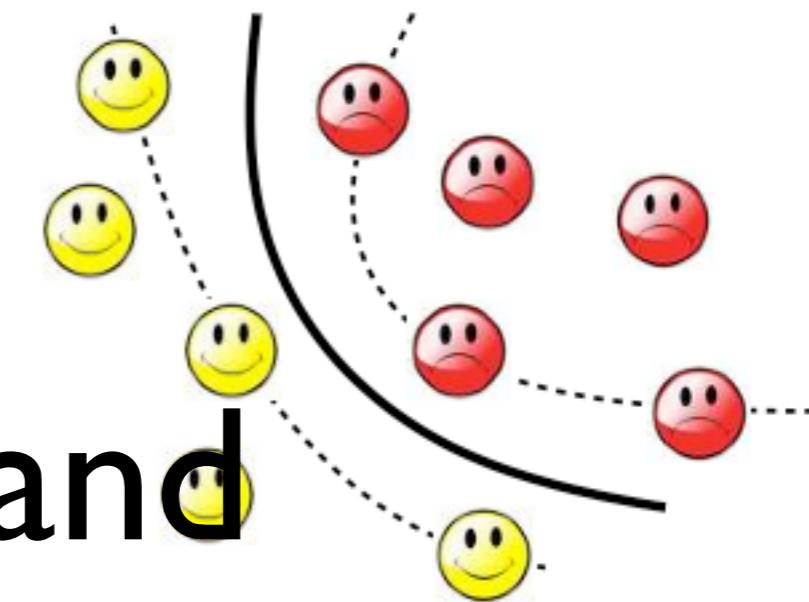




THE UNIVERSITY OF
SYDNEY

Machine Learning and Data Mining



Basic Matrix Analysis and Singular Value
Decomposition

Dr Tongliang Liu



THE UNIVERSITY OF
SYDNEY

What is Machine Learning?

Informally: Making predictions from data

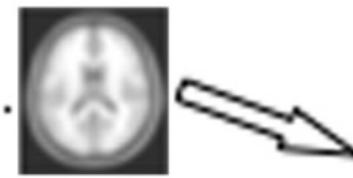
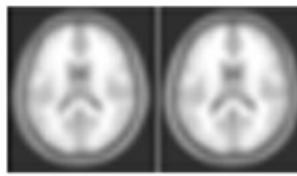
Formally: The construction of a statistical model that is an underlying distribution from which the data is drawn from.



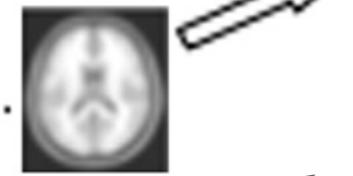
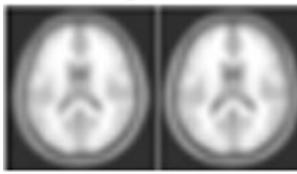
Elements of Machine Learning

Input training data

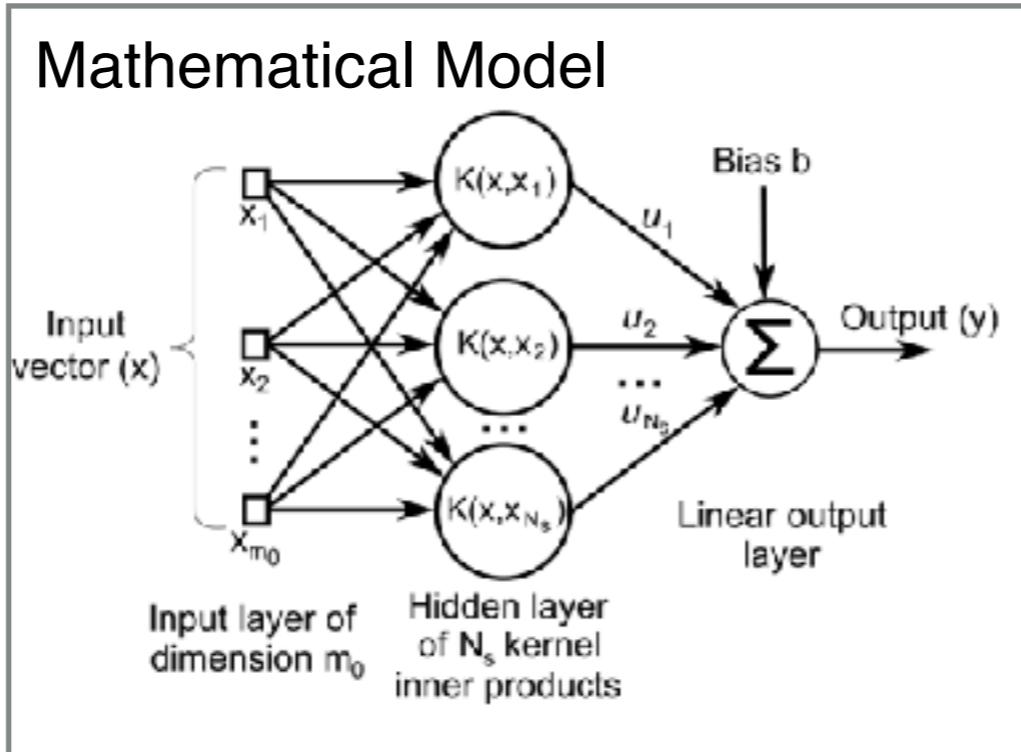
Group 1



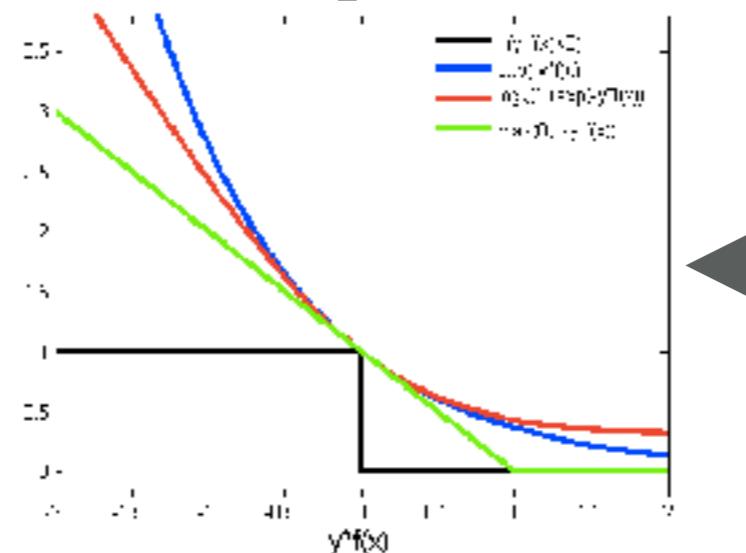
Group 2



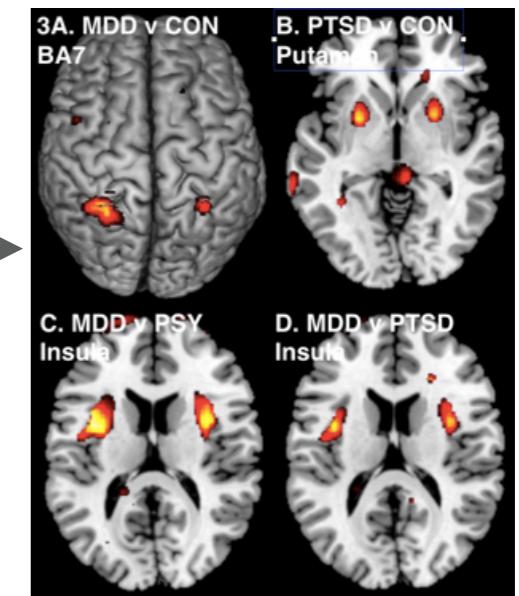
Data



Input predefined function class



Output hypothesis
Predictions/Patterns



Optimisation method



THE UNIVERSITY OF
SYDNEY

Iterations: 000,000 Learning rate: 0.03 Activation: ReLU Regularization: L2 Regularization rate: 0.001 Problem type: Regression

DATA

Which dataset do you want to use?



Ratio of training to test data: 40%



Noise: 20



Batch size: 10



REGENERATE

FEATURES

Which properties do you want to feed in?

+ -

4 neurons

x_1

x_2

x_1^2

x_2^2

$x_1 x_2$

$\sin(x_1)$

$\sin(x_2)$

+ -

3 HIDDEN LAYERS

+ -

4 neurons

x_1

x_2

x_1^2

x_2^2

$x_1 x_2$

$\sin(x_1)$

$\sin(x_2)$

+ -

2 neurons

x_1

x_2

x_1^2

x_2^2

$x_1 x_2$

$\sin(x_1)$

$\sin(x_2)$

+ -

2 neurons

x_1

x_2

x_1^2

x_2^2

$x_1 x_2$

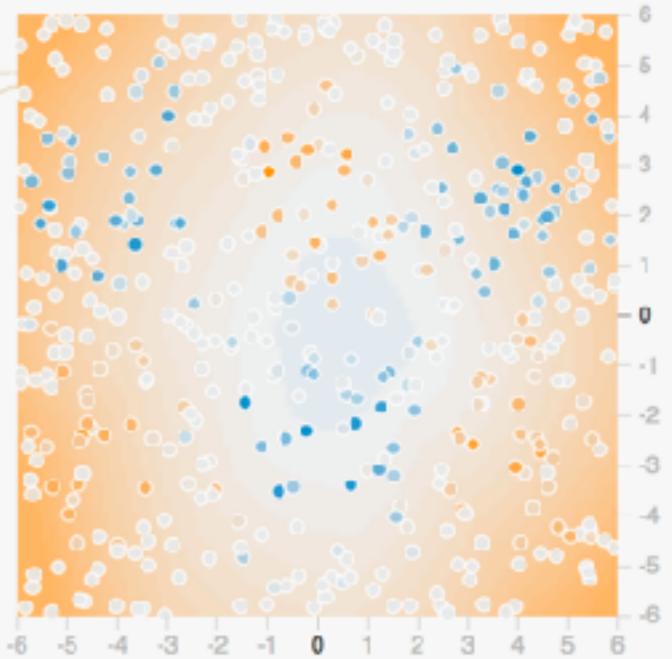
$\sin(x_1)$

$\sin(x_2)$

OUTPUT

Test loss 0.121

Training loss 0.134



Colors show data, neuron and weight values.



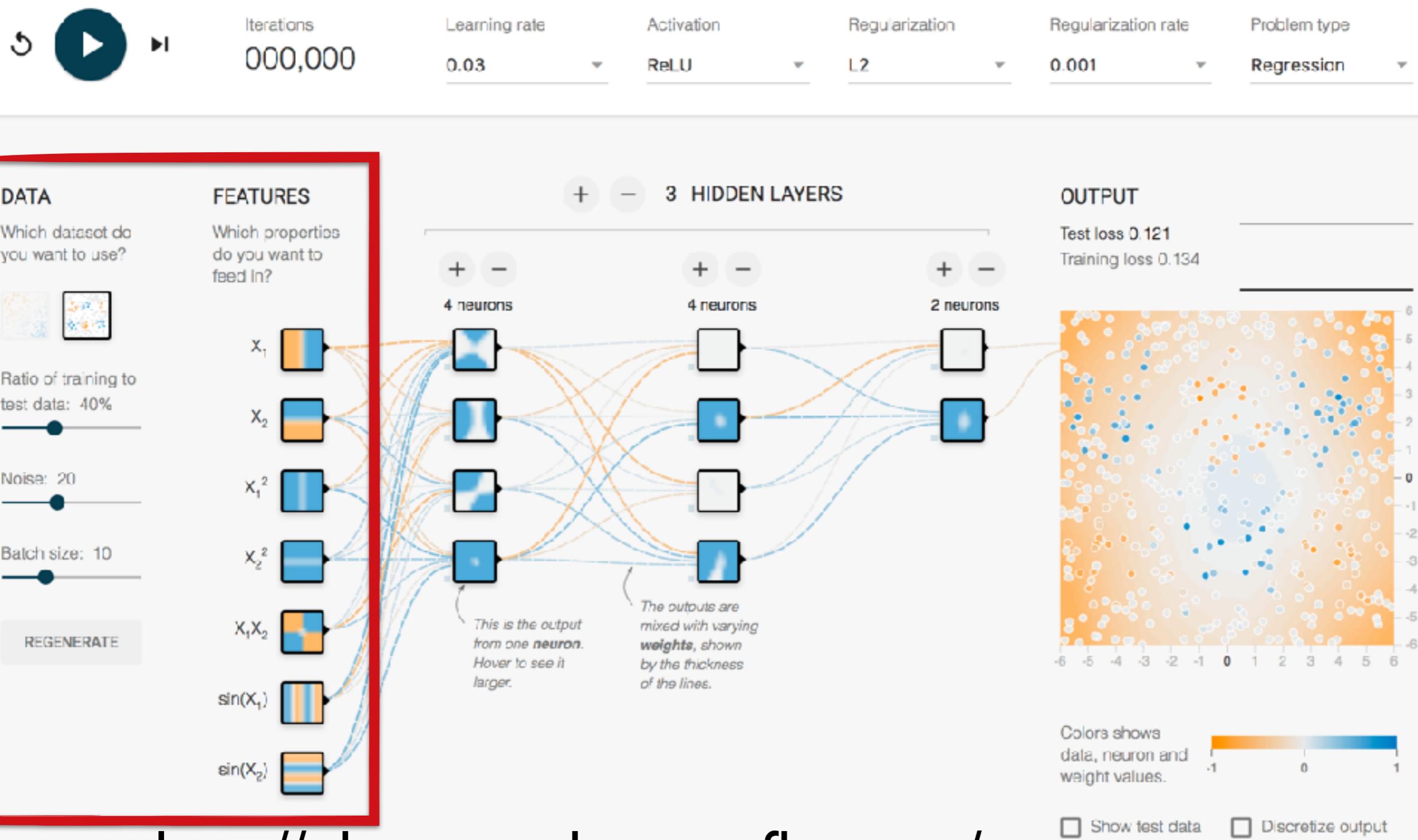
Show test data

Discretize output

Source: <http://playground.tensorflow.org/>



THE UNIVERSITY OF
SYDNEY



Source: <http://playground.tensorflow.org/>



Common representation

IMAGE/
VIDEO

TEXT/
COMMENT

TIME
SERIES

SYSTEM
LOGS

NETWORK

TABULAR/
RATING

Is there a common way to represent data
of different modalities ?



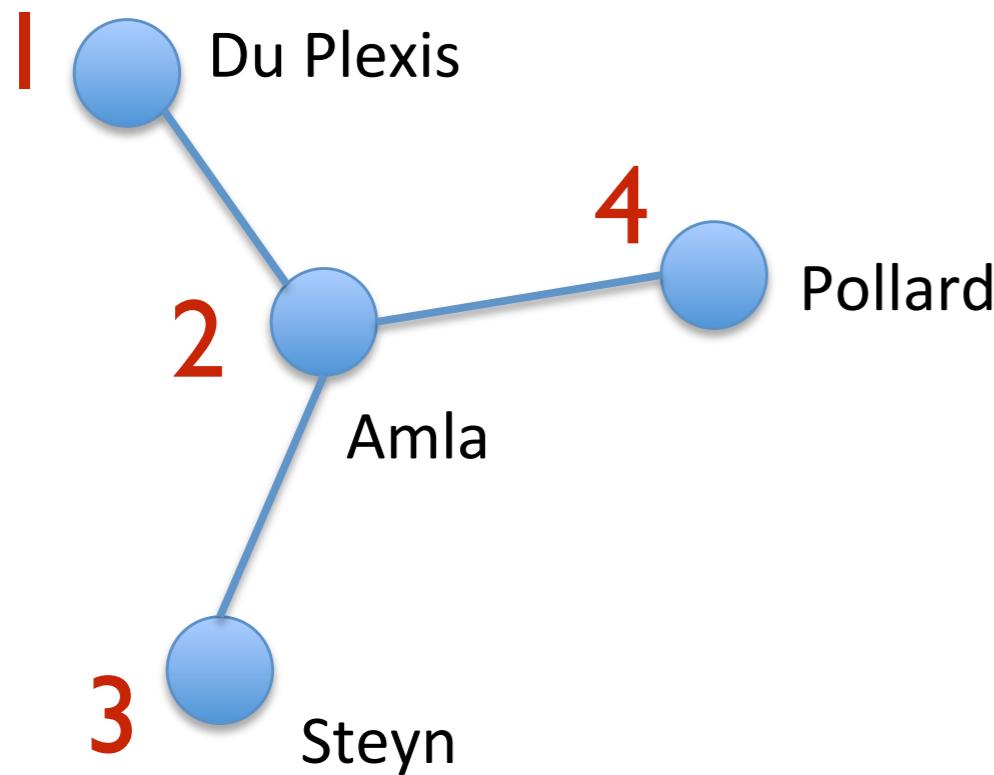
Text to matrix

- Document- Word Matrix
- Document 1:“AACCBBAAA”
- Document 2:“CCAABBDD”

$$\begin{bmatrix} A & B & C & D \\ 5 & 2 & 2 & 0 \\ 2 & 2 & 2 & 2 \end{bmatrix}$$



Network data



Nodes	Nodes	Nodes	Nodes
0	1	0	0
1	0	1	1
0	1	0	0
0	1	0	0



THE UNIVERSITY OF
SYDNEY

Image data



www.sydney.visitorsbureau.com.au



700 x 500

4	45	6
6	12	33
22	17	44



4	45	6	6	12	33	22	17	44
---	----	---	---	----	----	----	----	----



THE UNIVERSITY OF
SYDNEY

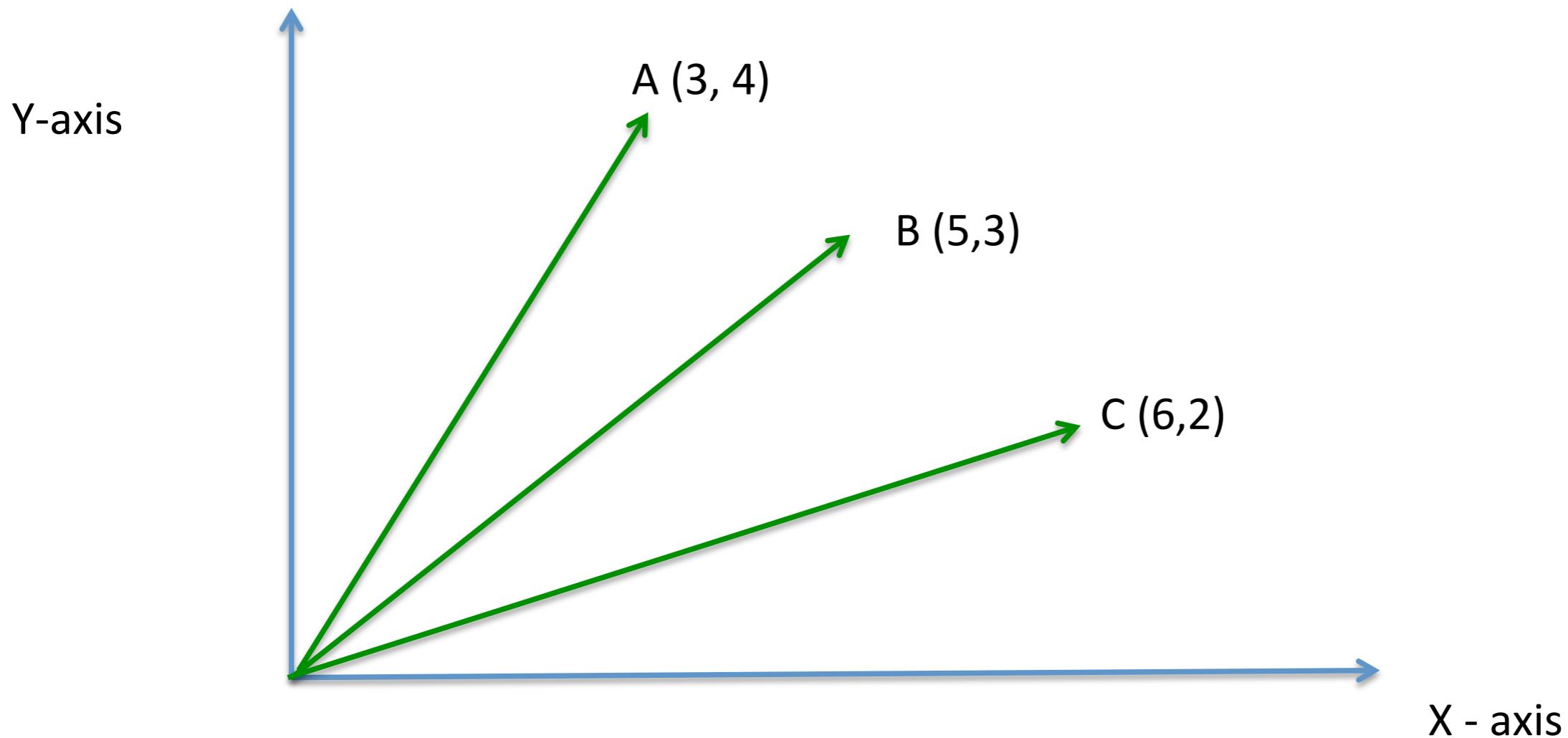
Similarity Computation

- We can now represent most data types as a matrix.
- A special case of a matrix is a vector.
- Now lets compute similarities with these objects.



Similarity Computation

How can we quantify similarity between A, B and C ?





Similarity Computation

- Dot product

$$x = (x_1, x_2, \dots, x_n); \quad y = (y_1, y_2, \dots, y_n);$$

$$x.y = (x_1y_1 + x_2y_2 + \dots + x_ny_n);$$

- Norm (length) of a vector

$$\|x\| = (x.x)^{1/2} = (x_1.x_1 + x_2.x_2 + x_n.x_n)^{1/2}$$



THE UNIVERSITY OF
SYDNEY

Similarity Computation

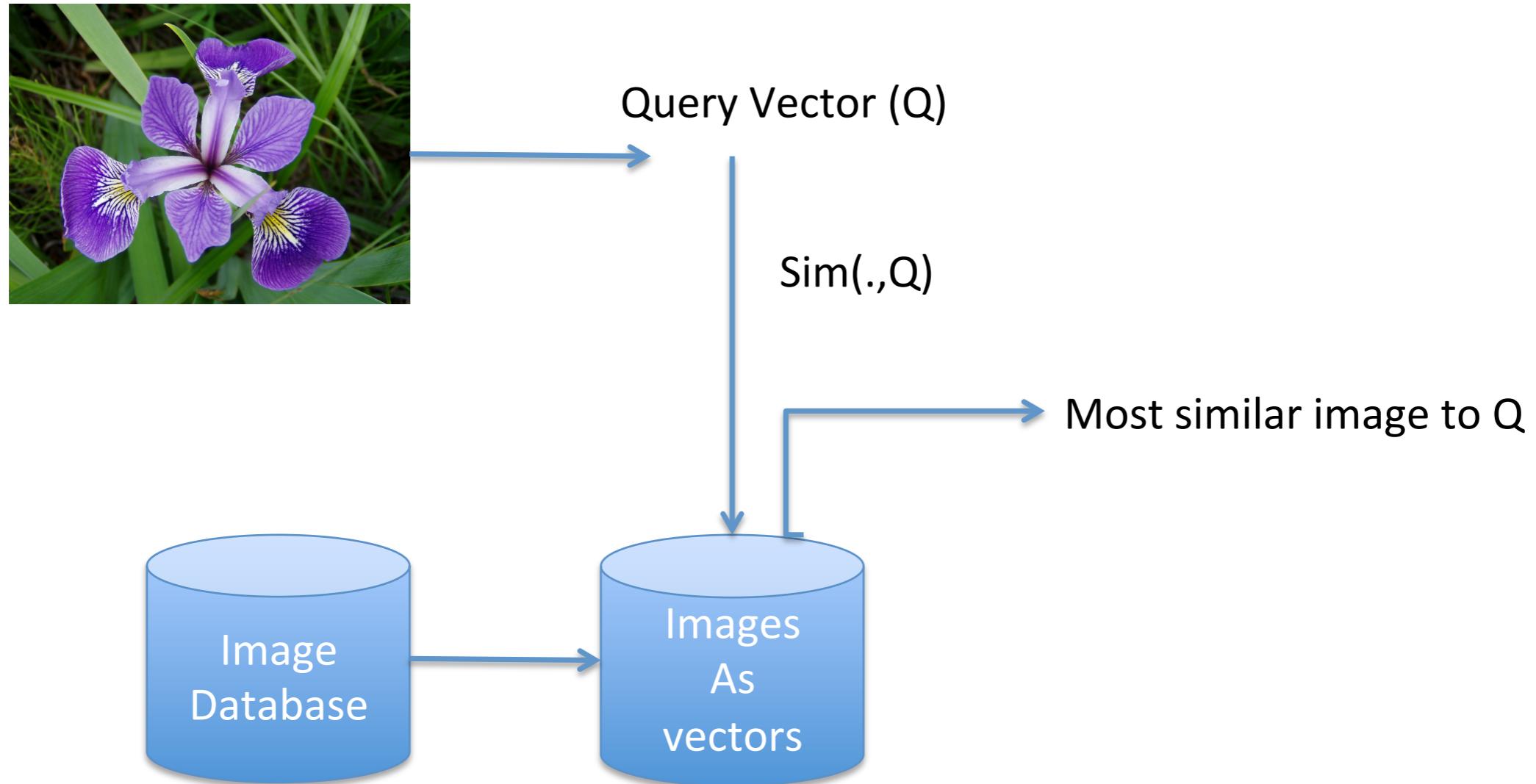
- The similarity between two vectors x and y is given by

$$sim(x, y) = x \cdot y / (\|x\| \|y\|)$$



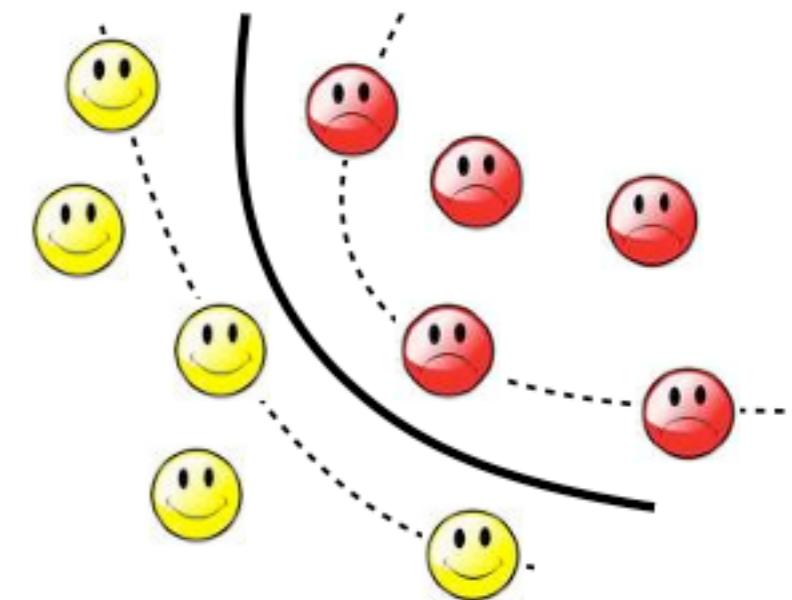
THE UNIVERSITY OF
SYDNEY

Image search engine





THE UNIVERSITY OF
SYDNEY



Matrix Algebra



Why Matrix Algebra?

- Data Mining: Computation for large data sets
- Key idea: Data Reduction leads to “Knowledge Discovery”
- Matrix algebra provides simple algorithms with great power for data management, e.g., data reduction or data summarisation.



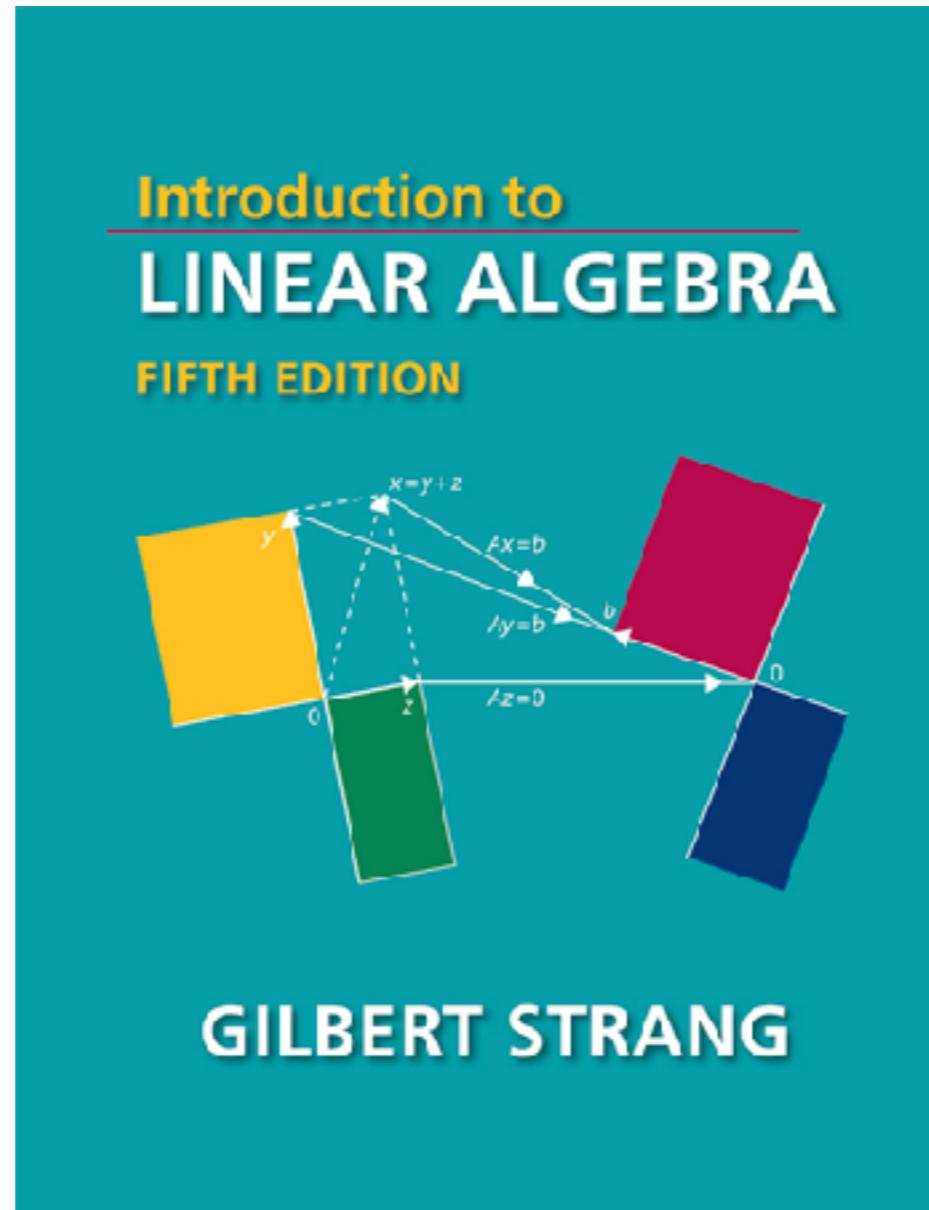
Why Matrix Algebra?

- In database management, the fundamental entity is a *relation*.
 - SQL (relational algebra) is a collection of operations on relations – Select, project, join, group-by
- In machine learning and data mining, the fundamental entity is a *matrix*
 - We therefore need to know the basic operations on matrices
 - One important application is to summarise data or extract patterns from data or compress data
 - In SIT several people apply data mining for different domains: Chinese medicine, bioinformatics, student learning, multimedia, image processing, quality of cloud computing service, text analysis, medical imaging, robotics



THE UNIVERSITY OF
SYDNEY

Suggested book



<http://math.mit.edu/~gs/linearalgebra/>



Numbers vs Matrices

Numbers

- Can add and subtract numbers
- Multiply numbers
- Divide two numbers a/b as long as b is not equal to 0.
- Can factorise positive numbers into product of primes

Matrices

- Can add and subtract compatible metrics
- Multiply and divide matrices
- Division of matrices is complicated
- **Can factorise any matrix to get data patterns (using a technique called Singular Value Decomposition)**



Linear Algebra

- Area in maths that deals with *vector spaces* and *linear mappings* between these spaces
- The generalisation of LA to infinite dimensions is known as *functional analysis*

2	4	7	3	6
---	---	---	---	---

$$D = 5$$

Linear algebra

2	4	7	3	6	9	...
---	---	---	---	---	---	-----

$$D = \infty$$

Functional analysis



Definition for vector space: 8 Axioms

- Associativity of addition: $(u + v) + w = u + (v + w)$
- Commutativity of addition: $u + v = v + u$
- Identity element of addition: $\exists 0, v + 0 = v$
- Inverse elements of addition: $\forall v, \exists -v, v + (-v) = 0$
- Compatibility of scalar multiplication with field multiplication: (a and b are scalars) $a(bv) = (ab)v$
- Identity element of scalar multiplication: $\exists 1, 1v = v$
- Distributivity of scalar multiplication with respect to vector addition: (a is a scalar)
$$a(u + v) = au + av$$
- Distributivity of scalar multiplication with respect to field addition: (a and b are scalars)
$$(a + b)v = av + bv$$



Basics I

$$S = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

- This is a 3×3 matrix.
- In general $m \times n$.
 - m rows and n columns
 - Square matrix when $m = n$
- Each row or column could represent one object. If rows are objects then columns are features/attributes/components



Basics II

- Identity matrix I

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

- If A is a square matrix, $AI = IA = A$
- I is an example of a **diagonal** matrix.
- If $A = [a_1, \dots, a_m]$ is matrix where a_i are the columns, then
 - A is orthogonal if $a_i \cdot a_j = 0$ for $i \neq j$
 - A is orthonormal if above and $a_i \cdot a_i = 1$



Basics III

- Every vector can be written as a linear combination of some finitely many “special” vectors.
- These are called basis-vectors.

$$S = \begin{bmatrix} 3 \\ 2 \\ 2 \end{bmatrix} = 3 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



Linear independence

- Intuitively, a set of vectors is linearly independent if any element of the set cannot be expressed as a linear combination of the others.
- The columns are not linearly independent:

$$S = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$



Exercise

Given $v = (1, 2, 3)$, calculate:

1. The length of v
2. A unit vector in the same direction as v
3. A set orthogonal bases for v
4. A vector perpendicular to v



Solution

Given $v = (1, 2, 3)$, calculate:

1. The length of v is 3
2. A unit vector in the same direction $(1/\sqrt{14}, 2/\sqrt{14}, 3/\sqrt{14})$
3. A set orthogonal bases $(0, 0, 1); (0, 1, 0); (0, 0, 1)$
4. A vector perpendicular to $v = (-1, -1, 1)$



Determinant

- The determinant of a square matrix is a single number telling us if the matrix is invertible, $\det A^{-1} = 1/(\det A)$
- For 2 by 2 matrix

$$|A| = \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc$$

- For 3 by 3 matrix

$$\begin{aligned} |A| &= \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = a \begin{vmatrix} e & f \\ h & i \end{vmatrix} - b \begin{vmatrix} d & f \\ g & i \end{vmatrix} + c \begin{vmatrix} d & e \\ g & h \end{vmatrix} \\ &= aei + bfg + cdh - ceg - bdi - afh. \end{aligned}$$



Rank of a matrix I

- Given a matrix M , the **rank** of a matrix is the maximum number of linearly independent columns.

- A rank 2 matrix:

$$S = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$



Rank of a matrix II

- Can we automatically determine the rank of a matrix

$$S = \begin{bmatrix} 3 & 1 & 3.5 \\ 2 & 2 & 3 \\ 4 & 2 & 5 \end{bmatrix}$$

- Why is rank important anyways...?



Transpose of a matrix

- The transpose of a matrix A , denoted A^T is another matrix B such that $B(i, j) = A(j, i)$ (just flip the matrix)
- The transpose of S is S^T

$$S = \begin{bmatrix} 3 & 3 & 1 \\ 0 & 2 & 4 \\ 0 & 0 & 0 \end{bmatrix} \quad S^T = \begin{bmatrix} 3 & 0 & 0 \\ 3 & 2 & 0 \\ 1 & 4 & 0 \end{bmatrix}$$



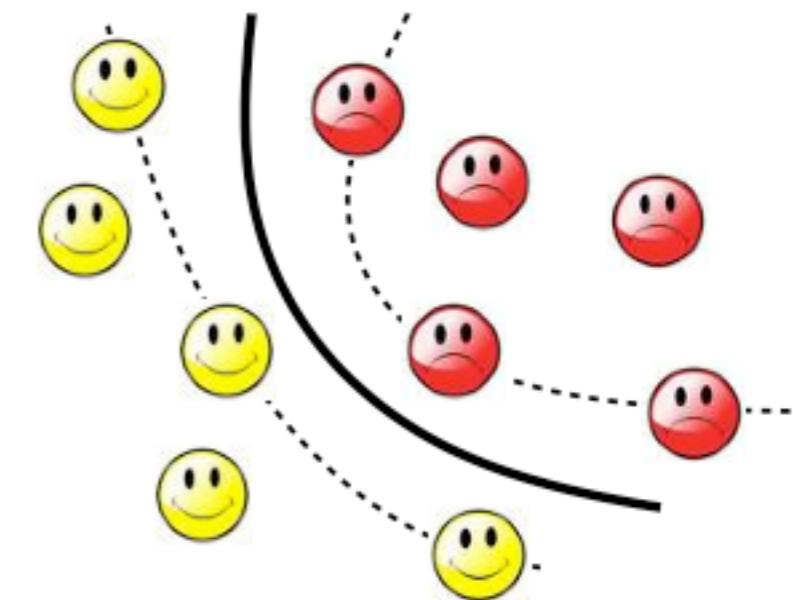
Exercise

- What is the rank of the following matrix:

$$\begin{bmatrix} 1 & 2 & 3 \\ 3 & 4 & 5 \\ 5 & 4 & 3 \\ 0 & 2 & 4 \\ 1 & 3 & 5 \end{bmatrix}$$



THE UNIVERSITY OF
SYDNEY



Matrix Decompositions



Eigen Decomposition

- For any square matrix A we say that λ is an eigenvalue and \mathbf{u} is its eigenvector if

$$A\mathbf{u} = \lambda\mathbf{u}, \quad \mathbf{u} \neq 0.$$

- Stacking up all eigenvectors/values gives

$$AU = U\Lambda = \begin{bmatrix} & & & \\ | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}$$



Eigen Decomposition

- If A is symmetric, all its e-vals are real, and all its e-vecs are orthonormal, $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$
- Hence $U^T U = U \underbrace{U^T}_{n} = I$, $|U| = 1$.
- and $A = U \Lambda U^T = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^T$

$$\begin{aligned} A &= \left[\begin{array}{cccc|c} | & & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_n & | \\ | & | & & & | \end{array} \right] \left[\begin{array}{ccccc} \lambda_1 & & & & \\ & \lambda_2 & & & \\ & & \ddots & & \\ & & & \lambda_n & \end{array} \right] \left[\begin{array}{ccc|c} | & & & \mathbf{u}_1^T \\ \mathbf{u}_1^T & \mathbf{u}_2^T & \dots & | \\ | & & & | \\ \vdots & & & | \\ \mathbf{u}_n^T & & & | \end{array} \right] \\ &= \lambda_1 \left[\begin{array}{c|c} | & \\ \mathbf{u}_1 & | \\ | & \end{array} \right] \left[\begin{array}{ccc|c} | & \mathbf{u}_1^T & & | \end{array} \right] + \dots + \lambda_n \left[\begin{array}{c|c} | & \\ \mathbf{u}_n & | \\ | & \end{array} \right] \left[\begin{array}{ccc|c} | & & \mathbf{u}_n^T & | \end{array} \right] \end{aligned}$$



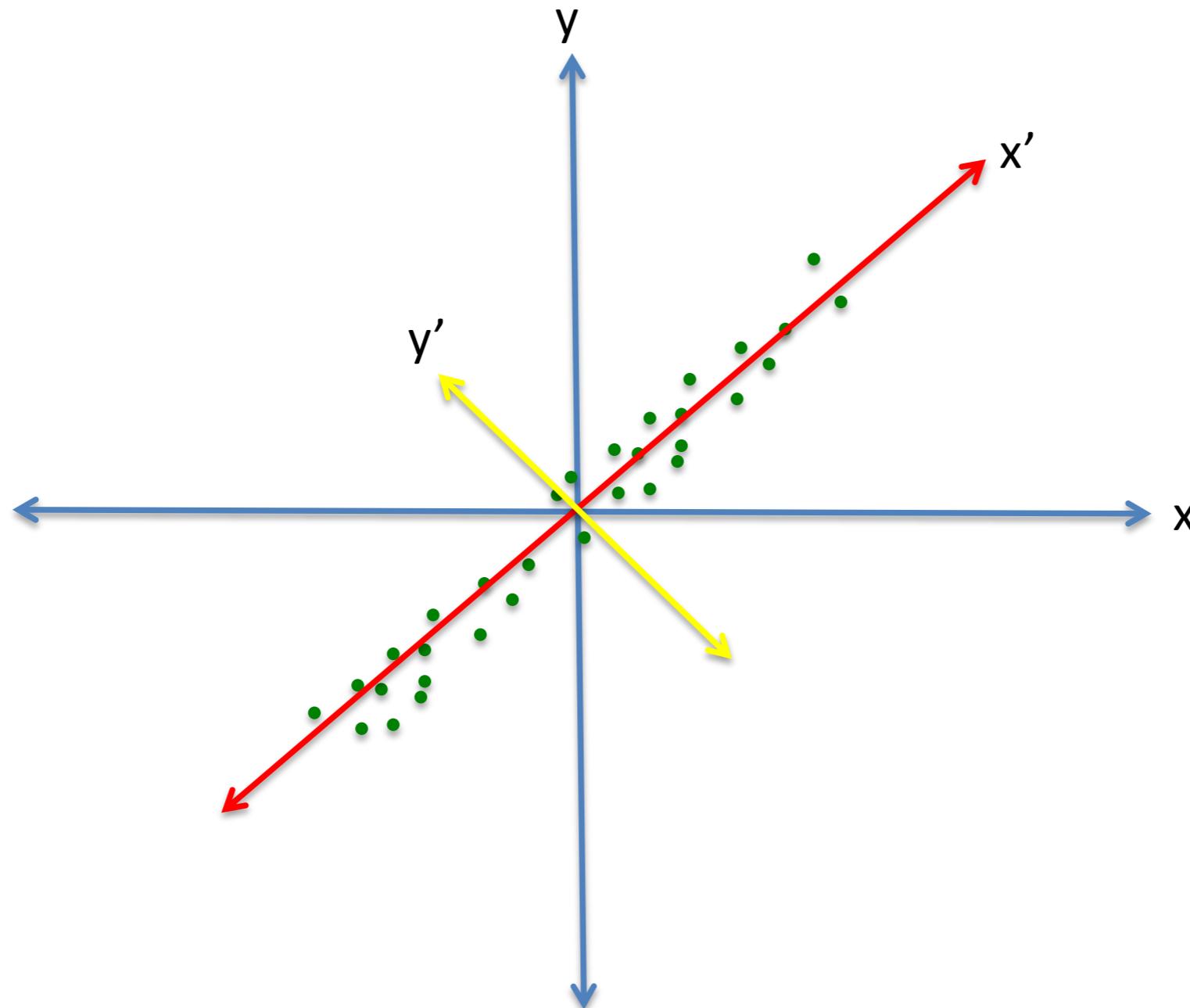
Principal component

- Given a data matrix $X \in \mathbb{R}^{n \times d}$, the principle components of X are the eigenvectors of $X^\top X$
- Principal component analysis (PCA) for X is to find the eigenvectors and eigenvalues of the matrix $X^\top X$.



THE UNIVERSITY OF
SYDNEY

Geometric intuition





Exercise

- Find the eigenvectors and eigenvalues for the following matrix

$$\begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix}$$



Singular Value Decomposition

- Given **any** real matrix X of size (m, n) it can be expressed as:

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

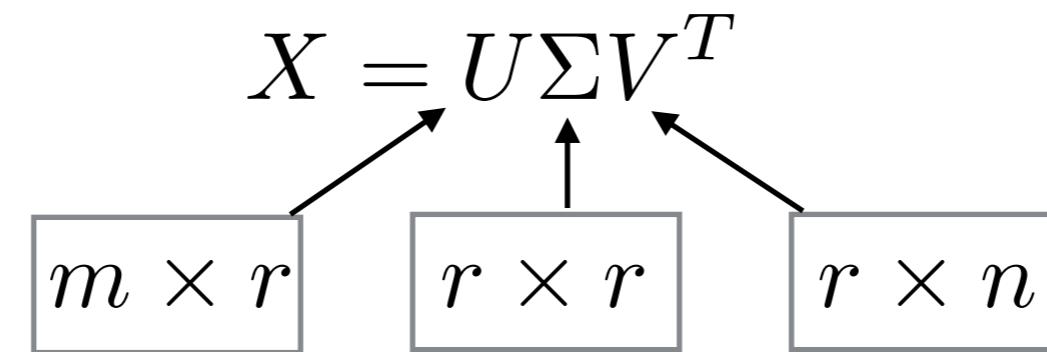
The diagram illustrates the dimensions of the matrices in the SVD formula. It shows three boxes with dimensions: $m \times r$ (left), $r \times r$ (middle), and $r \times n$ (right). Arrows point from each dimension box to its corresponding term in the formula: the $m \times r$ box points to U , the $r \times r$ box points to Σ , and the $r \times n$ box points to V^T .

- r is the rank of matrix X
- U is a (m, r) column-orthonormal matrix
- V is a (n, r) column-orthonormal matrix
- Σ is diagonal $r \times r$ matrix



THE UNIVERSITY OF
SYDNEY

Singular Value Decomposition



$$X = \lambda_1 \mathbf{u}_1 \times \mathbf{v}_1^t + \lambda_2 \mathbf{u}_2 \times \mathbf{v}_2^t + \cdots + \lambda_r \mathbf{u}_r \times \mathbf{v}_r^t$$



SVD in Python

- One line command:

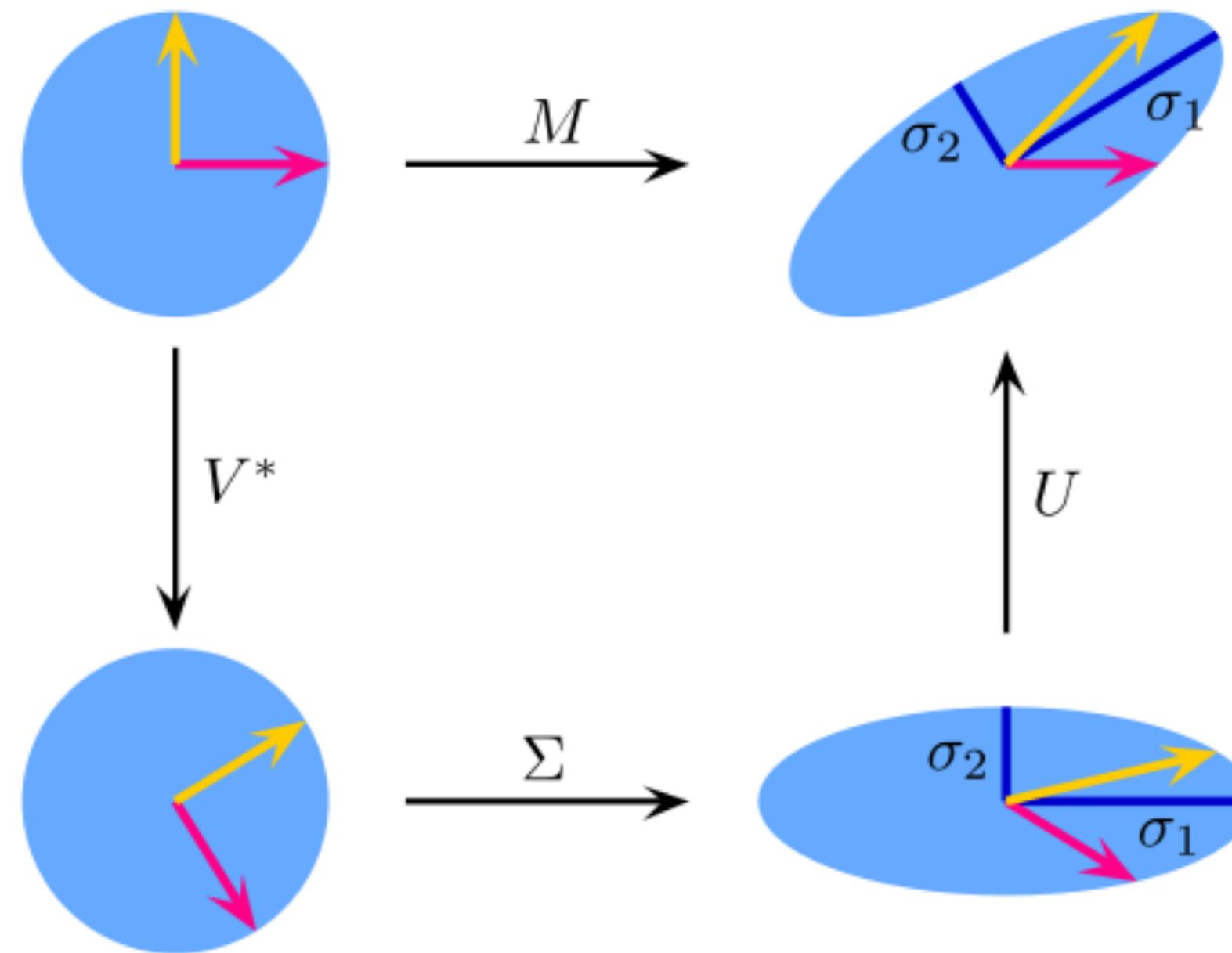
```
>>> U, s, V = np.linalg.svd(a, full_matrices = True)
```

$$X =$$

$$X = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$



SVD as a sequence of operations



$$M = U \cdot \Sigma \cdot V^*$$



Qualitative use of SVD

$$X = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

- Two kinds of customers (businesses and individuals)
- Two kinds of days (weekday and weekends)
- U is a customer-pattern matrix
- V is a day-pattern matrix
- $V(1,2) = 0$ means Wednesday has zero similarity with the “weekend pattern.”

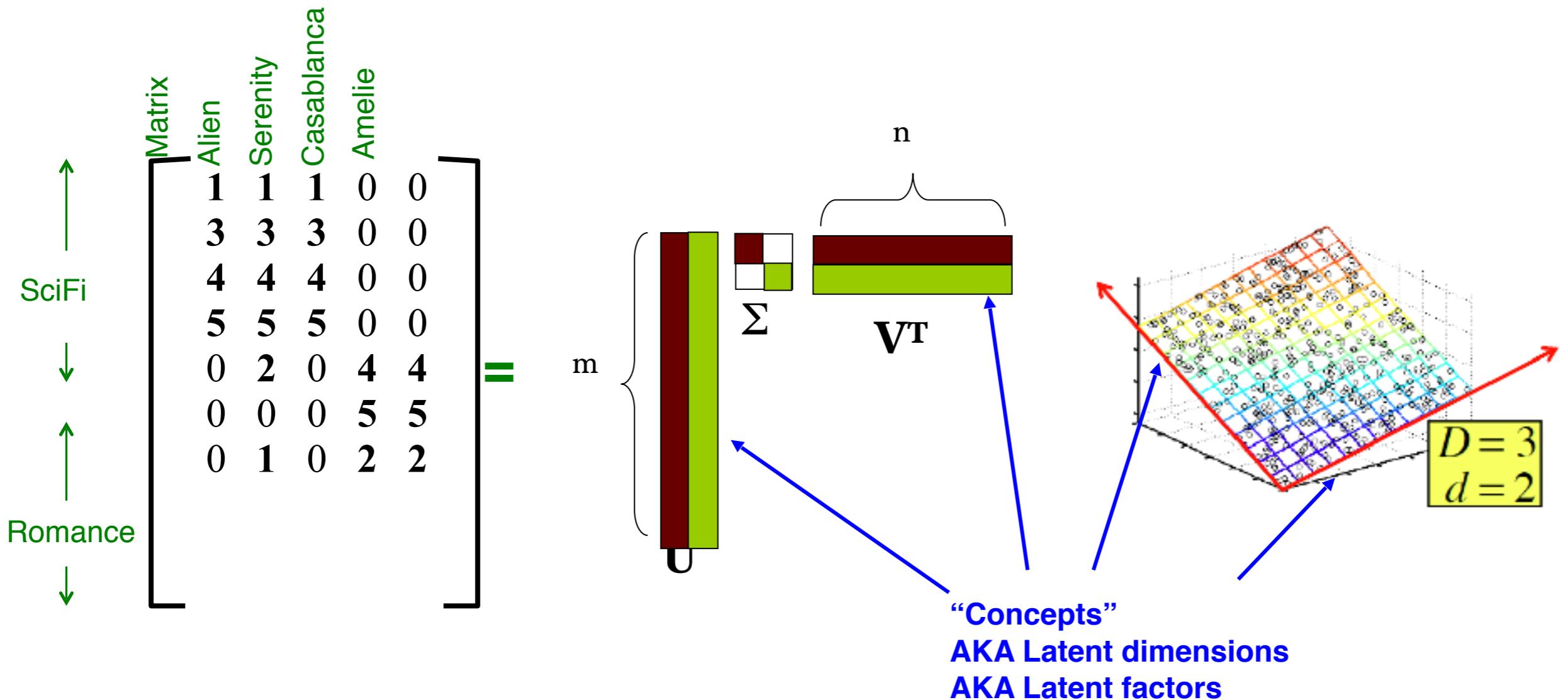


Data reduction with SVD

- Now how can we use an SVD decomposition...
- The first way is **qualitative**...
 - Customer-day pattern; or Customer-Pattern-day matrix...
- The second way is **quantitative**...
 - When X is very large, we can compress X into a smaller matrix but still retain important information...

Example: Users to Movies

- $A = U\Sigma V^T$ - example: Users to Movies





Example: Users to Movies

- $A = U\Sigma V^T$ - example: Users to Movies

U is “user-to-concept” similarity matrix

Matrix

	Alien	Dr Strange	American	Capt	Amelie
1	1	1	0	0	
3	3	3	0	0	
4	4	4	0	0	
5	5	5	0	0	
0	2	0	4	4	
0	0	0	5	5	
0	1	0	2	2	

SciFi

Romance

\equiv

	SciFi-concept	Romance-concept	
0.13	0.02	-0.01	
0.41	0.07	-0.03	
0.55	0.09	-0.04	
0.68	0.11	-0.05	
0.15	-0.59	0.65	
0.07	-0.73	-0.67	
0.07	-0.29	0.32	

“strength” of the SciFi-concept

\times

12.4	0	0	
0	9.5	0	
0	0	1.3	

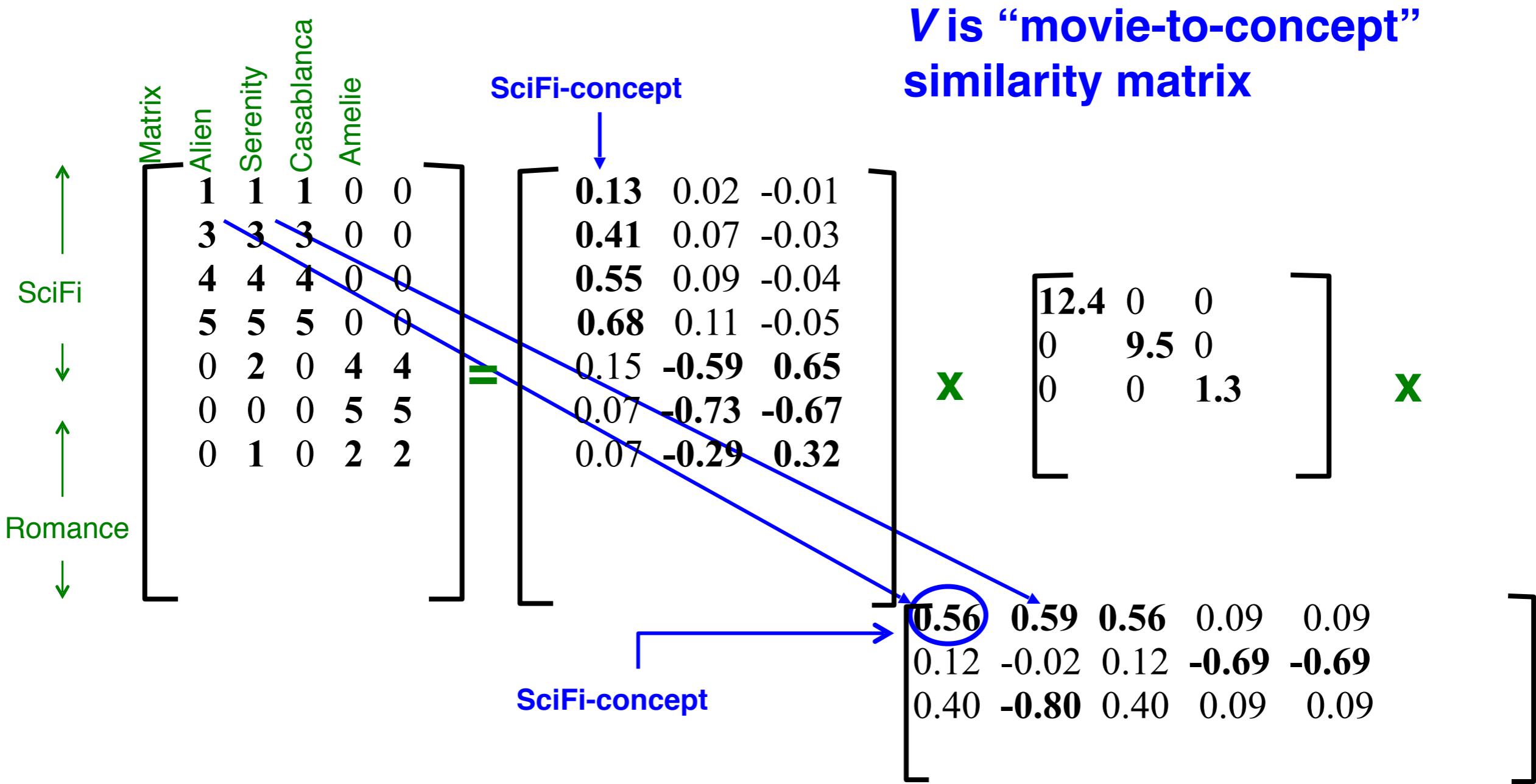
\times

0.56	0.59	0.56	0.09	0.09
0.12	-0.02	0.12	-0.69	-0.69
0.40	-0.80	0.40	0.09	0.09



Example: Users to Movies

- $A = U\Sigma V^T$ - example:

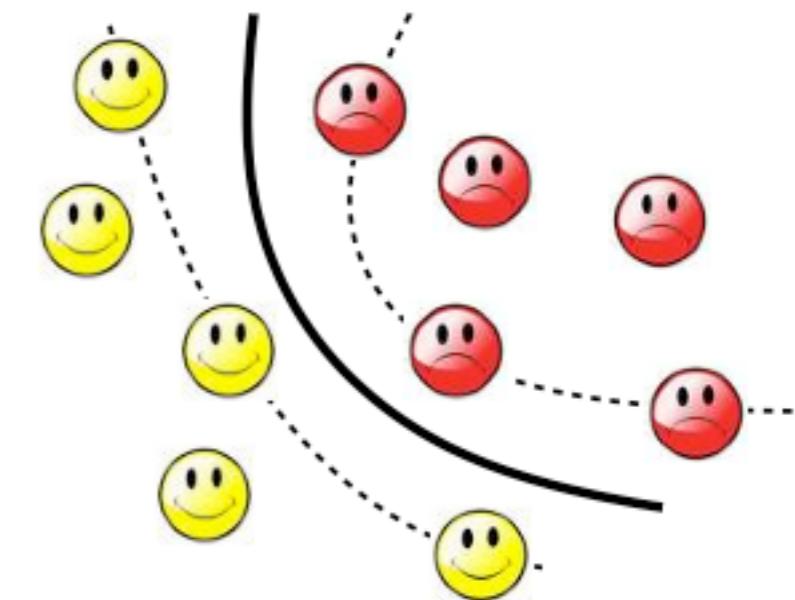




Interpretation

‘movies’, ‘users’ and ‘concepts’:

- U : user-to-concept similarity matrix
- V : movie-to-concept similarity matrix
- Σ : its diagonal elements:
‘strength’ of each concept



Data Reduction and Matrix compression



Advice Querying on X

- Suppose X was very very large... (all mobile users in China)
- Now the query is “Find the amount spent by Jack Ma on Friday”
 - The original answer $X(\text{JackMa}, \text{Friday})$
 - The approximate answer $\hat{X}(\text{JackMa}, \text{Friday})$
 - Is it a good approximation?



Spectral representation

- Matrix X can also be written as:

$$X = \lambda_1 \mathbf{u}_1 \times \mathbf{v}_1^t + \lambda_2 \mathbf{u}_2 \times \mathbf{v}_2^t + \cdots + \lambda_r \mathbf{u}_r \times \mathbf{v}_r^t$$

- The above is called **spectral** representation.

$$X = 9.64 \times \begin{bmatrix} 0.18 \\ 0.36 \\ 0.18 \\ 0.90 \\ 0 \\ 0 \\ 0 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \end{bmatrix} + 5.29 \times \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0.53 \\ 0.80 \\ 0.27 \end{bmatrix} \times \begin{bmatrix} 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$



Example: compression of X

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

The diagram illustrates the dimensions of the matrices in the Singular Value Decomposition (SVD) of matrix X . It shows three boxes with dimensions: $m \times r$, $r \times r$, and $r \times n$. Arrows point from these boxes to the corresponding terms in the equation $X = U\Sigma V^T$.

- Size of $X = mn$
- Size of $U + \Sigma + V$ is $mr + r + nr$
- Thus compression ratio is

$$\frac{mr + r + nr}{mn} = \frac{r(m + 1 + n)}{mn} \approx \frac{rm}{mn} = \frac{r}{n}$$

- Can we do better?



Example: compression of X

- To get better compression we should look at the λ values.
 - These are called **singular values**.
- We can arrange them in descending order:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$$

- Now recall....

$$X = \lambda_1 \mathbf{u}_1 \times \mathbf{v}_1^t + \lambda_2 \mathbf{u}_2 \times \mathbf{v}_2^t + \dots + \lambda_r \mathbf{u}_r \times \mathbf{v}_r^t$$



Example: compression of X

- Now a compact way of writing the spectral representation is:

$$X = \sum_{i=1}^r \lambda_i \mathbf{u}_i \times \mathbf{v}_i^t$$

- However, can approximate it as:

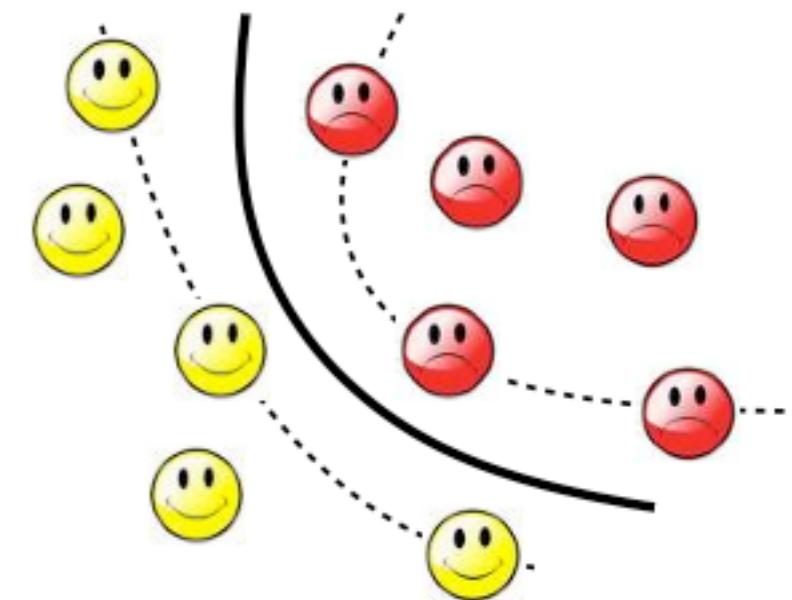
$$\hat{X} = \sum_{i=1}^k \lambda_i \mathbf{u}_i \times \mathbf{v}_i^t$$

- This new compression ratio is:

$$\frac{mk + k + nk}{mn} = \frac{k(m+1+n)}{mn} \approx \frac{km}{mn} = \frac{k}{n} \leqslant \frac{r}{n}$$



THE UNIVERSITY OF
SYDNEY



Thanks!