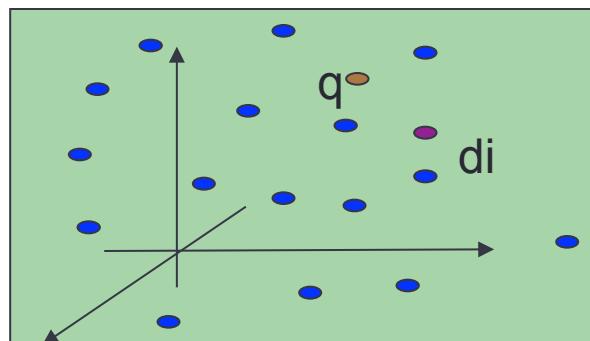
How about query terms??

- Given a set of keywords, each query is also represented as a vector:

$$q = \langle w_{q1}, w_{q2}, w_{q3}, \dots, w_{qn} \rangle$$

where

- $w_{qj} = 0$, if the keyword does not occur in d_q
- $w_{qj} > 0$, if the keyword occurs in d_q



How about query terms??

- They need to capture how
 - good the term (feature) is in describing the query
 - differentiating the term (feature) is..
 - Salton&Buckley suggests..



$$tfidf = \left(0.5 + 0.5 \frac{\frac{n}{K}}{\max freq} \right) \log\left(\frac{N}{m}\right)$$

Relevance based feature significance

- D : all documents
- R : the set of objects known to be relevant to the user
- I : the set of objects known to be irrelevant to the user

$$\log \left(\frac{p(f_k|R)(1 - p(f_k|I))}{p(f_k|I)(1 - p(f_k|R))} \right) \times (p(f_k|rel) - p(f_k|\overline{rel}))$$

assuming $I = D - R$:

$$\log \left(\frac{r_k(|D - R| - (d_k - r_k))}{(d_k - r_k)(|R| - r_k)} \right) \times \left| \frac{r_k}{|R|} - \frac{d_k - r_k}{|D - R|} \right|$$

r_k is the number of objects in R with f_k
 d_k is the number of objects in D with f_k

Entropy

- Total (or expected) information content
 - Uncertainty!!

$$H(X) \equiv \sum_{x \in \mathcal{A}_X} P(x) \log \frac{1}{P(x)}$$

Information content of
the event

Entropy (example)

- Total information content (uncertainty)

$$H(X) \equiv \sum_{x \in \mathcal{A}_X} P(x) \log \frac{1}{P(x)}$$

$p(A)$	$p(B)$	$-log_2 p(A)$	$-log_2 p(B)$	$\frac{p(A)}{log_2 p(A)}$	$\frac{p(B)}{log_2 p(B)}$	$H(E)$
0.05	0.95	1.3	0.022	0.216	0.07	0.29
0.5	0.5	1	1	0.5	0.5	1
0.95	0.05	0.022	1.3	0.07	0.216	0.29

Entropy (example)

- Total information content (uncertainty)

$$H(X) \equiv \sum_{x \in \mathcal{A}_X} P(x) \log \frac{1}{P(x)}$$

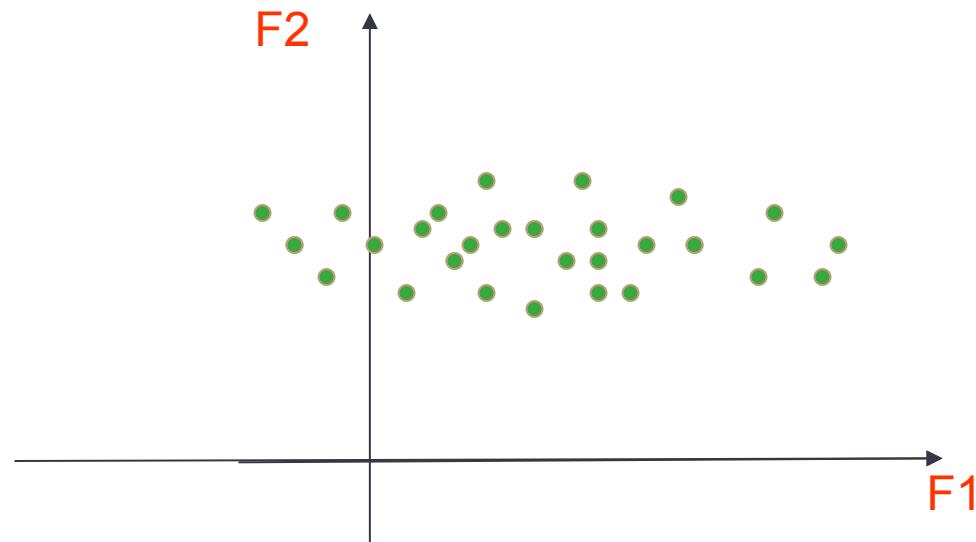
$P(a) = 0.5, P(b) = 0.5 \longrightarrow H = 1$ more uncertain
more information

$P(a) = 0.95, P(b) = 0.05 \longrightarrow H = 0.29$ less uncertain
less information

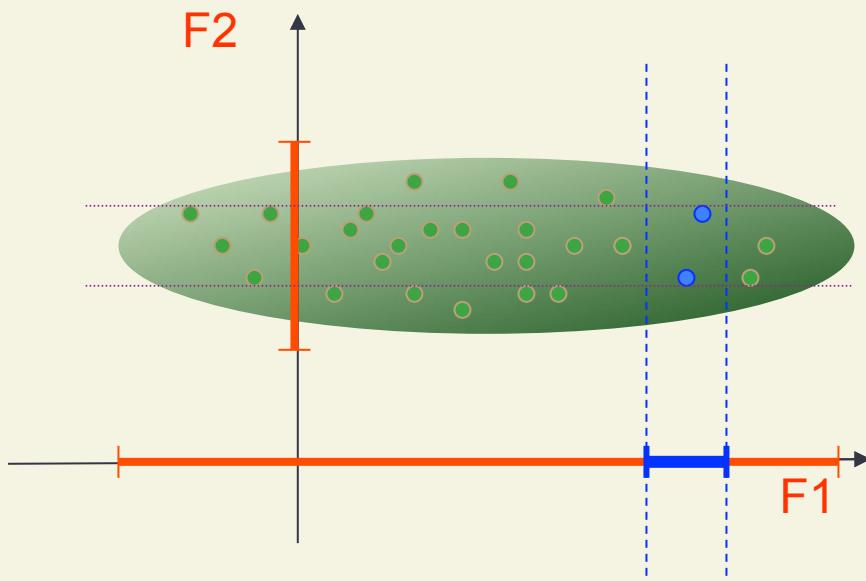
Conditional Entropy

- X takes values from $\{ x_1, \dots, x_n \}$
 - $\sum_i p(X=x_i) = 1$
- Y takes values from $\{ y_1, \dots, y_m \}$
 - $\sum_i p(Y=y_i) = 1$
- Conditional entropy of X given $Y=y_j$ is:
 - $H(X | Y=y_j) = \sum_i p(X=x_i | Y=y_j) * \lg 1/p(X=x_i | Y=y_j)$
- Conditional entropy of X given Y is:
 - $H(X | Y) = \sum_j p(Y=y_j) \sum_i p(X=x_i | Y=y_j) * \lg 1/p(X=x_i | Y=y_j)$

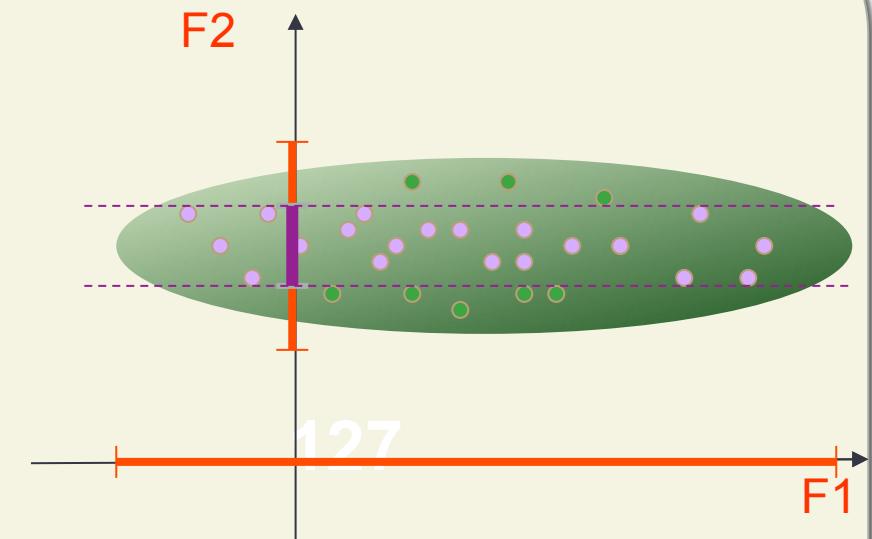
Feature selection...



Which feature is better?

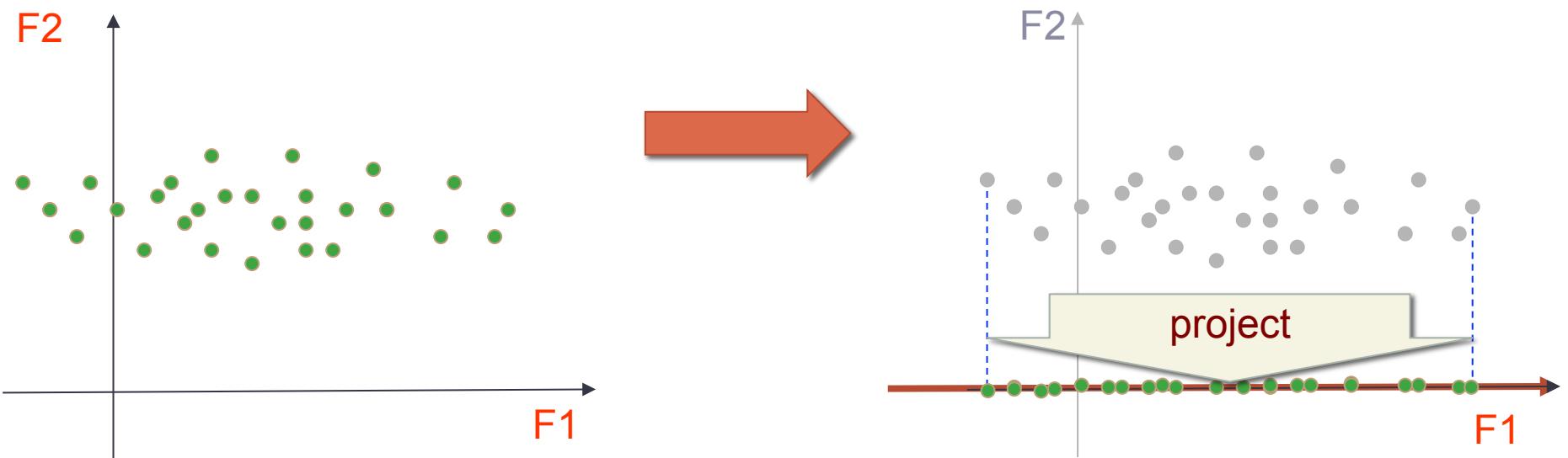


F1: good discrimination

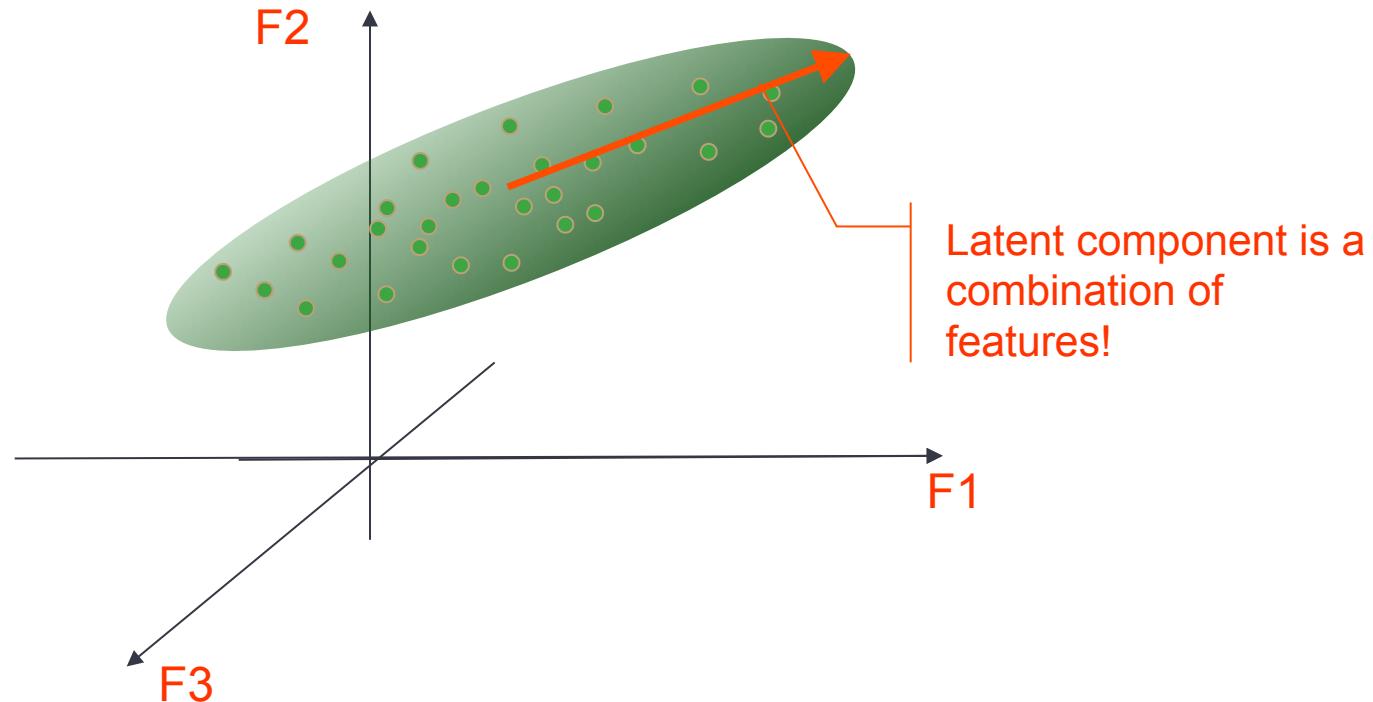


F2: poor discrimination

Feature selection...

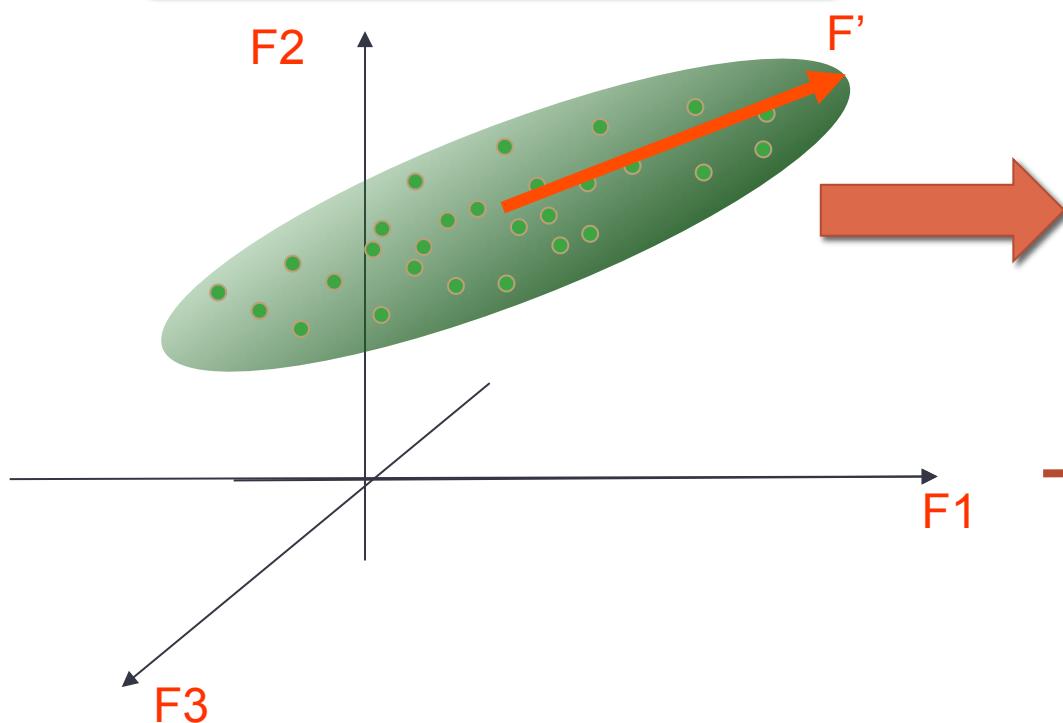


Latent component analysis

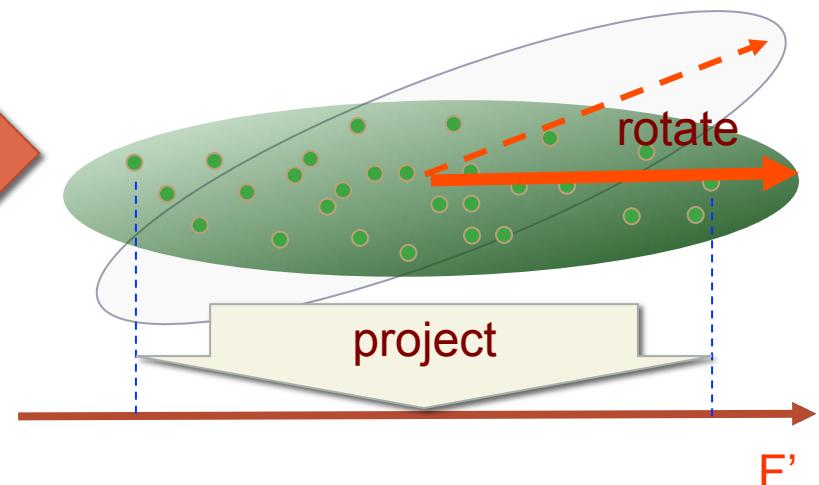


Latent feature...

Information is distributed across features



Information is concentrated on one feature



Through an appropriate transformation, the latent component can be used as a “discriminating” feature of the data

Search for latent features

- Given a database,

- $D = \{o_1, \dots, o_n\}$ of objects,
- originally described in terms of a feature set, $F = \{f_1, \dots, f_m\}$,

identify

- a set of “latent” features, $Z = \{z_1, \dots, z_k\}$,

such that

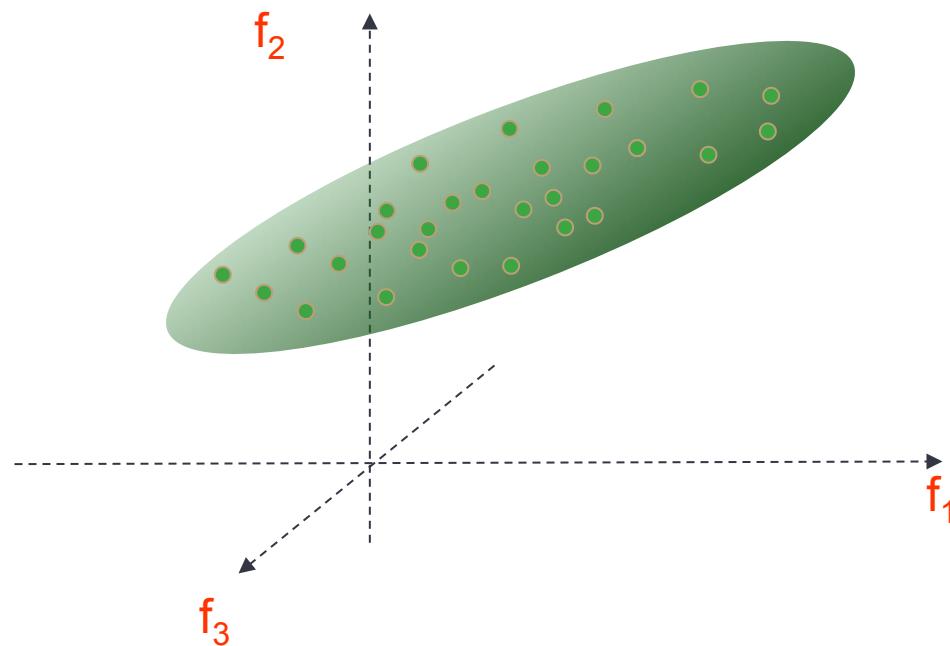
- $k << m$ and
- the transformed database $D' = \{o'_1, \dots, o'_n\}$ (described in the terms of latent features Z instead of F) is not very different from the original database D' i.e.,

$$D \sim D'$$

Search for latent features

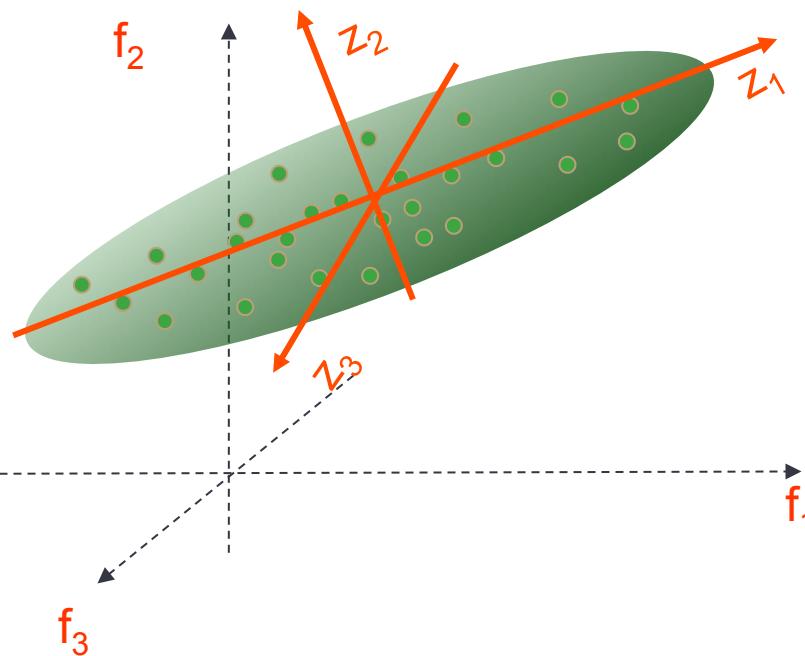
- Alternative #1: Feature selection
 - $Z \subseteq F$
- Alternative #2: Representative object selection
 - $Z \subseteq D$
- Alternative #3: Ortho-normal basis transformation
 - $\forall j \ z'_j = \Theta_j(f_1, \dots, f_m)$, where
 - $\forall j () \ \Theta_j()$ is a “normal” transformation
 - $\forall i,j () \ \Theta_i()$ and $\Theta_j()$ are mutually orthogonal
- Alternative #4: Arbitrary transformation

Ortho-normal basis transformations



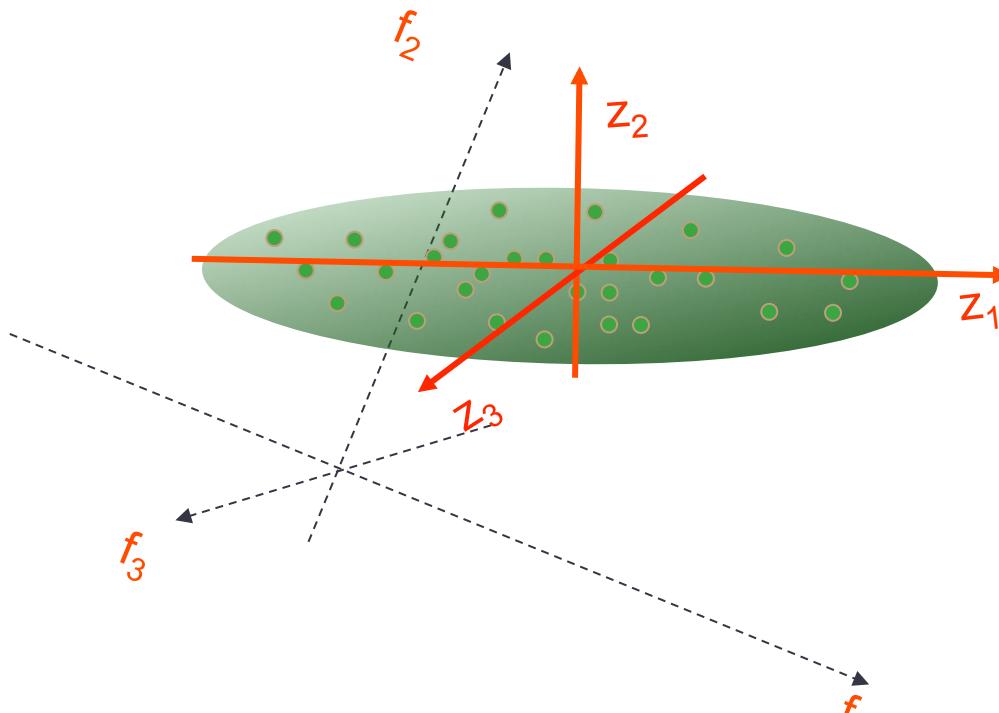
Let "D" denote distances among the objects

Ortho-normal basis transformations



Let "D" denote distances among the objects

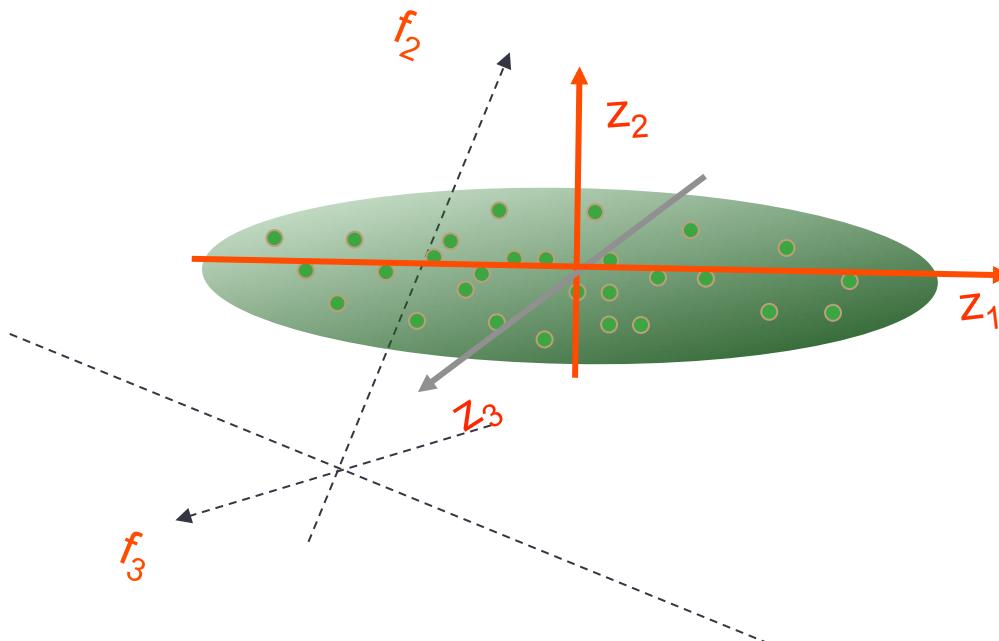
Ortho-normal basis transformations



Distances and angles among the objects are preserved

$$D' = D$$

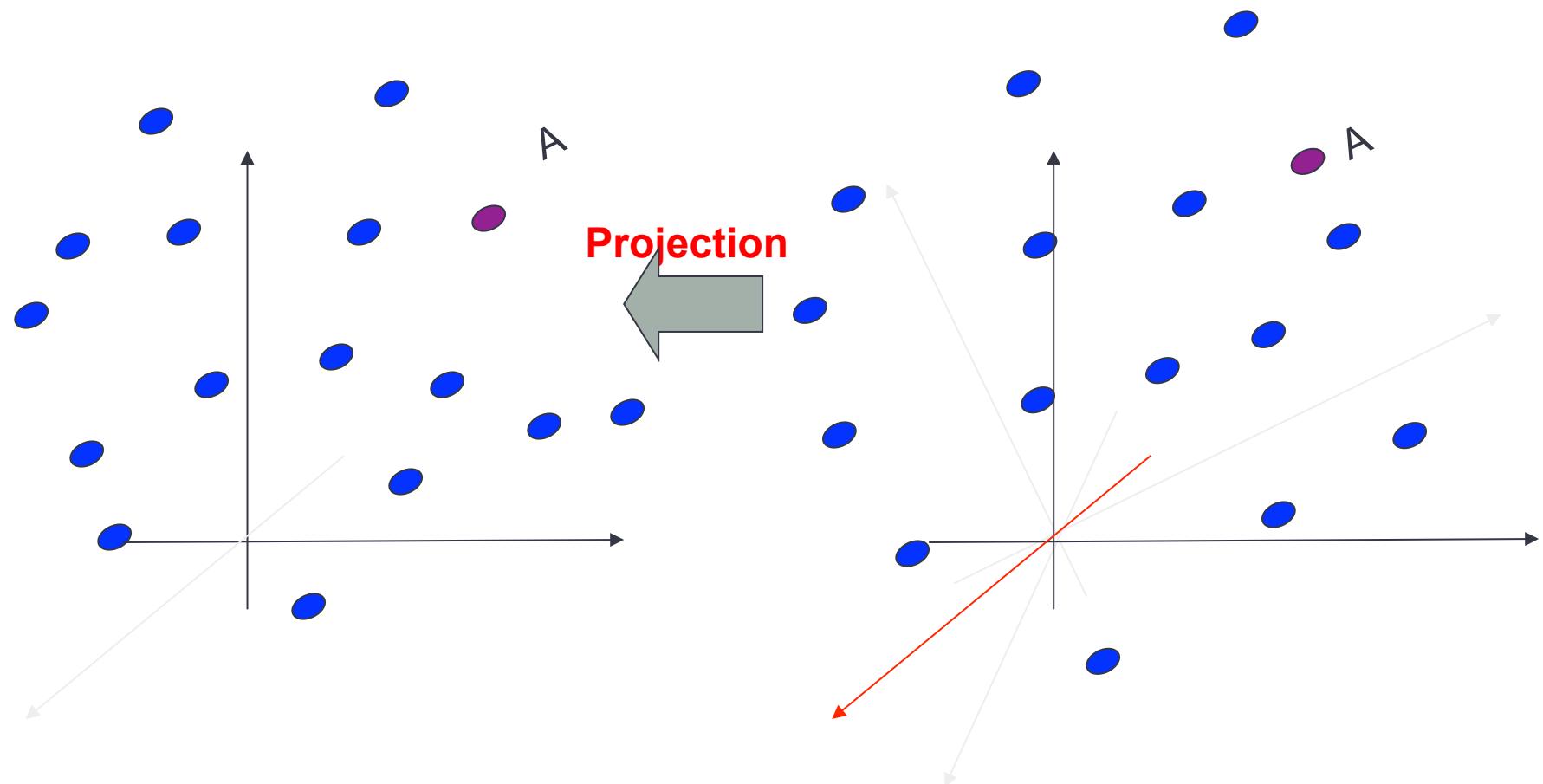
Ortho-normal basis transformations



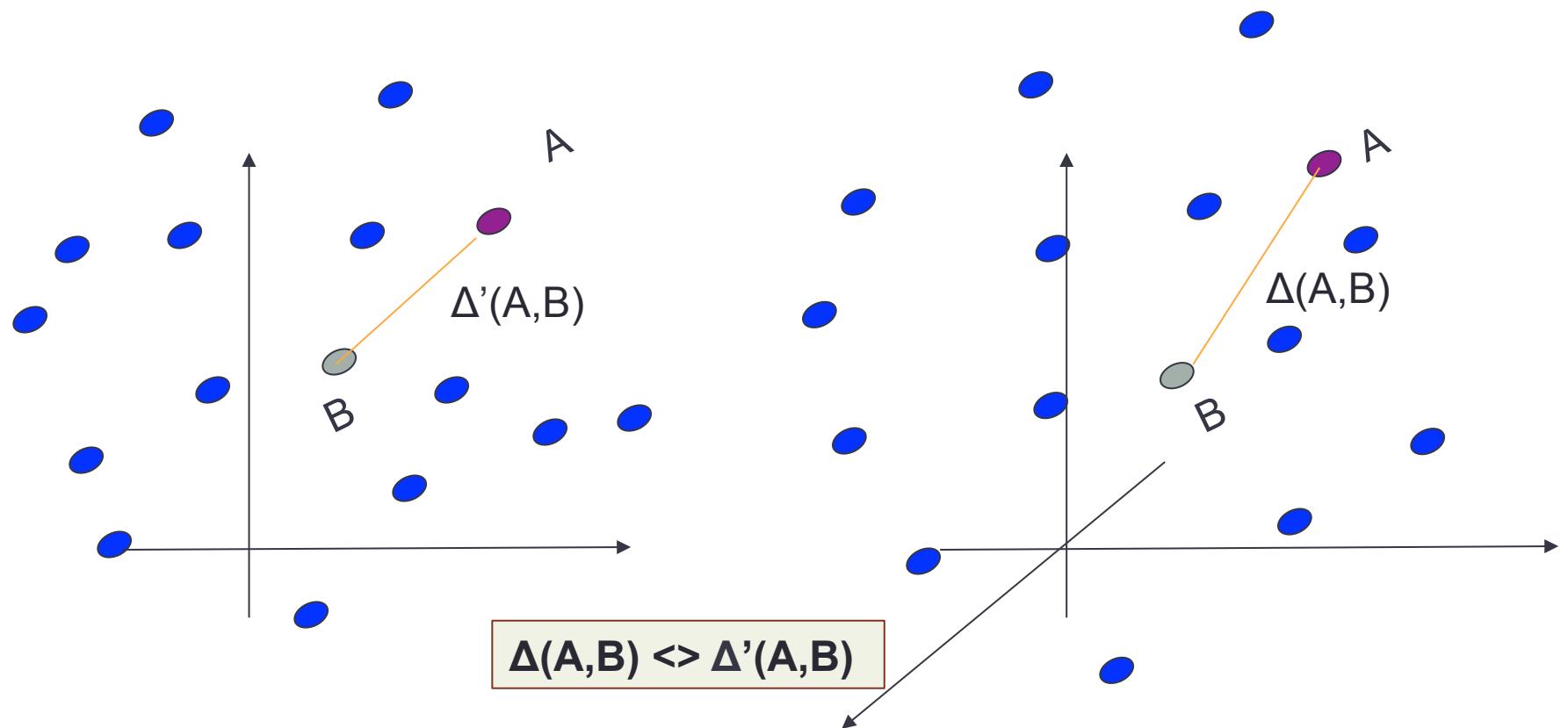
Distances and angles among the objects are “approximately” preserved

$$D' \sim D$$

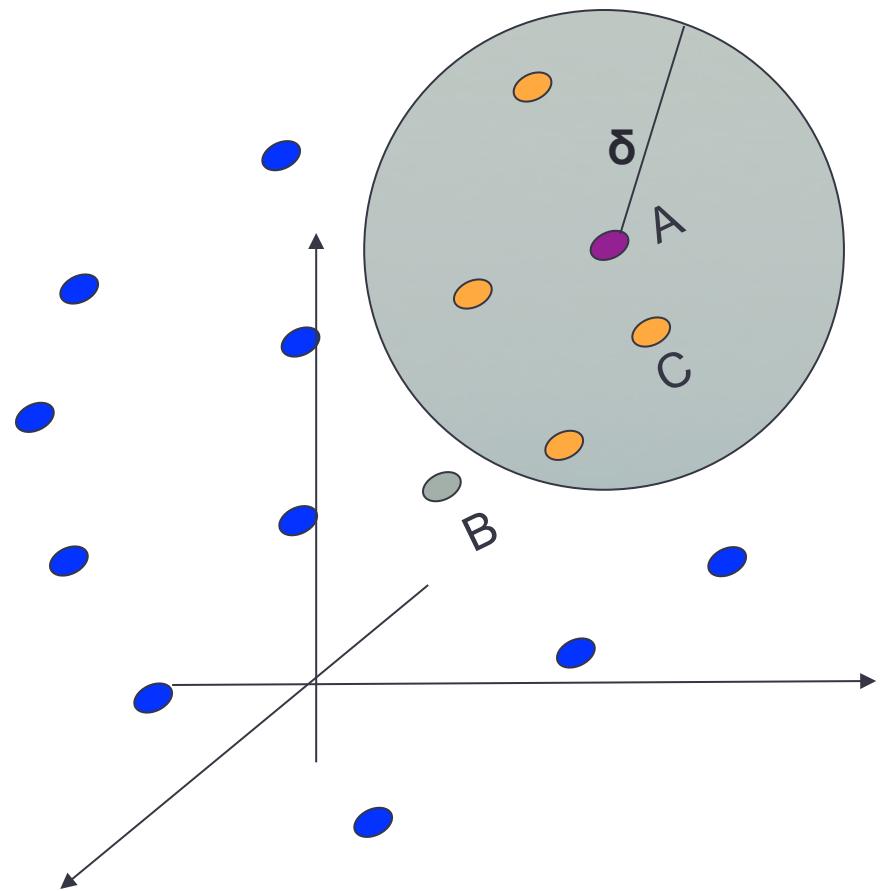
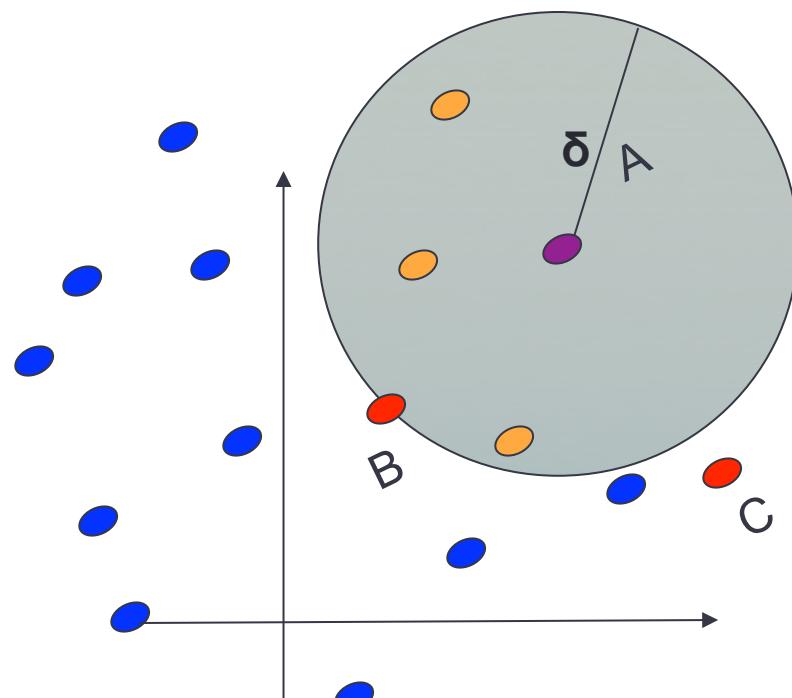
Transform + Projection (Compression or Feature selection)



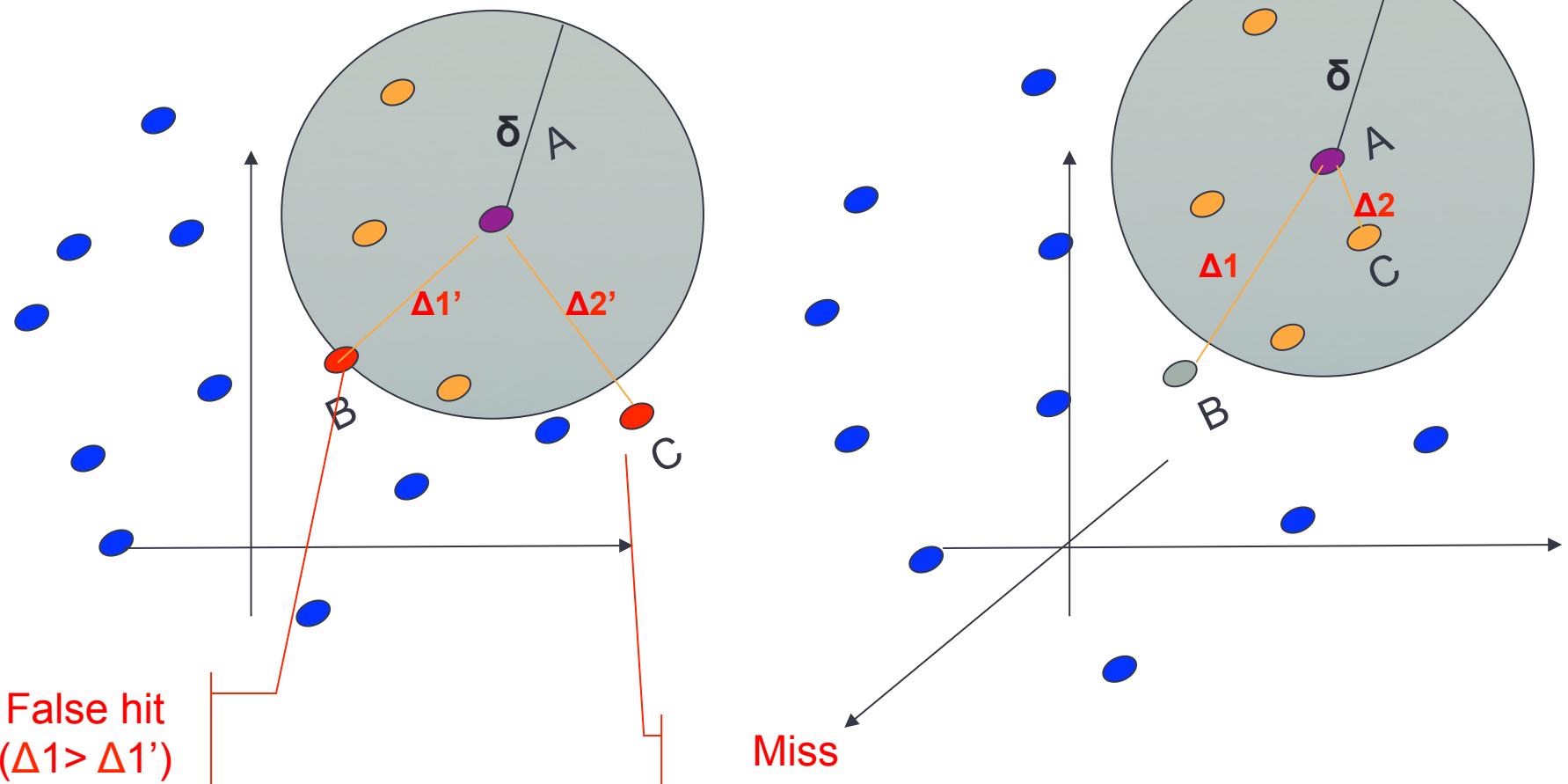
What happens to distances???



What happens to distances???

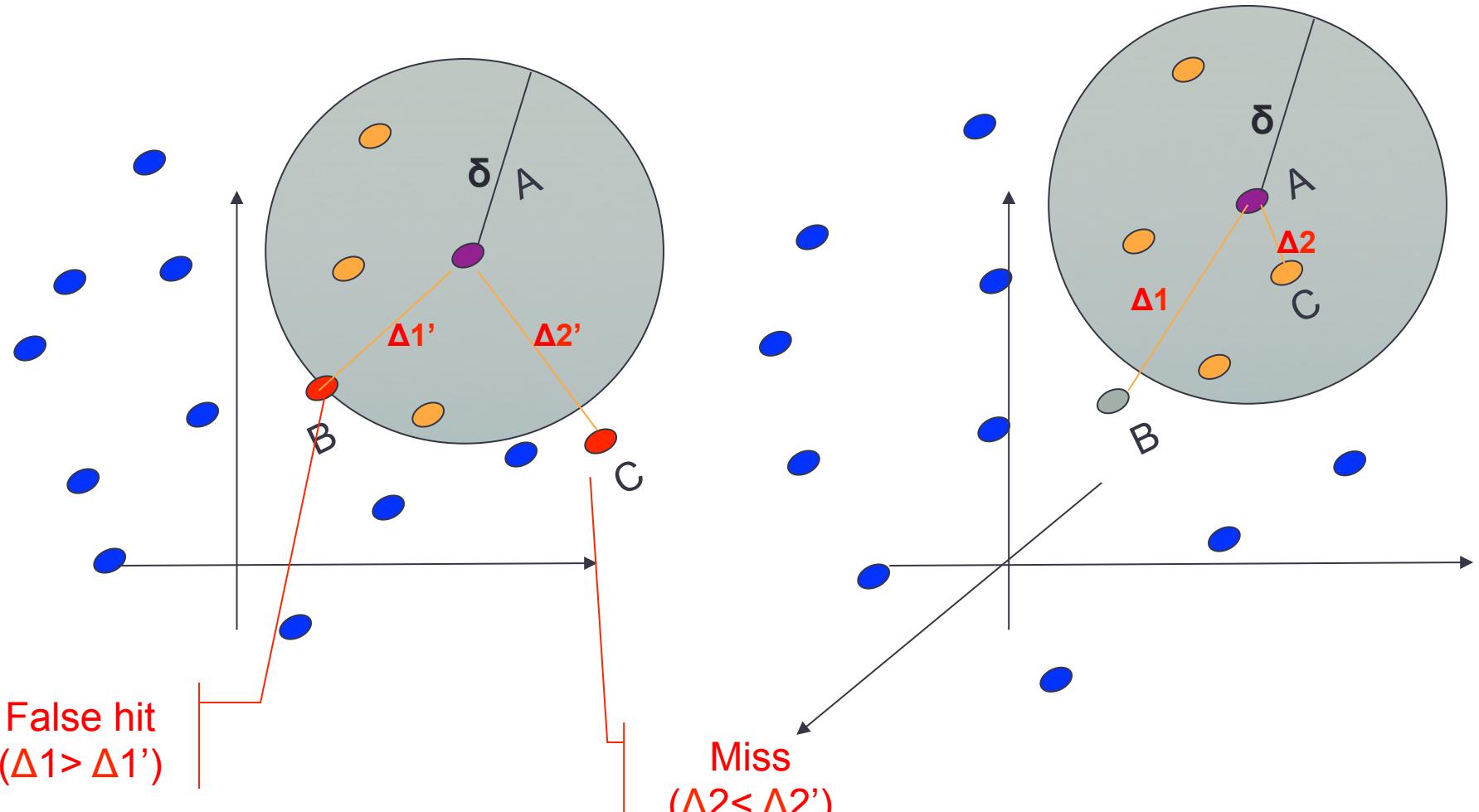


What happens to distances???



Misses are not desirable!
Can not be eliminated with post-processing

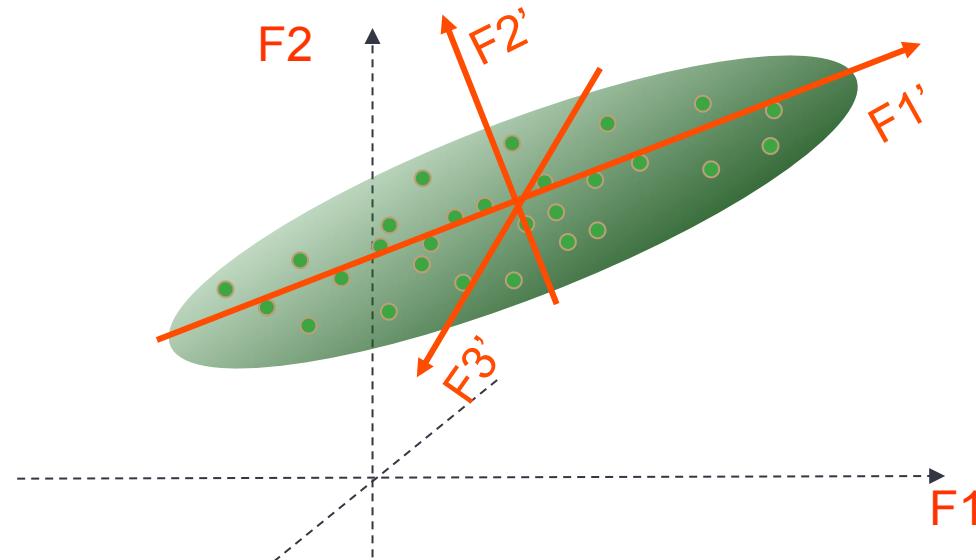
What happens to distances???



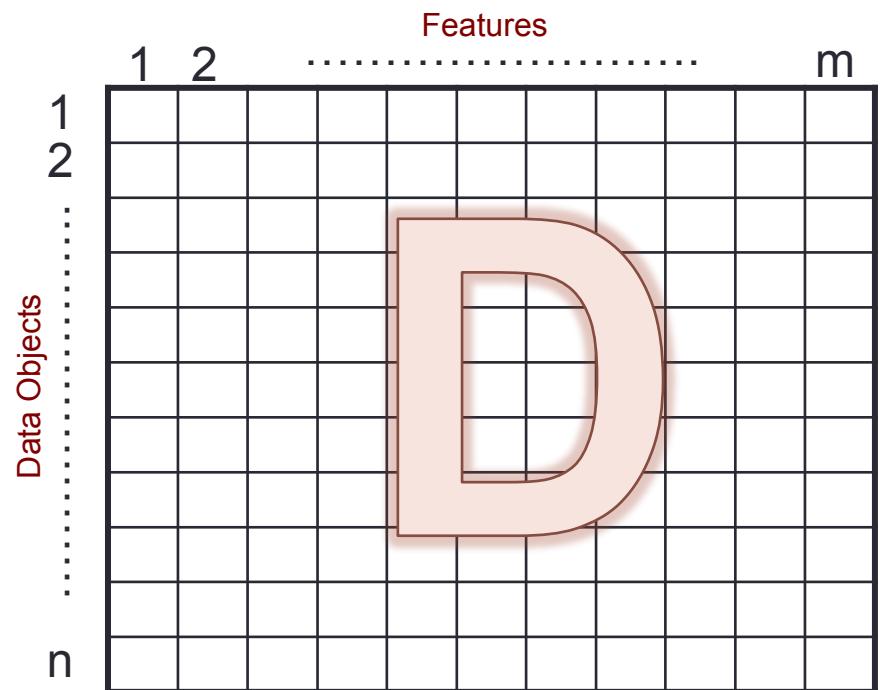
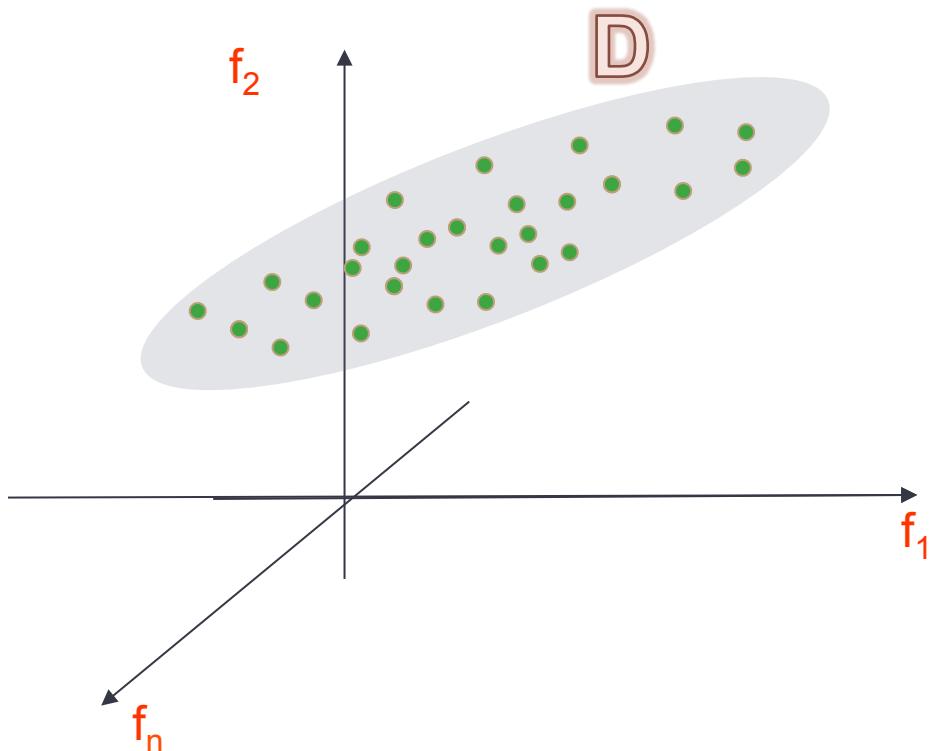
So we need data transforms where

....

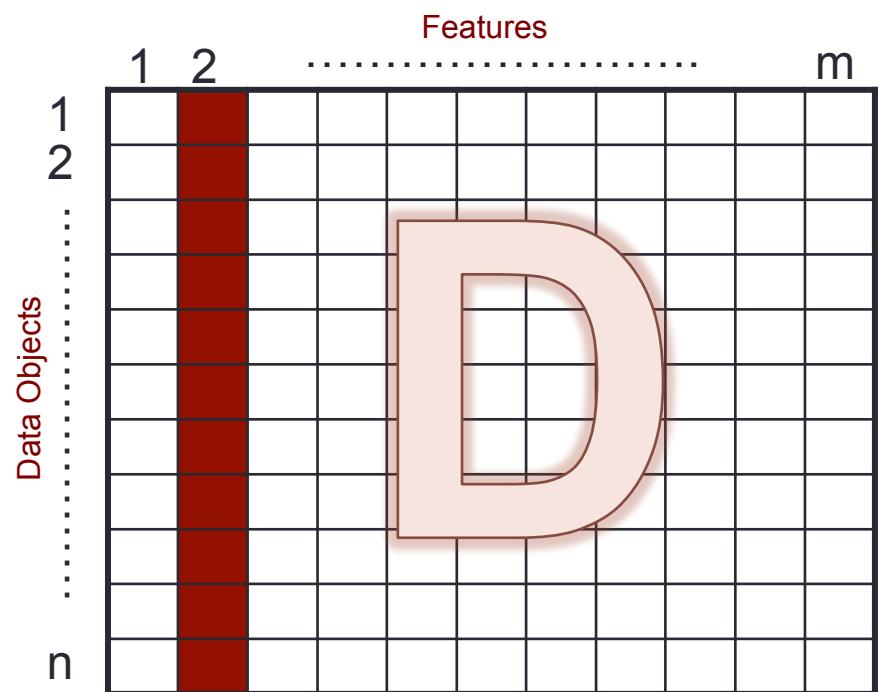
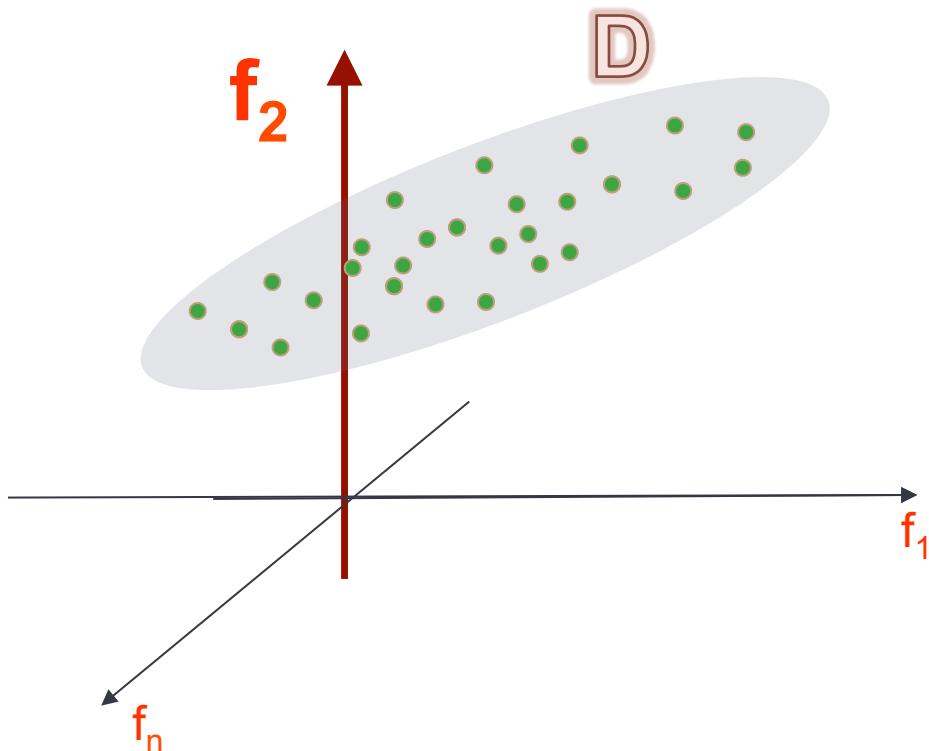
- distances and angles are preserved (e.g. linear, orthonormal transforms).
- the discriminatory power is associated with a few dimensions of the new space (**the dimensions are as little correlated as possible**)



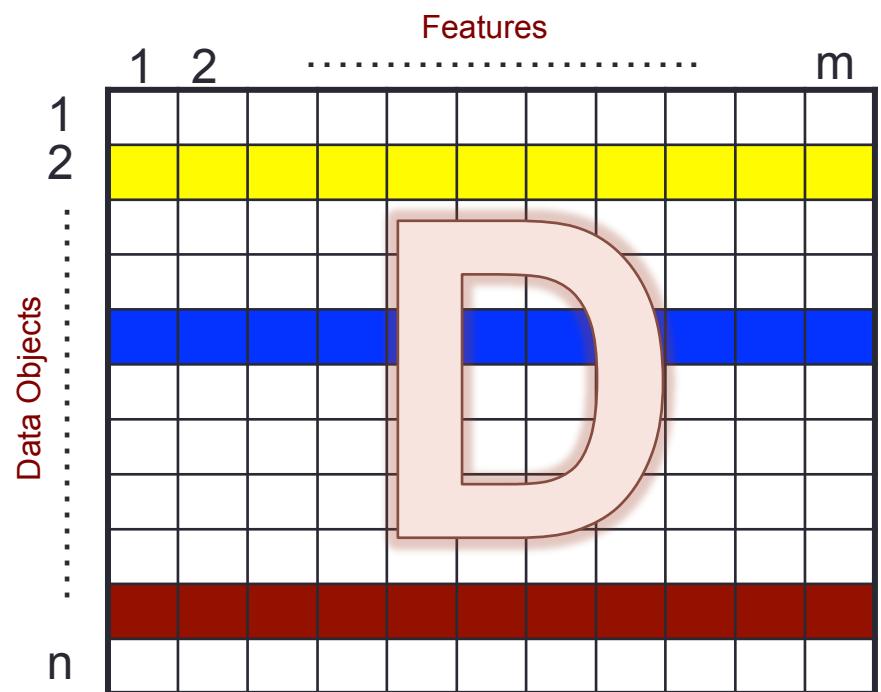
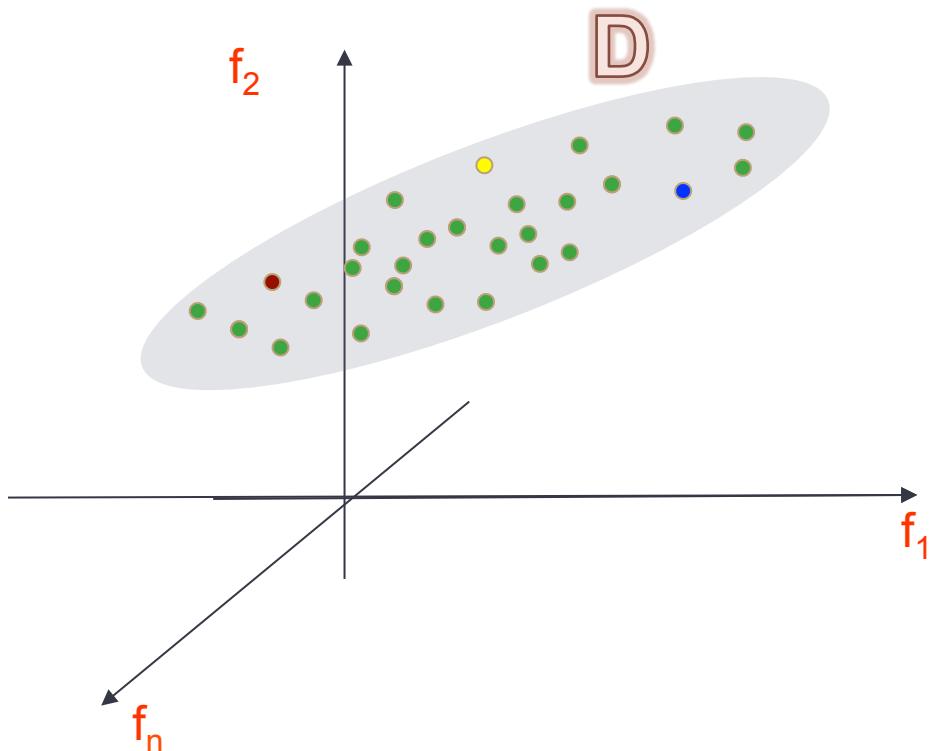
Singular valued decomposition (SVD)



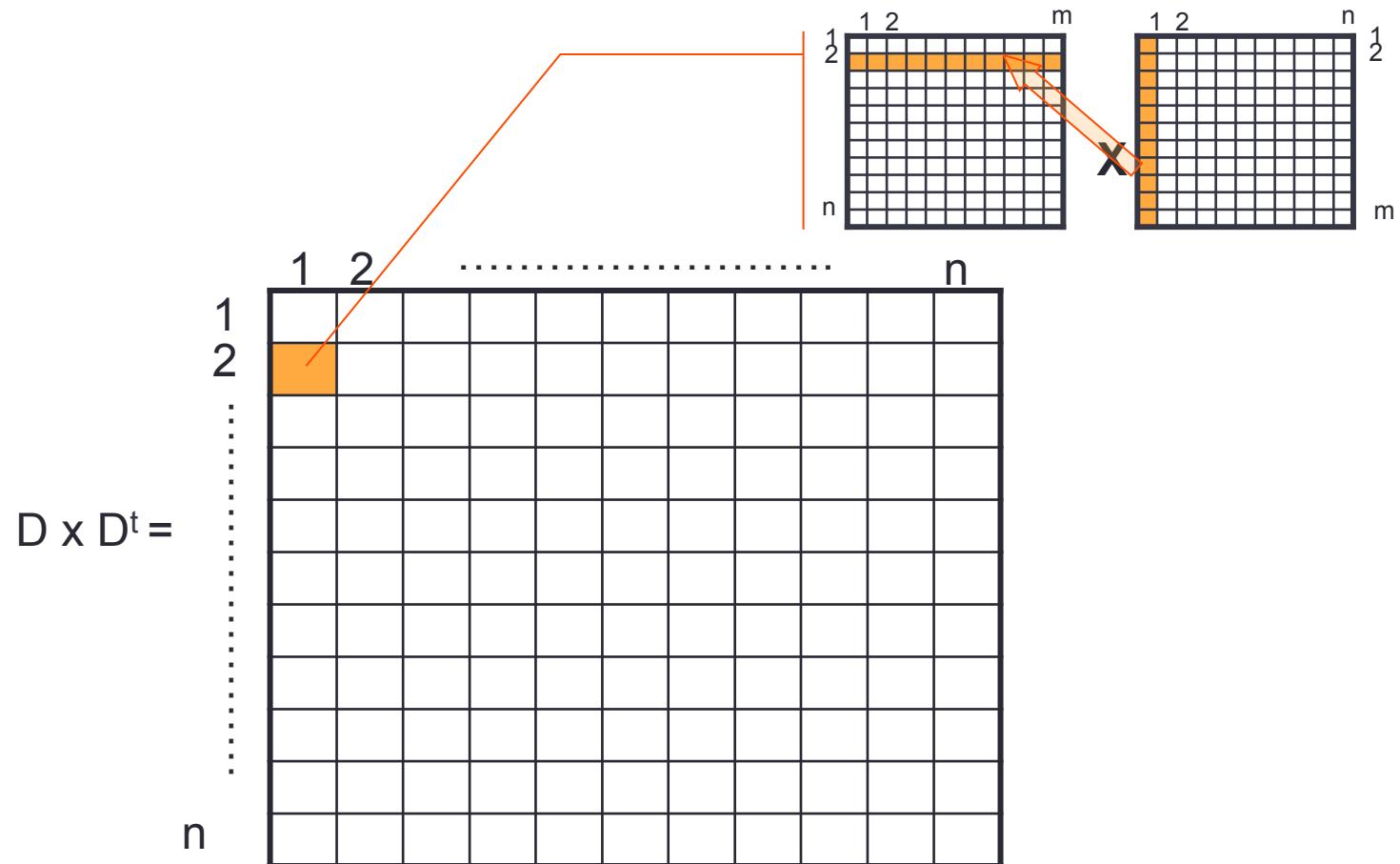
Singular valued decomposition (SVD)



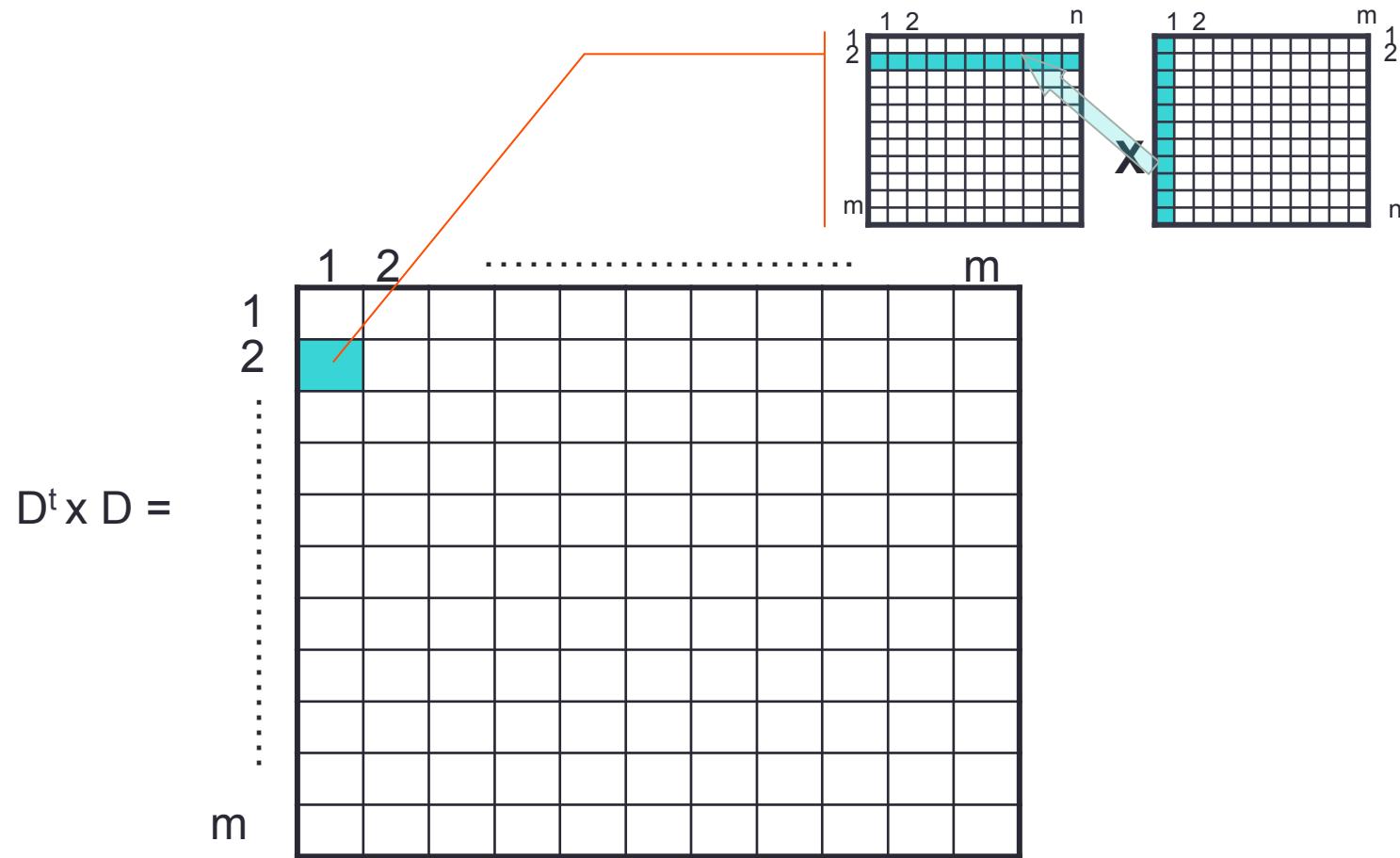
Singular valued decomposition (SVD)



Obj-obj similarity matrix!!!!



Feature-feature similarity matrix!!!!



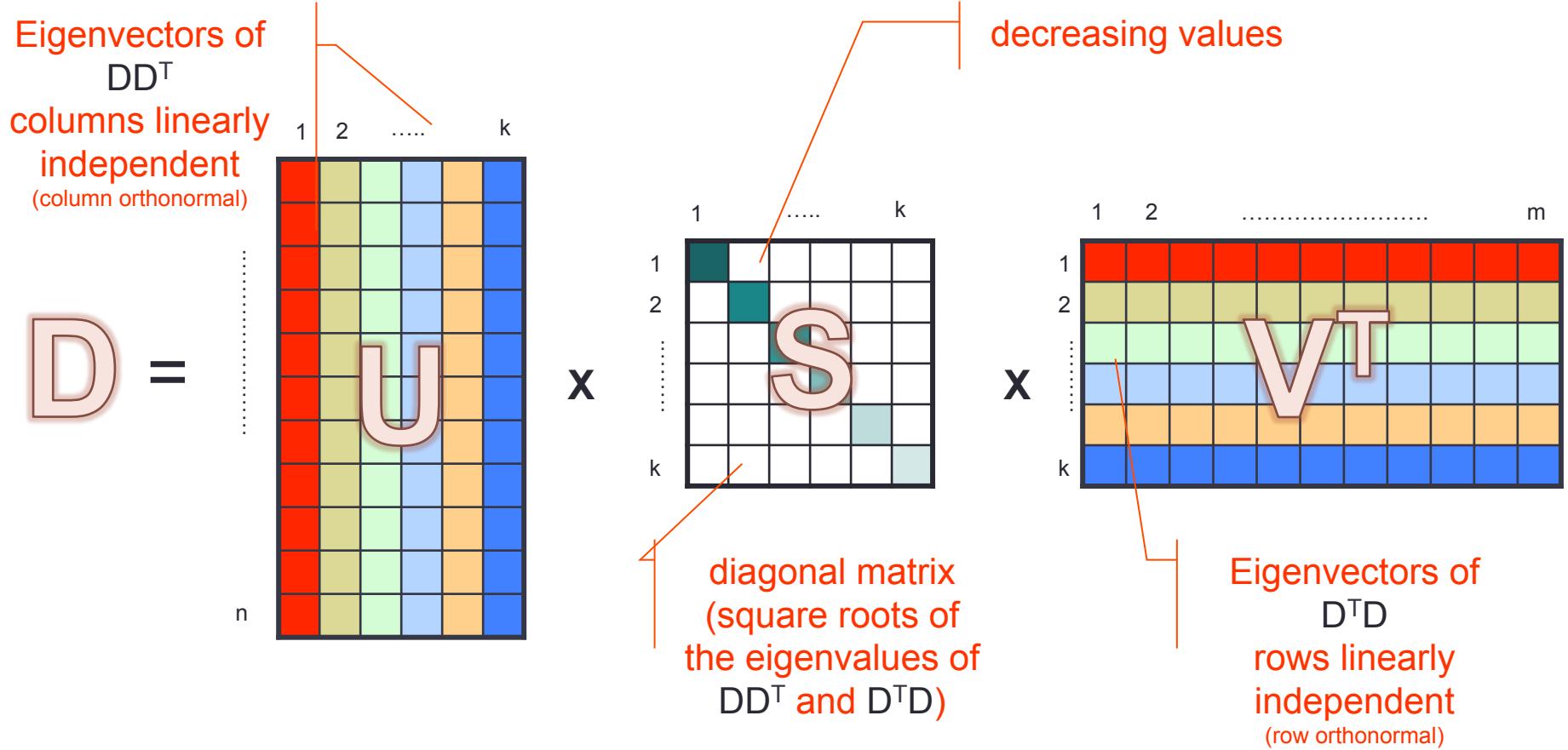
K. Selcuk Candan (CSE515)

Singular valued decomposition

$$D = \begin{matrix} & \begin{matrix} 1 & 2 & \dots & k \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ n \end{matrix} & \begin{matrix} \text{U} \\ \times \\ \text{S} \\ \times \\ \text{V}\text{T} \end{matrix} \end{matrix}$$

The diagram illustrates the Singular Value Decomposition (SVD) of a matrix D. Matrix D is shown as a grid with columns labeled 1, 2, ..., k and rows labeled 1, 2, ..., n. The matrix is decomposed into three components: U, S, and V^T. Matrix U is a tall matrix with n rows and k columns. Matrix S is a diagonal matrix with k rows and columns, containing singular values. Matrix V^T is a wide matrix with k rows and m columns.

Singular valued decomposition



...reminder

- Eigenvalue and eigenvector
- Given a matrix A , let c (scalar) and x (vector) be such that

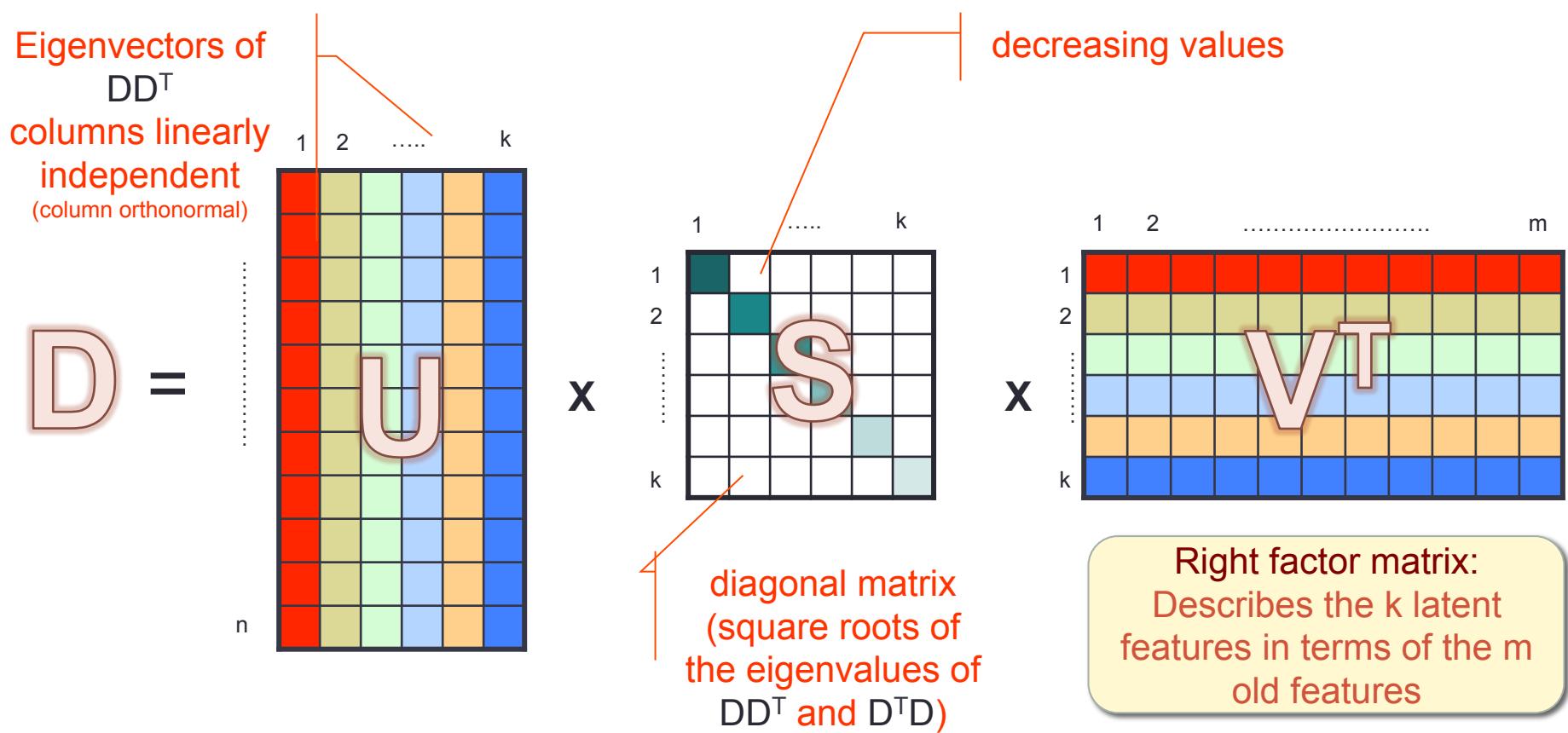
$$c x = A x$$

The diagram illustrates the relationship between eigenvalues and eigenvectors. It features a central equation $c x = A x$. Two red arrows originate from the words "Eigenvalue" and "Eigenvector" located below the equation. One arrow points to the scalar c , and the other points to the vector x .

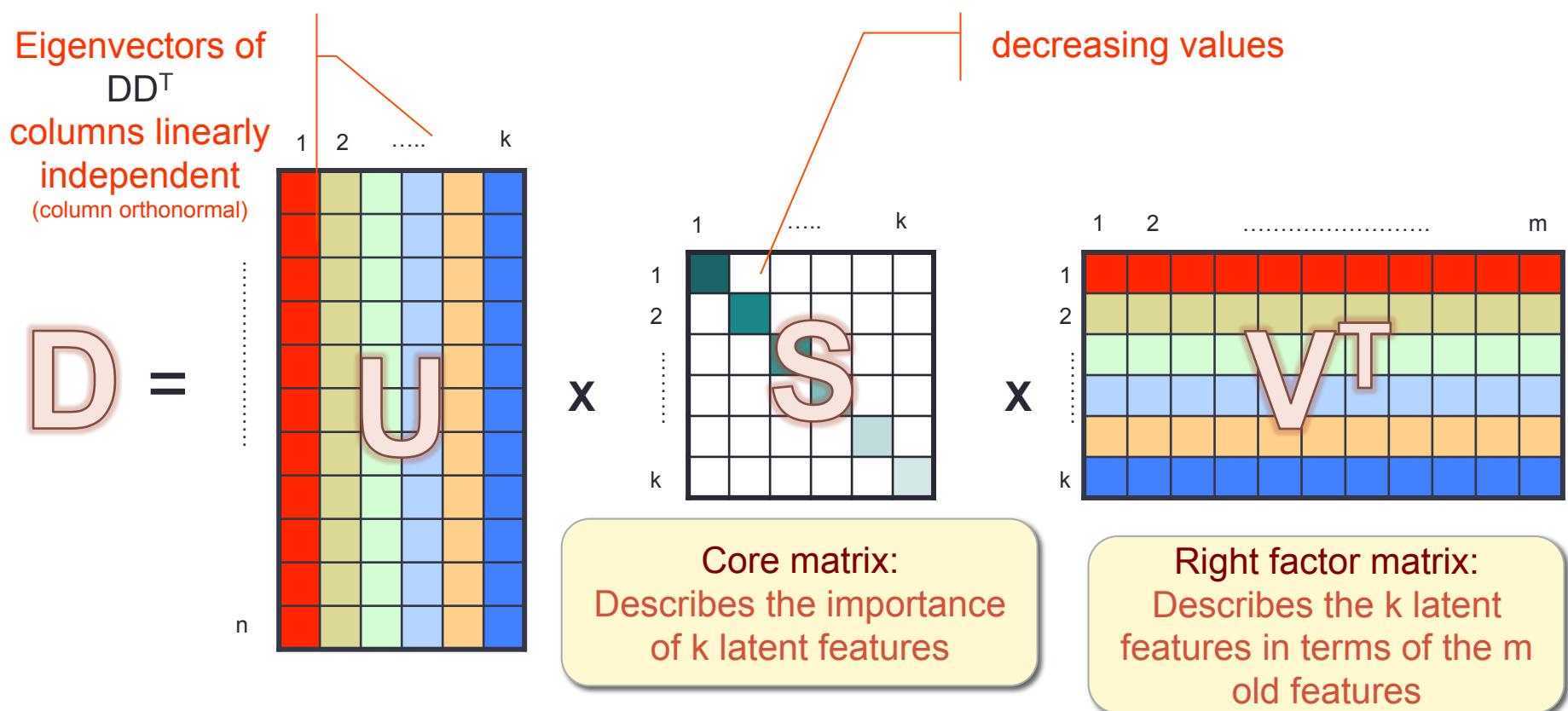
Properties of Eigenvectors

- Suppose that A is an $n \times n$ square matrix
 - if the eigenvalues, c_1, \dots, c_k , are distinct, then eigenvectors v_1, \dots, v_k are a set of k **linearly independent** vectors.
 - thus they can be used as the basis of the space!!!
 - The value of c_i describes the **contribution** of v_i in A.

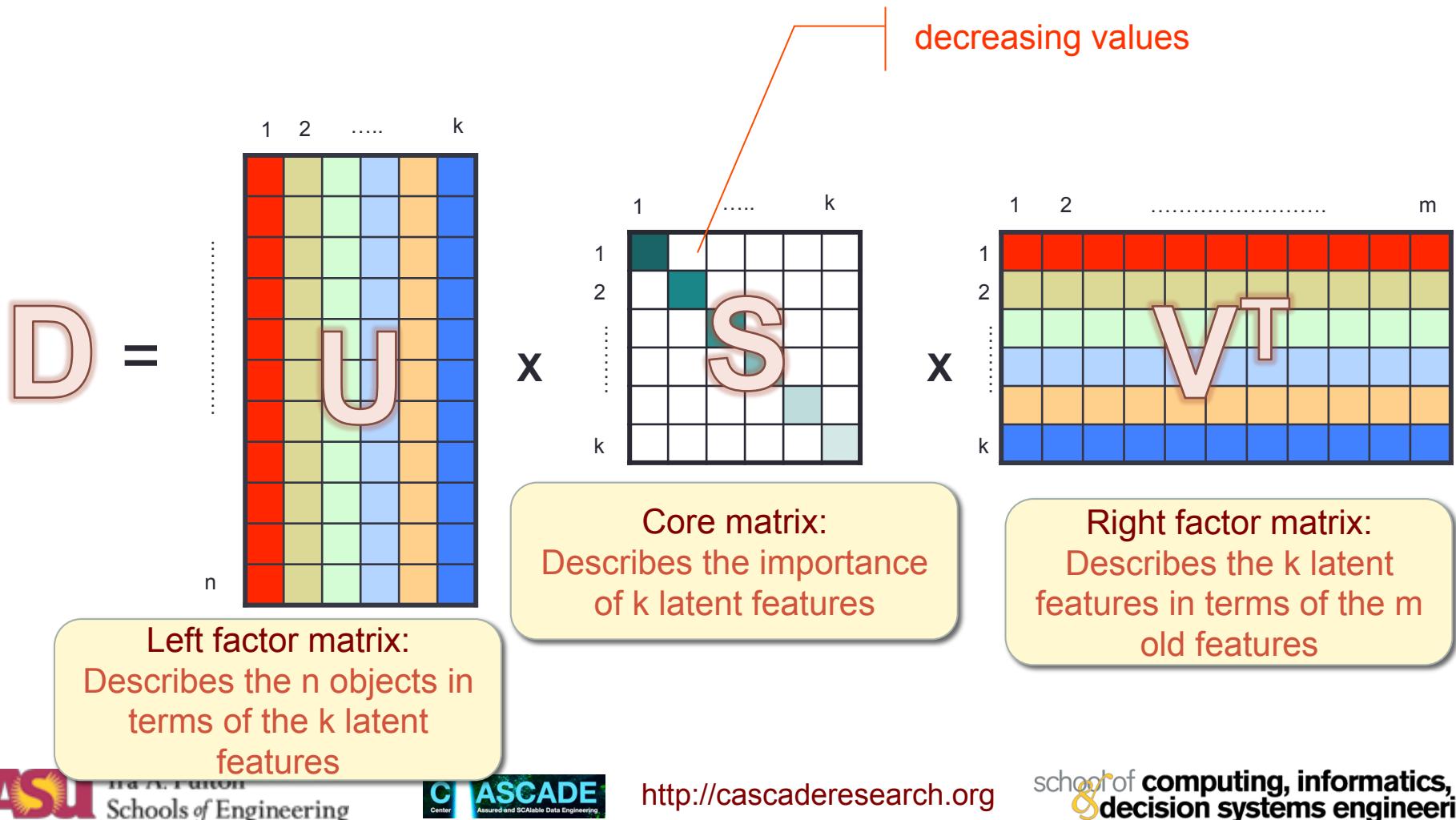
Singular valued decomposition



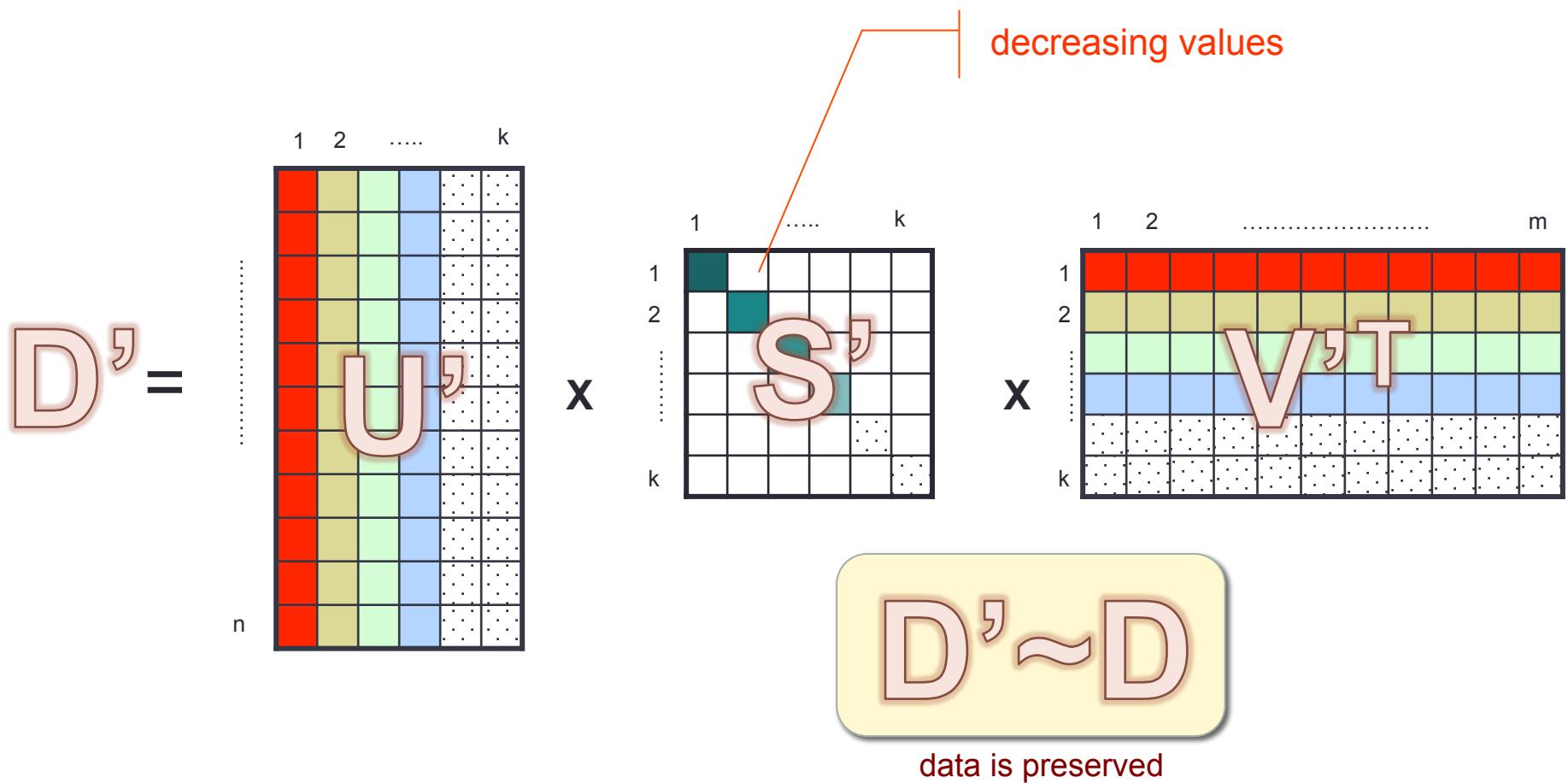
Singular valued decomposition



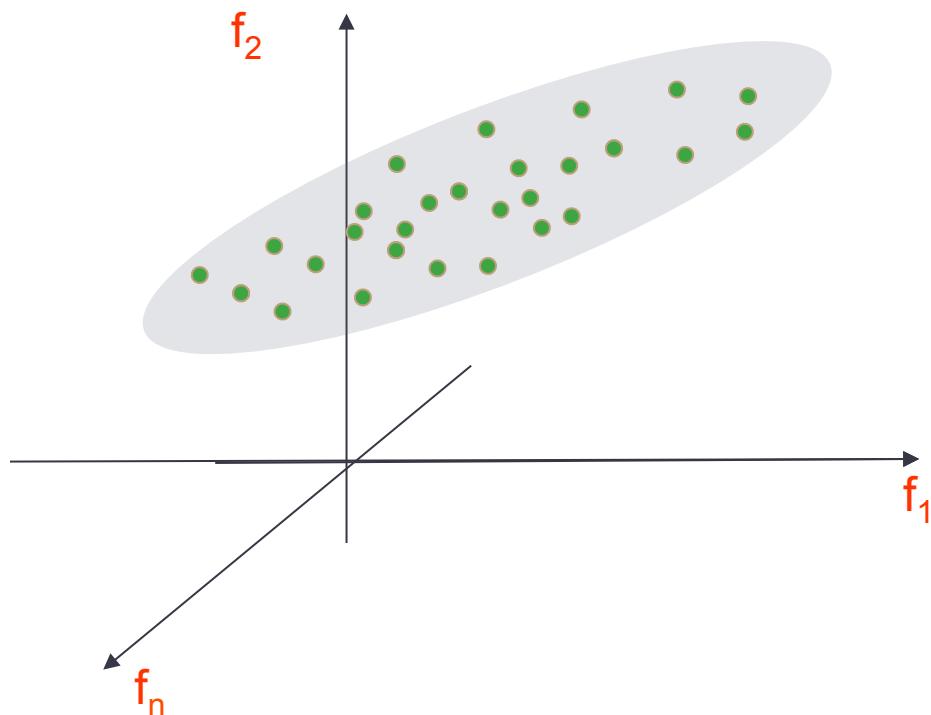
Singular valued decomposition



Singular valued decomposition



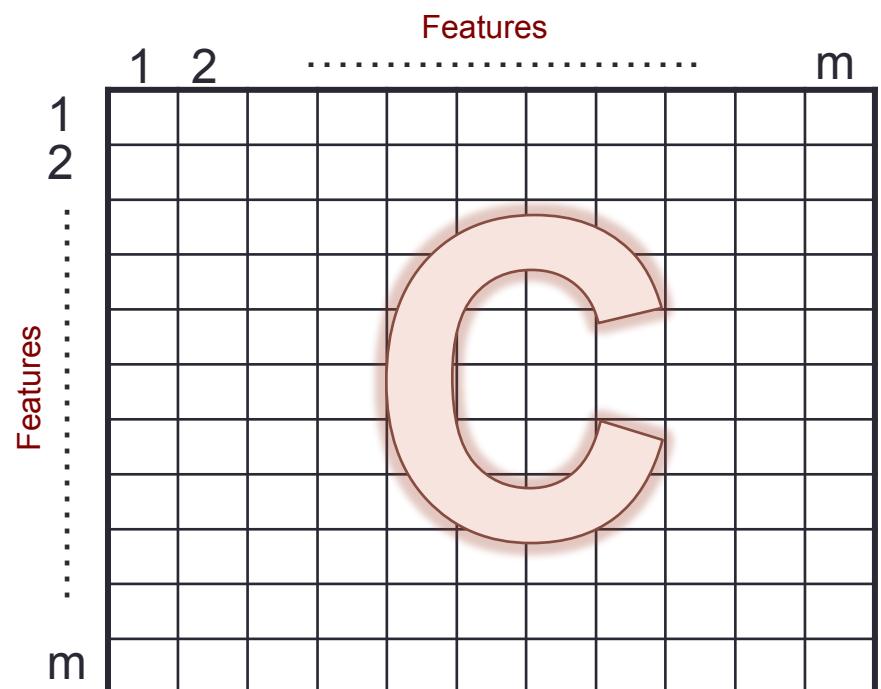
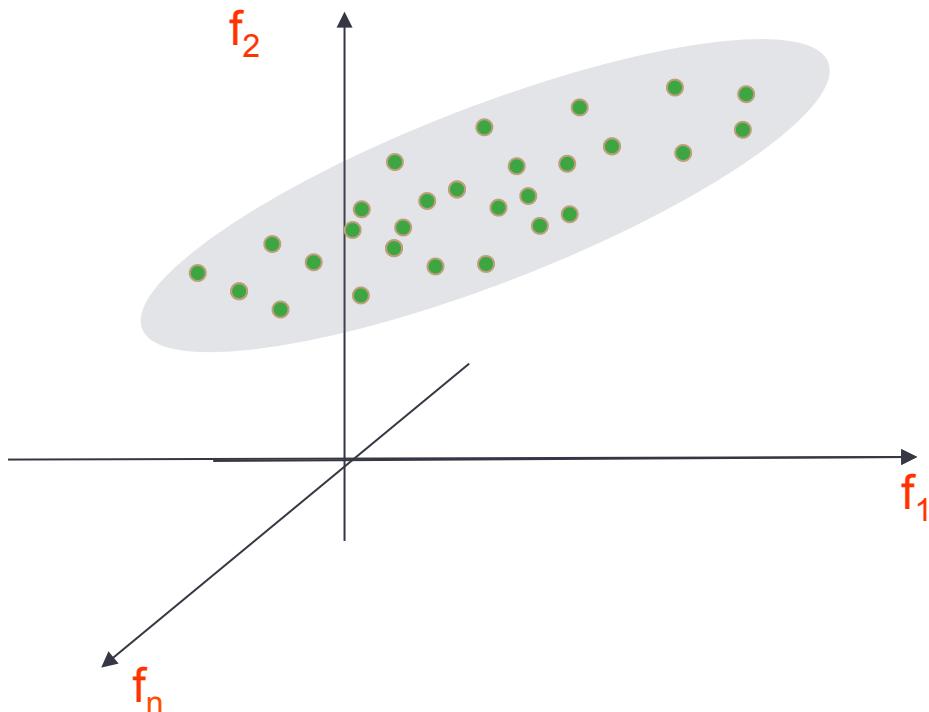
Principle Component Analysis (PCA)



Reminder: The higher the variance is, the better (features are more discriminating)

Principle Component Analysis (PCA)

- ..also known as Karhunen-Loeve Transform
 - ..a linear transform that optimally de-correlates the input.

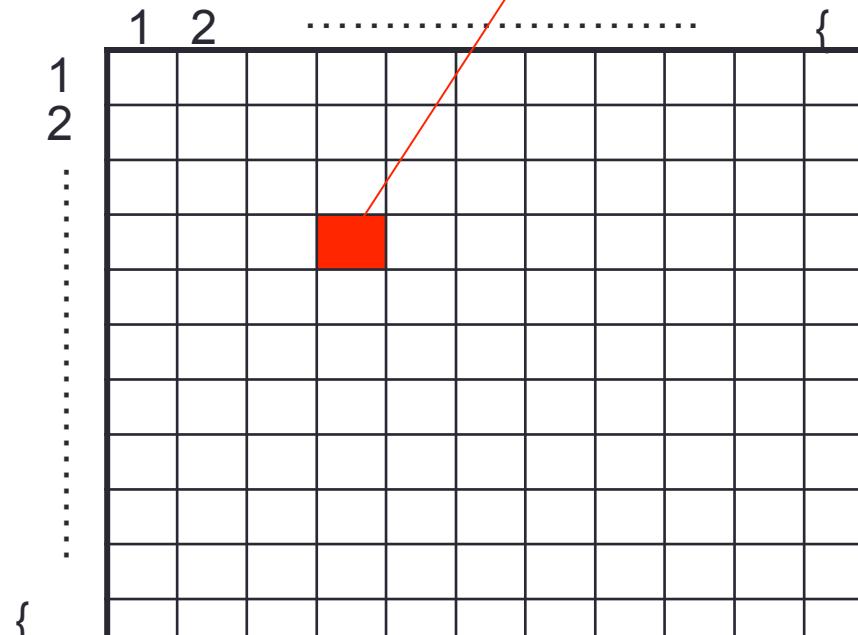


Covariance Matrix
(square symmetric)

$$C[i, j] = Cov(i, j) = E((i - \mu_i)(j - \mu_j))$$

Eigen decomposition

C =

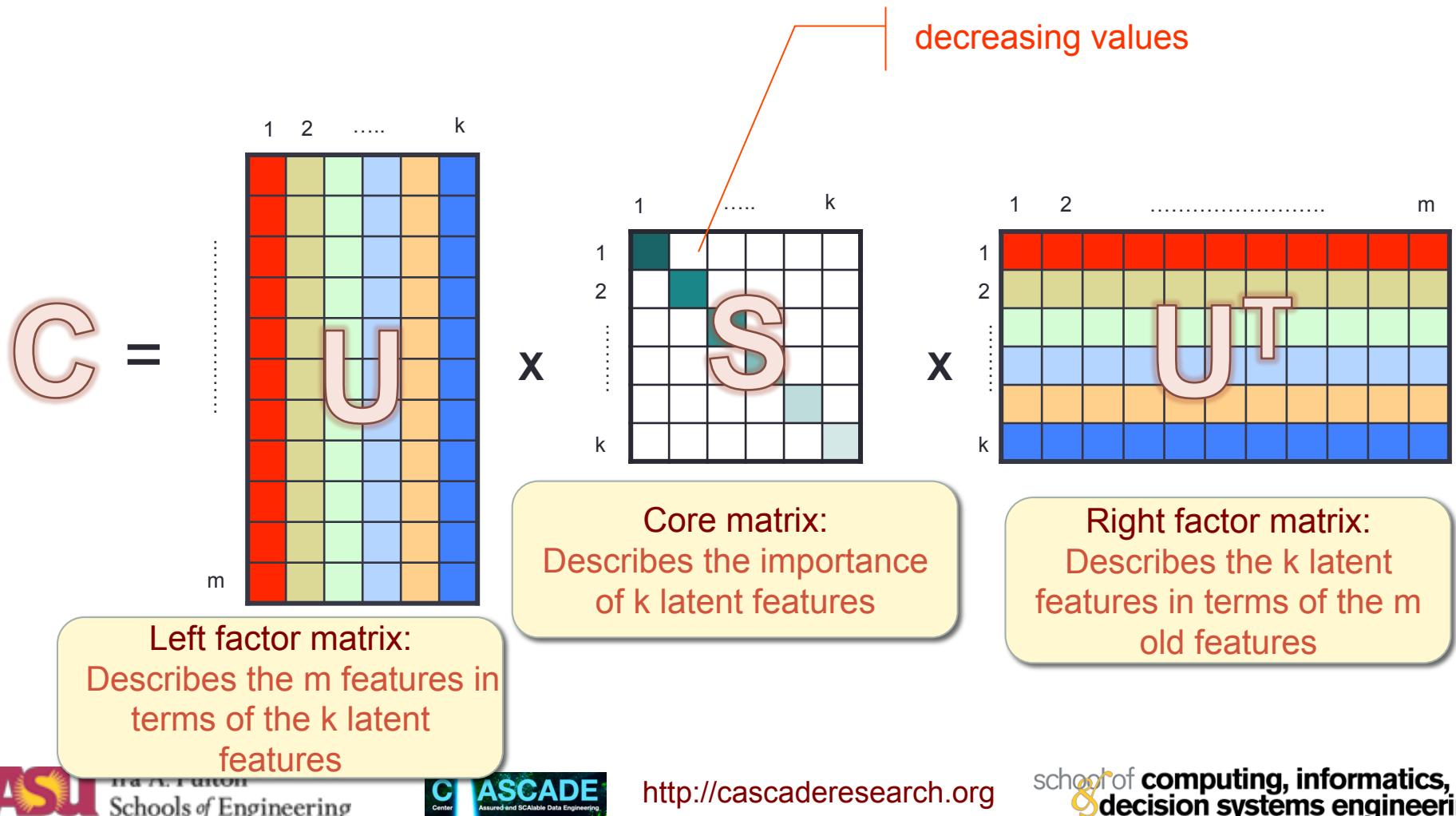


Covariance matrix

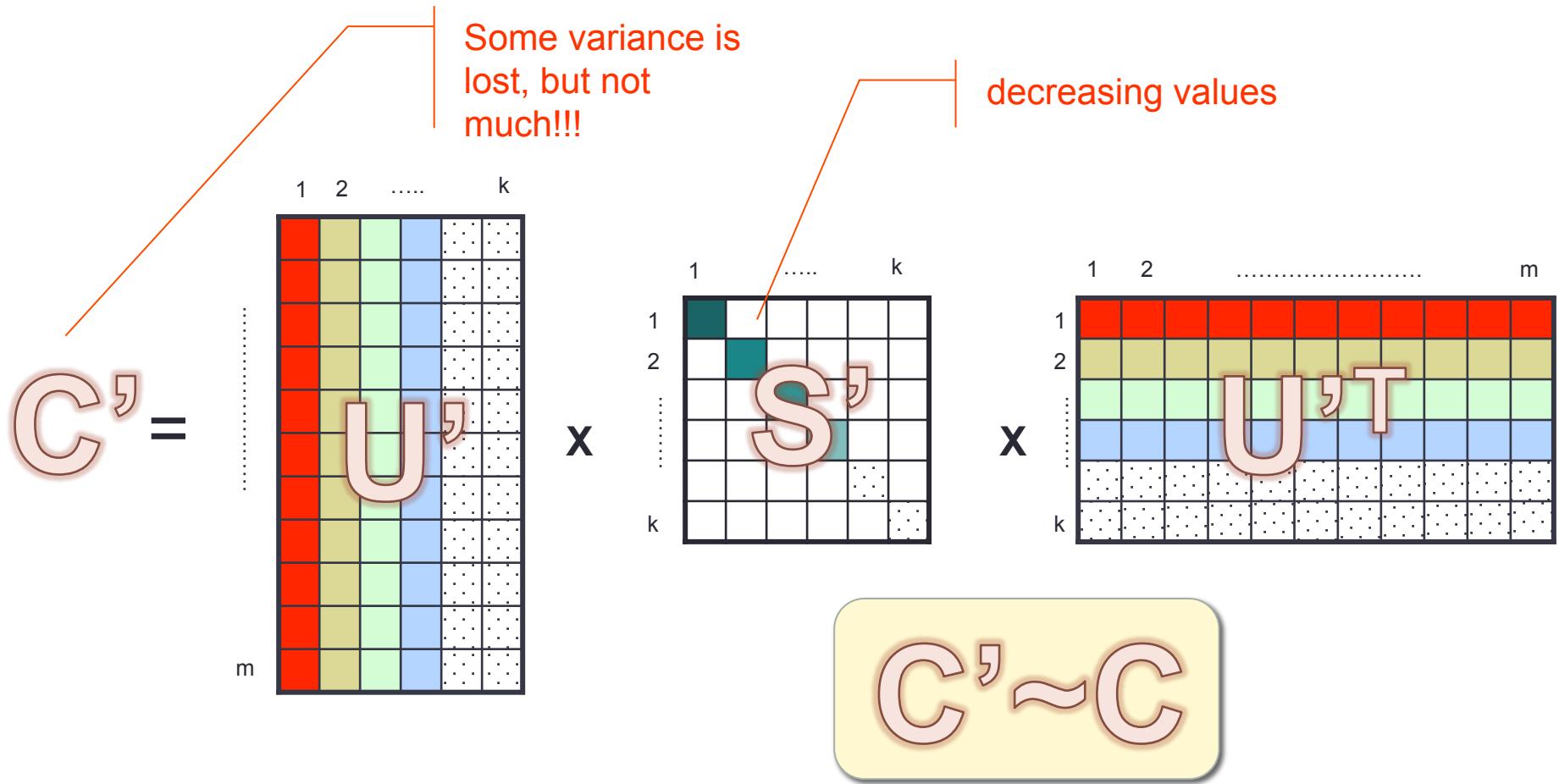
Reminder: Properties of Eigenvectors

- If \hat{O} is an $n \times n$ square covariance matrix
 - if the eigenvalues, c_1, \dots, c_k , are distinct, then eigenvectors v_1, \dots, v_k are a set of k **linearly independent** vectors.
 - thus they can be used as the basis of the space!!!
 - The value of c_i describes the contribution of v_i in A . Thus
 - if we pick a C that describes the variation of data,
 - v_i will describe **the directions along which variance is high**

Eigen decomposition of a symmetric covariance matrix



Principle Component Analysis (PCA)



How many eigenvectors shall we maintain?

- **Mean eigenvalue:** use only the dimensions whose eigenvalues are greater than or equal to the mean eigenvalue.
- **Kaiser rule:** keep only those eigenvectors whose eigenvalues are greater than 1 (applied when A is the correlation matrix)
 - eigenvalue is the amount of variance explained by one or more latent semantics,
 - don't add a semantics that explains less variance than is contained in one dimension when all eigenvalues of a correlation matrix are exactly one
- **Parallel analysis:**
 - analyze a random covariance matrix,
 - Plot cumulative eigenvalues for both random and intended matrices;
 - Find where the two curves intersect.
- **Scree test:** plot the successive eigenvalues to find a point where the plot levels off.
- **Variance explained:** keep enough dimensions to account for 95% of the initial variance