



THE UNIVERSITY OF
SYDNEY

Machine Learning and Data Mining (COMP 5318)

Latent Variable Models

Fabio Ramos
Roman Marchant



THE UNIVERSITY OF
SYDNEY

Announcements

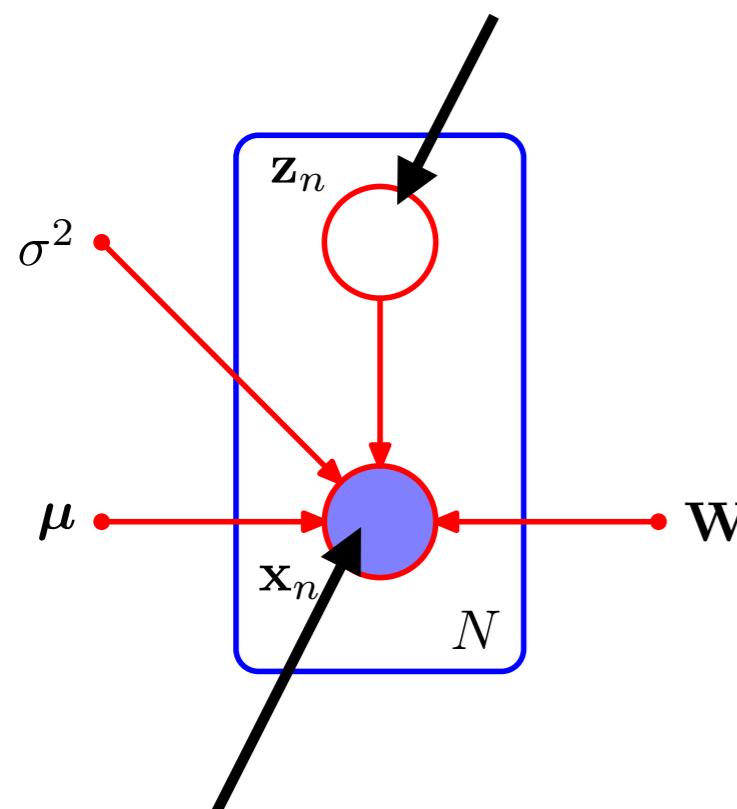
Assignment 2 will be available
at the end of today's lecture.



THE UNIVERSITY OF
SYDNEY

What are latent variable models?

Latent



Essentially: models where some variables are not observed

Example: what is the low dimensional representation for a given dataset?
In this lecture, we will use latent variable models for dimensionality reduction

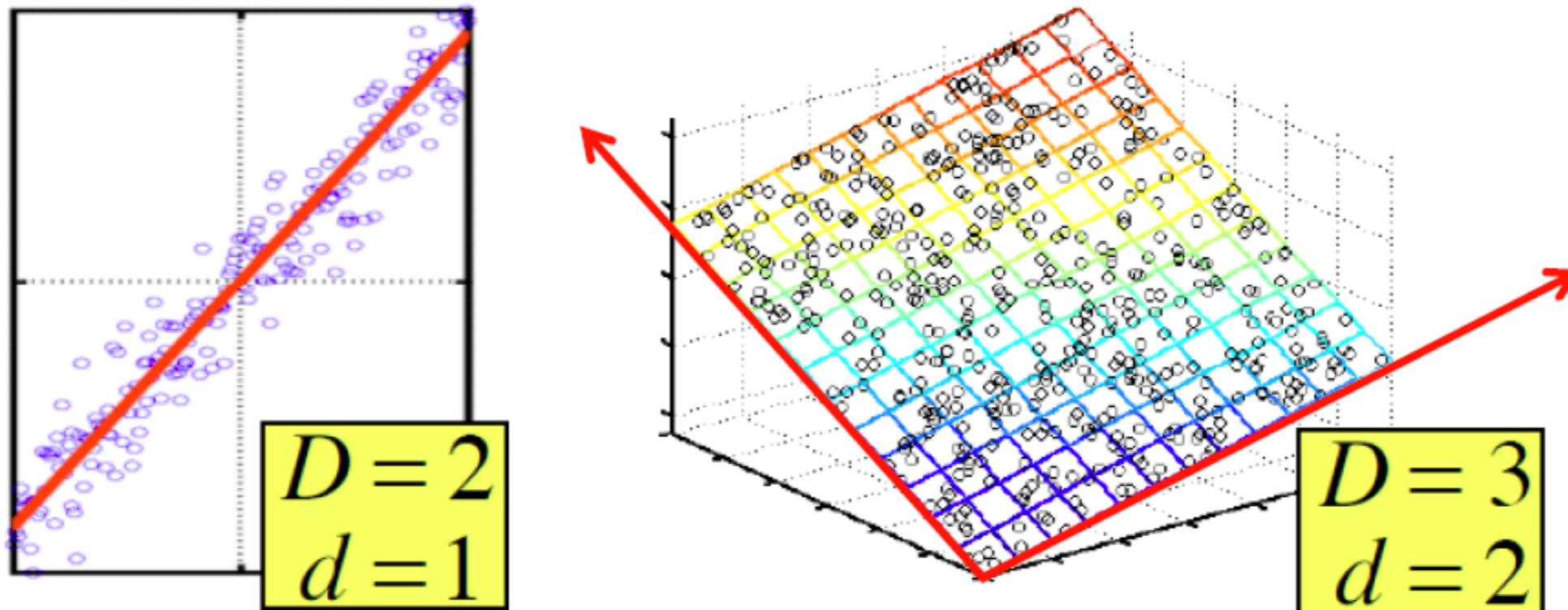
Observed

They can also be used for clustering, recommender systems and others



THE UNIVERSITY OF
SYDNEY

Dimensionality Reduction



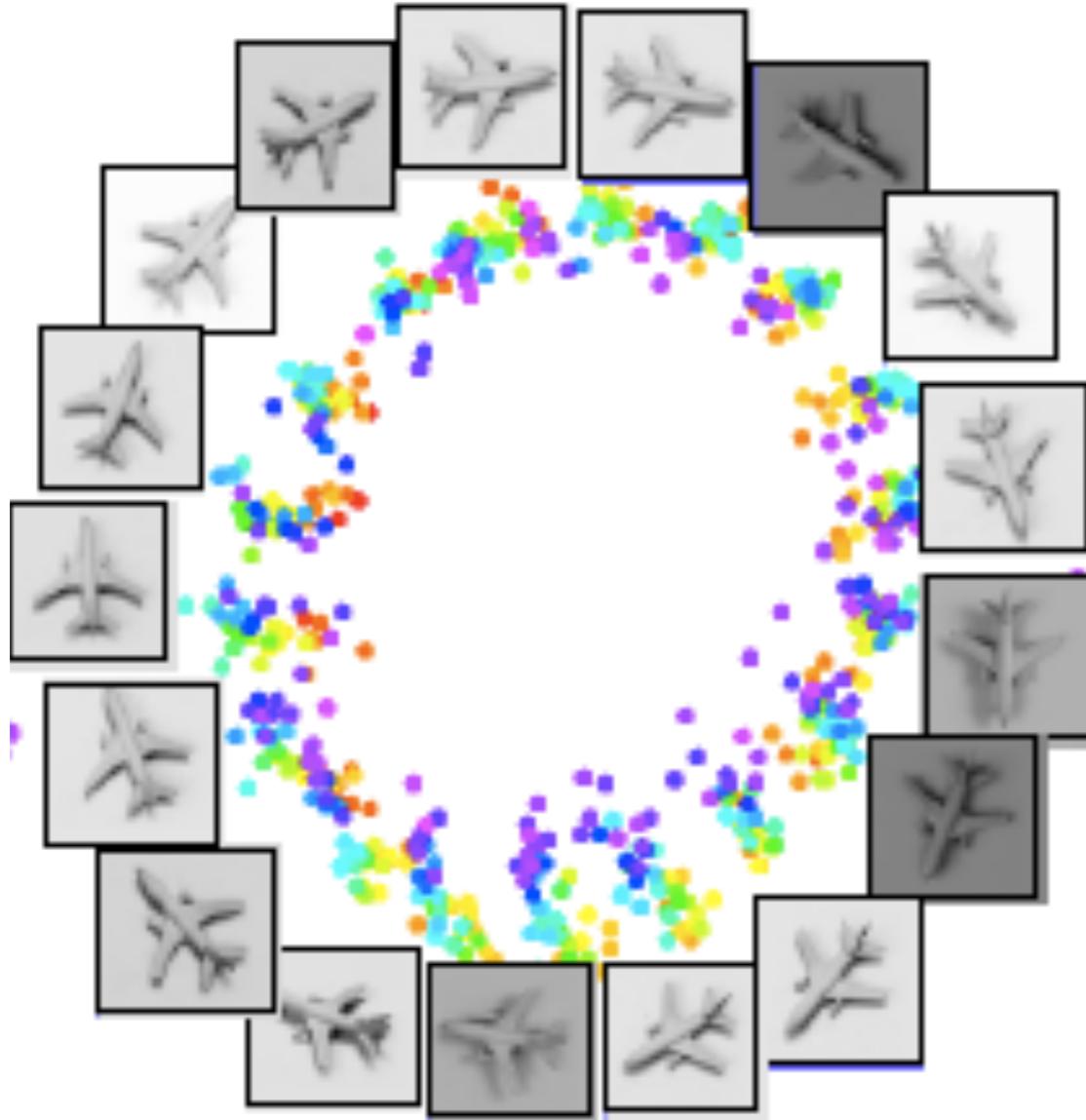
Assumption: Data lies on or near a low d -dimensional subspace

Axes of this subspace are effective representation of the data



THE UNIVERSITY OF
SYDNEY

Dimensionality Reduction

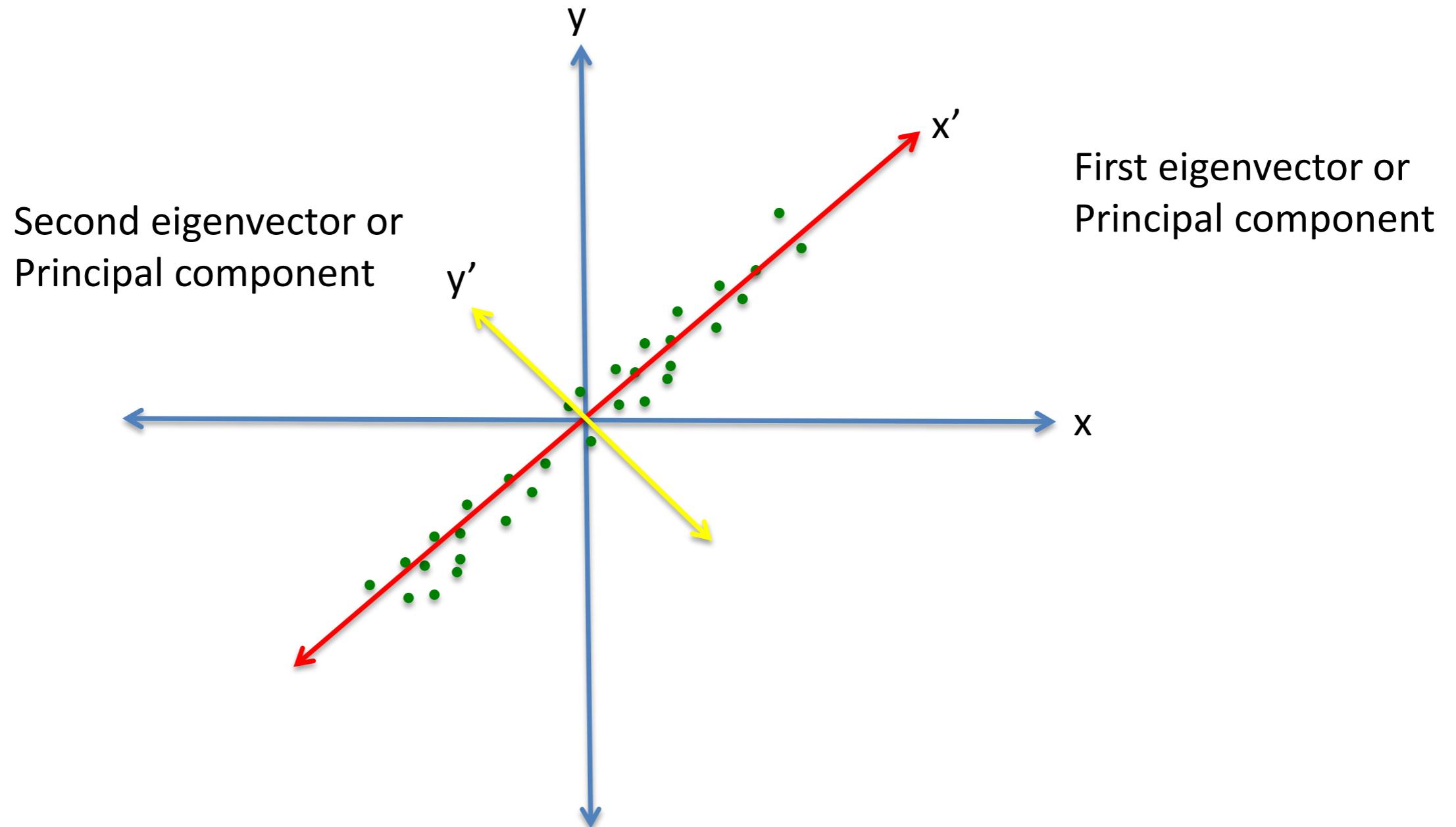


How to find
complicated
data structures?



THE UNIVERSITY OF
SYDNEY

Example: compression of X

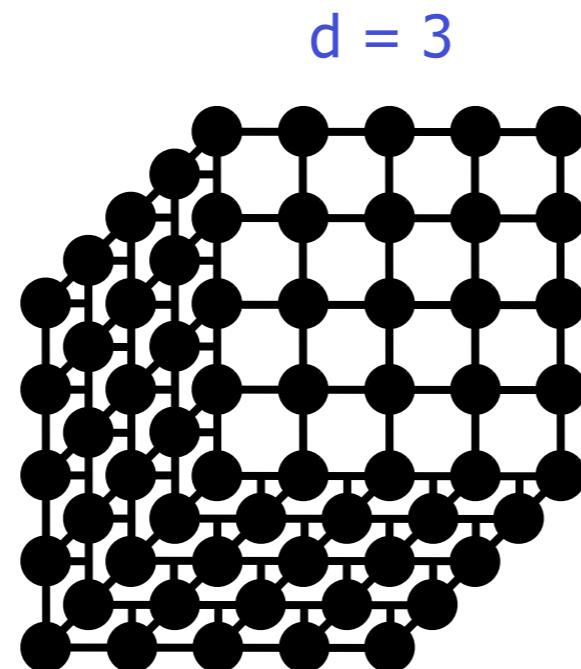
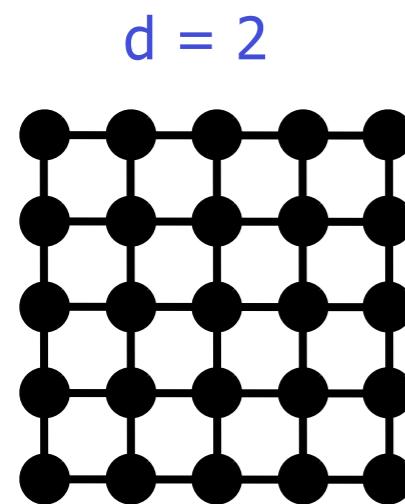
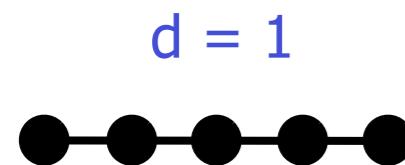




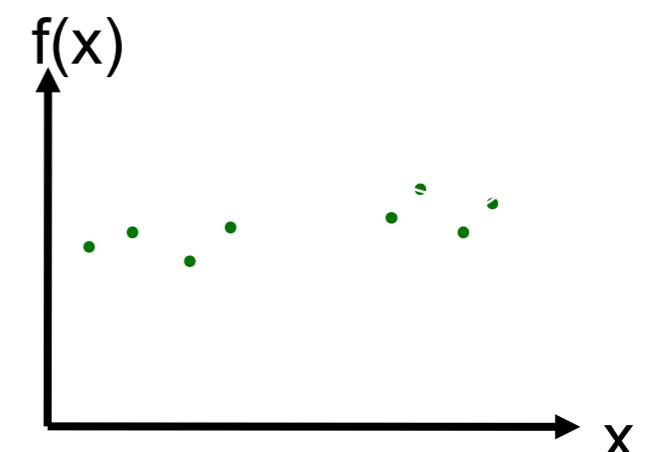
THE UNIVERSITY OF
SYDNEY

The empty space phenomenon

How to “fill” a (compact) space regularly, for a given precision



Number of points on a grid increases **exponentially** with d



In high dimension d :

- never enough data
- never sure to interpolate



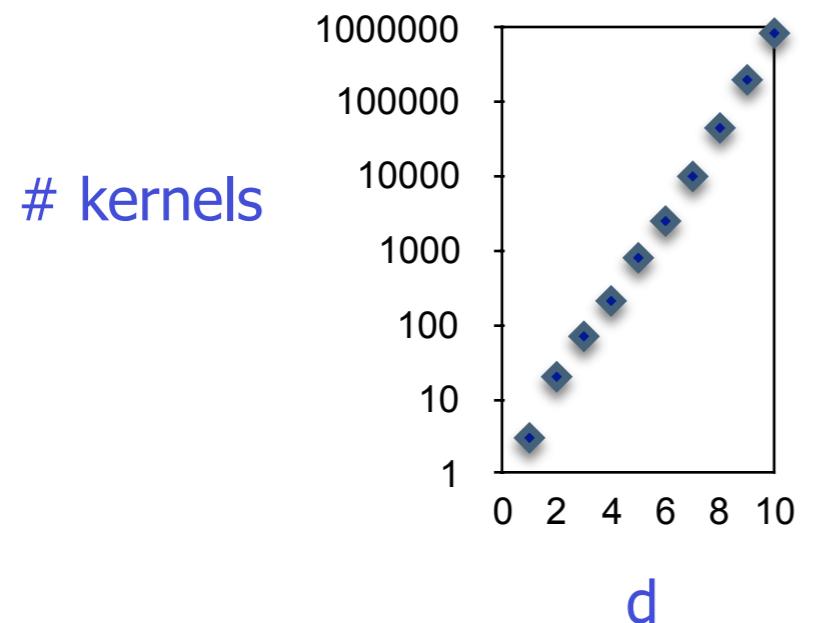
THE UNIVERSITY OF
SYDNEY

The curse of dimensionality

All surprising, or unexpected phenomena, encountered in high-dimensional spaces

Historical example: Silverman (1986)

Number of Gaussian kernels necessary to approximate a (Gaussian) distribution in dimension d , at a given precision



- clear exponential increase
- already high for small d !



THE UNIVERSITY OF
SYDNEY

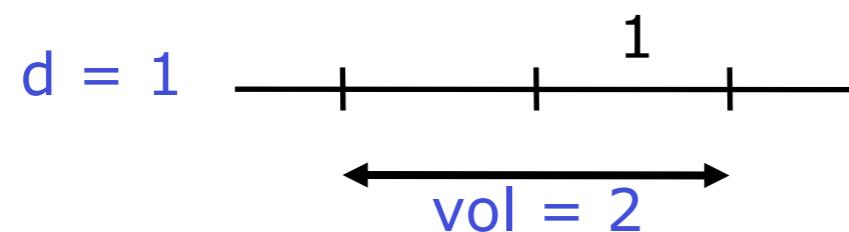
Surprising facts: volume of the sphere

Some facts that are intuitive in low d are less intuitive in high d !

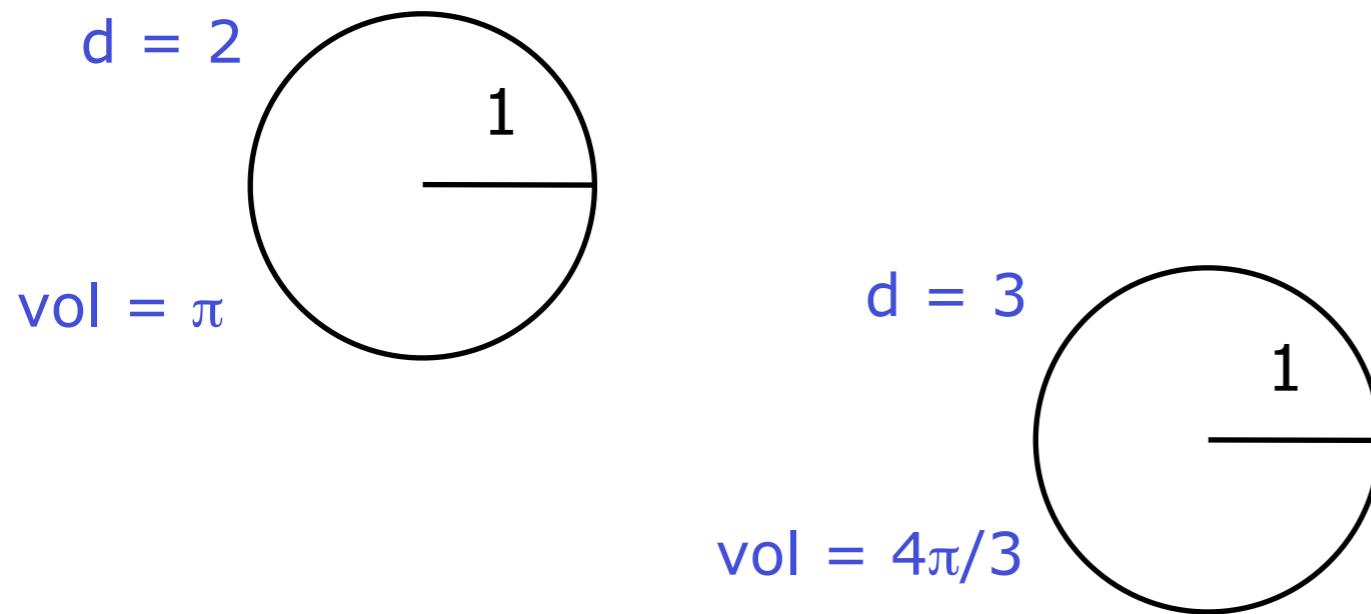
Example: volume of the sphere

Constant radius ($=1$)

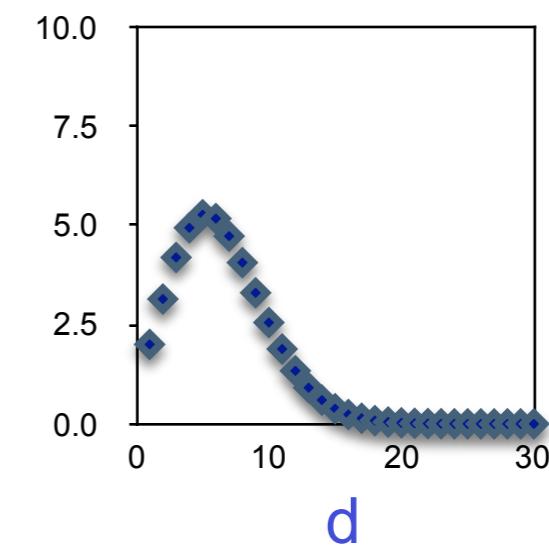
Increasing space dimension



$$v(d) = \frac{\pi^{d/2}}{\Gamma(d/2 + 1)} r^d$$



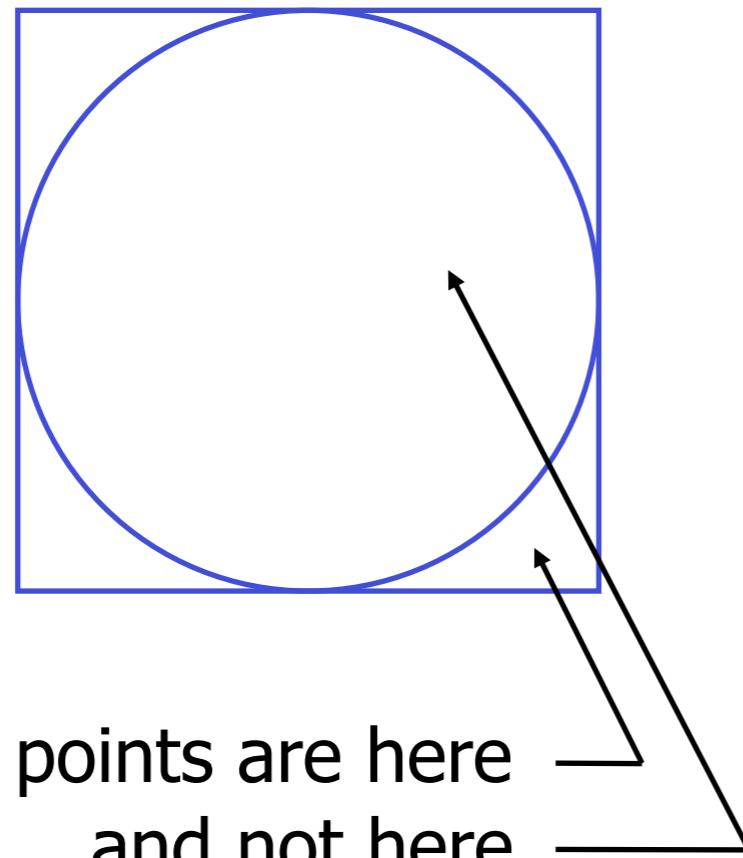
volume



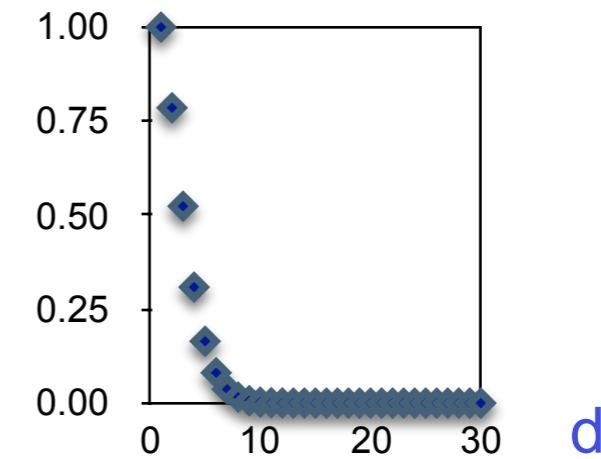


THE UNIVERSITY OF
SYDNEY

Ratio volume sphere / cube



ratio

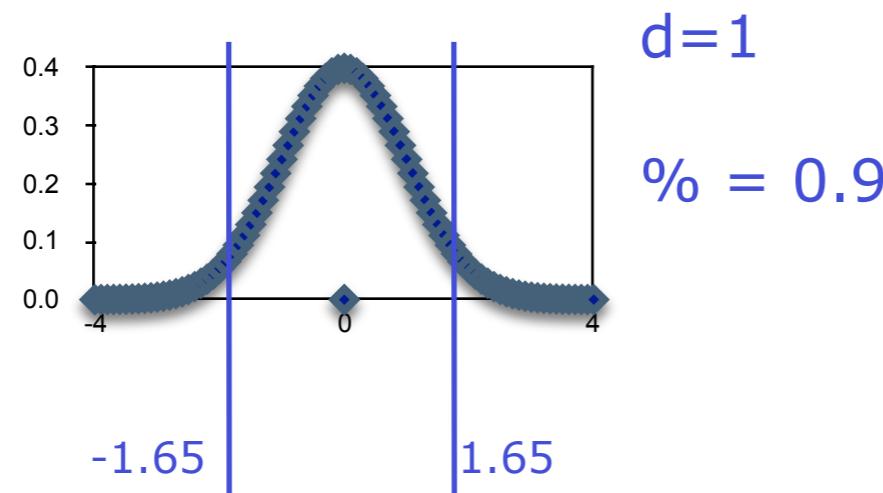




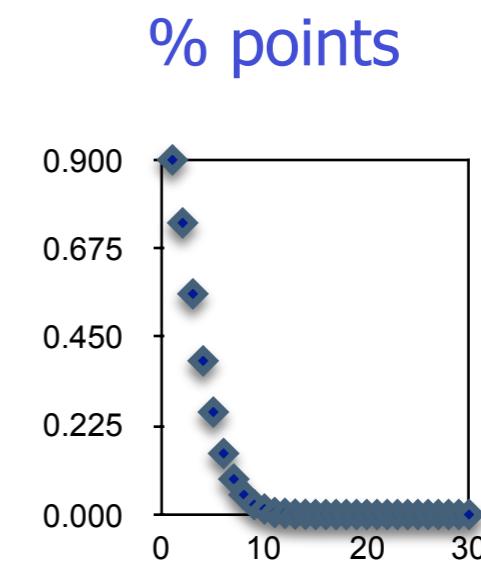
THE UNIVERSITY OF
SYDNEY

High-dimensional Gaussians

% points whose distance to centre is < 1.65



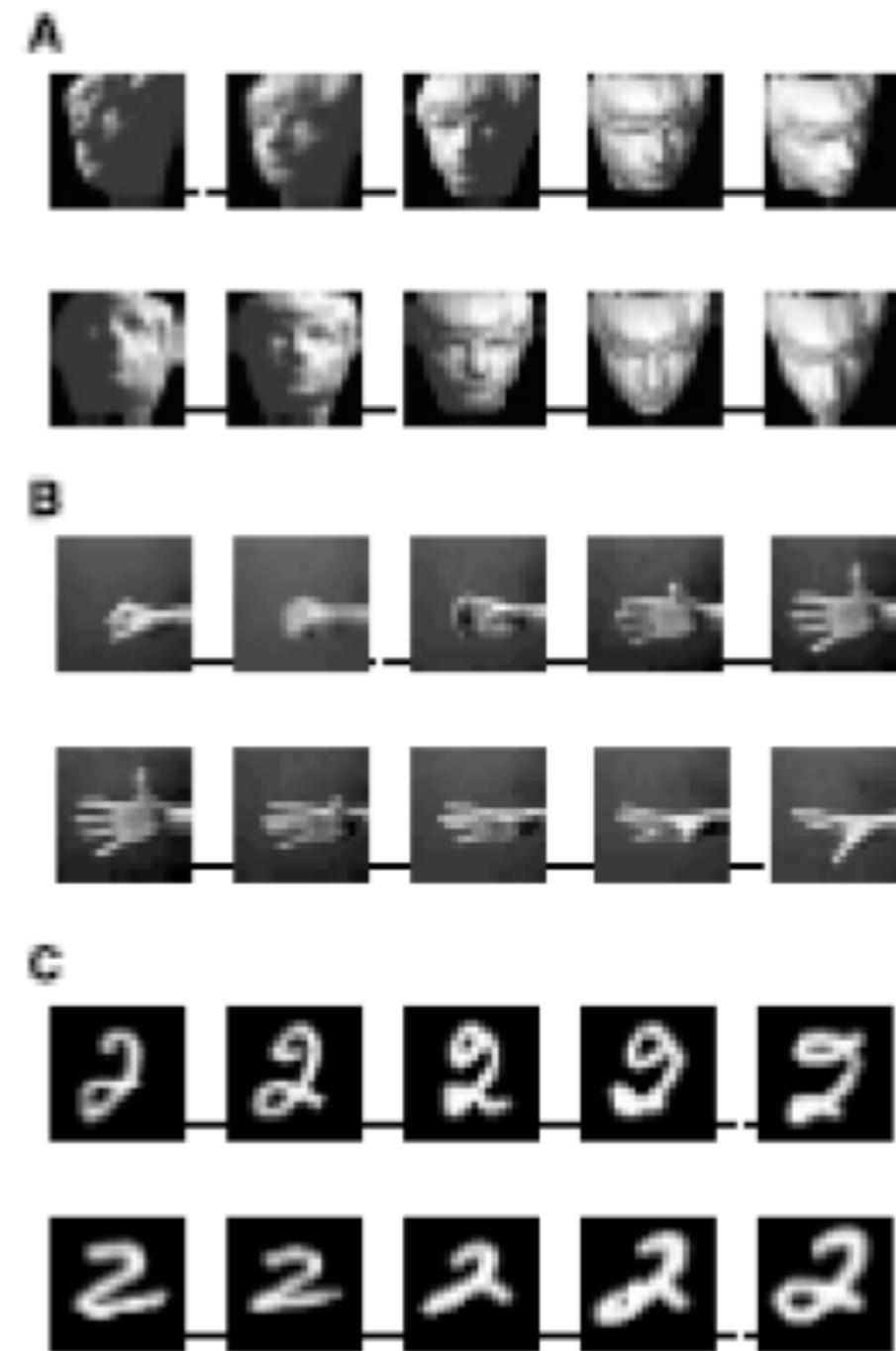
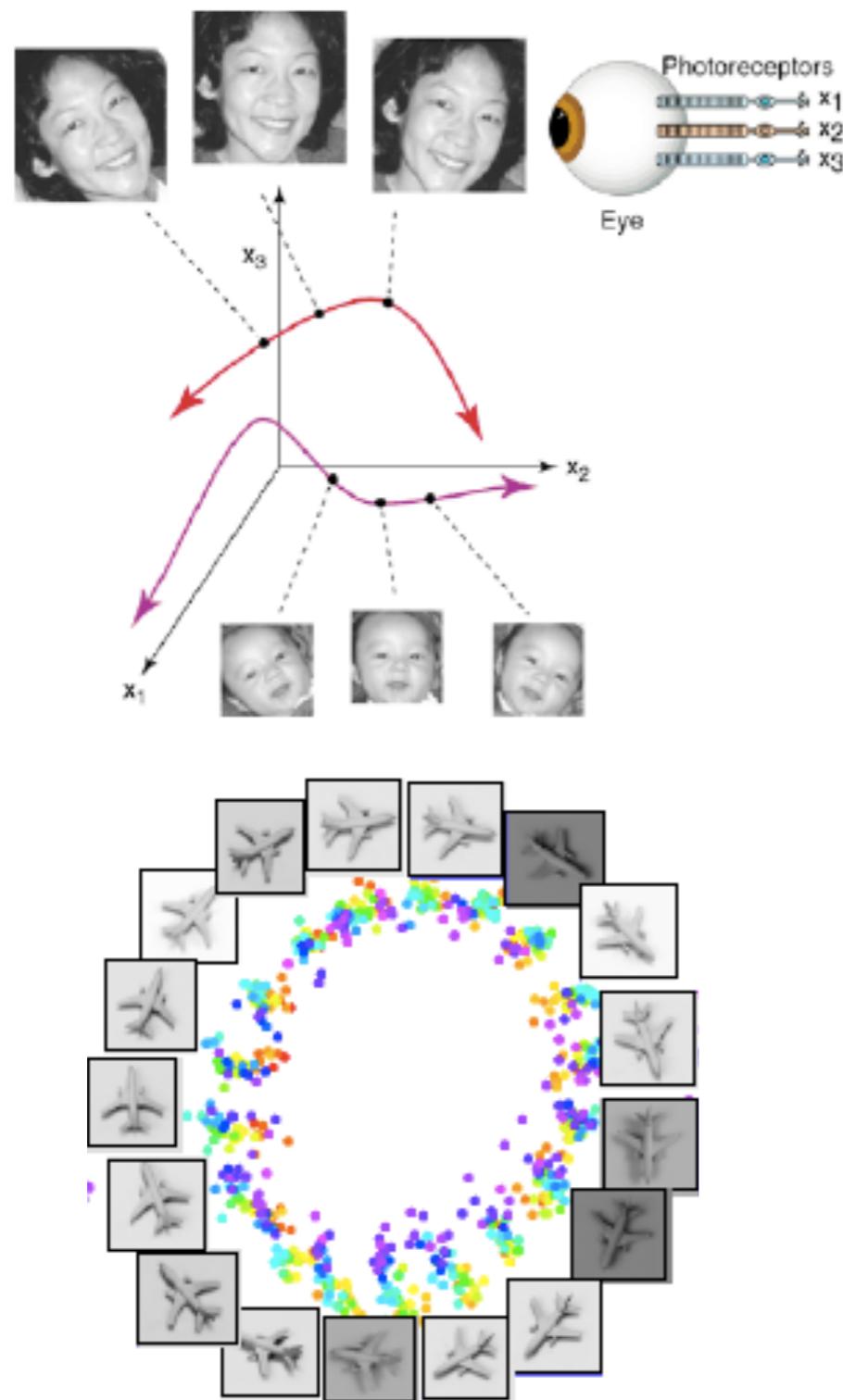
The Gaussian is not local anymore!





THE UNIVERSITY OF
SYDNEY

But there is hope!





THE UNIVERSITY OF
SYDNEY

Principal Component Analysis

C. Bishop, *Pattern Recognition and Machine Learning*, Chapter 12: Continuous Latent Variables
Springer New York, 2006



PCA

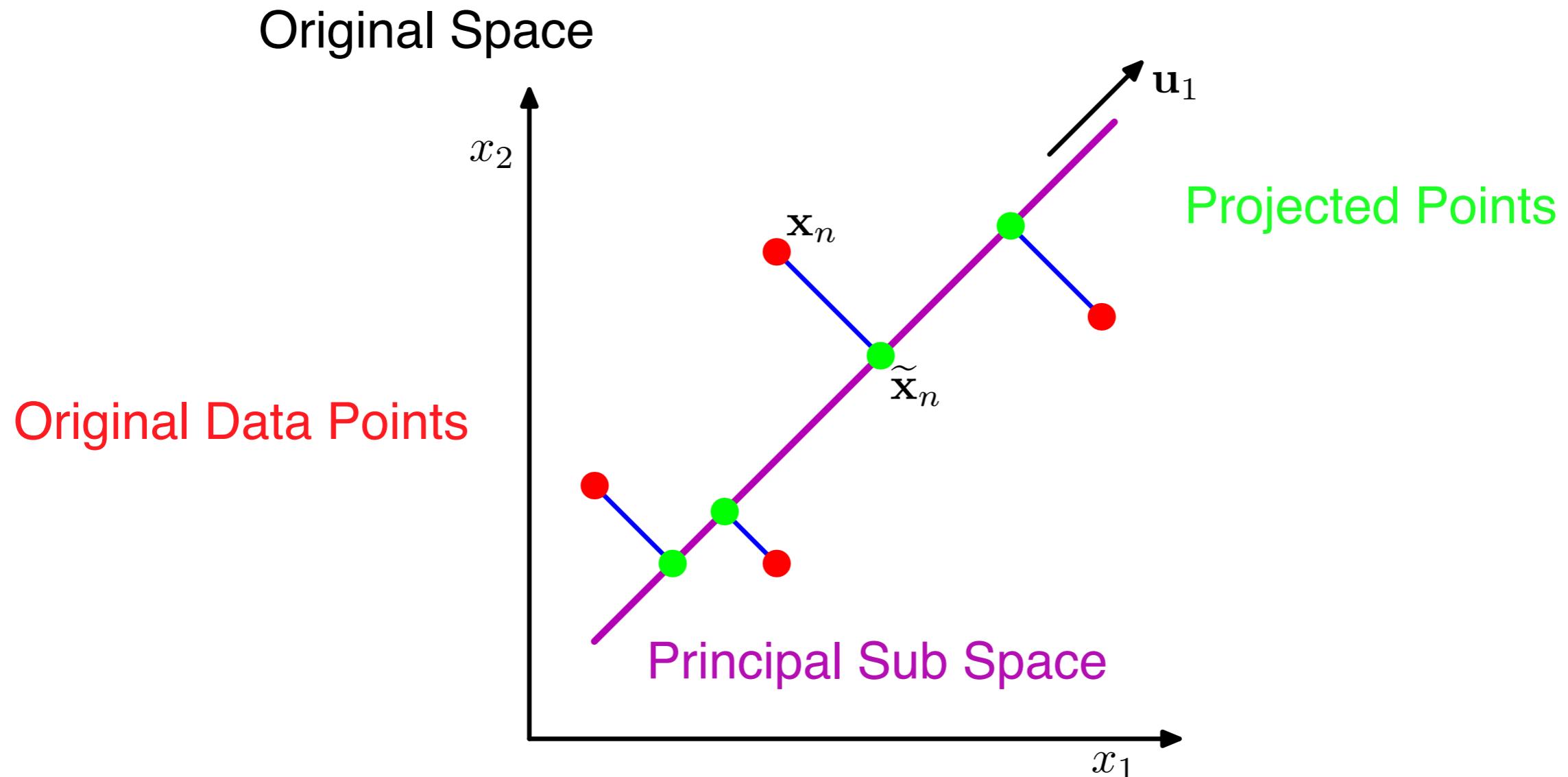
THE UNIVERSITY OF
SYDNEY

*The central idea of principal component analysis (PCA) is to reduce the dimensionality of a data set consisting of a large number of interrelated variables, while retaining as much as possible of the variation present in the data set. This is achieved by transforming to a new set of variables, the principal components (PCs), which are uncorrelated, and which are ordered so that the first few retain most of the variation present in all of the original variables.
[Jolliffe, Principal Component Analysis, 2nd edition]*



THE UNIVERSITY OF
SYDNEY

Geometric intuition





THE UNIVERSITY OF
SYDNEY

PCA Formulation

Let $\{\mathbf{x}_n\}_{n=1,\dots,N}$ be a dataset of observations with dimensionality D

Goal, to represent the data in a lower dimensional space M (known).

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

$$X = [(\mathbf{x}_1 - \bar{\mathbf{x}}) \ (\mathbf{x}_2 - \bar{\mathbf{x}}) \ (\mathbf{x}_3 - \bar{\mathbf{x}}) \ \dots]$$

Translating the coordinate system to the location of the mean.



$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^T$$

$$= \frac{1}{N} [(\mathbf{x}_1 - \bar{\mathbf{x}}) \ (\mathbf{x}_2 - \bar{\mathbf{x}}) \ (\mathbf{x}_3 - \bar{\mathbf{x}}) \ \dots] \begin{bmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})^T \\ (\mathbf{x}_2 - \bar{\mathbf{x}})^T \\ (\mathbf{x}_3 - \bar{\mathbf{x}})^T \\ \dots \end{bmatrix}$$

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

S is the covariance matrix. (scatter matrix)

S is the square and symmetric

S can be very large (in vision, N is the number of pixels)



PCA Derivation

Let us assume, for now, that $M = 1$

\mathbf{u}_1 is the direction of this sub-space

\mathbf{u}_1 is a D dimensional vector.

\mathbf{u}_1 is a unit vector, i.e. $\mathbf{u}_1^T \mathbf{u}_1 = 1$

If we project the data point \mathbf{x}_n into a scalar value $\mathbf{u}_1^T \mathbf{x}_n$.

The mean of the projected data becomes

$$\mathbf{u}_1^T \bar{\mathbf{x}} = \frac{\mathbf{u}_1^T}{N} \sum_{n=1}^N \mathbf{x}_n$$

The variance becomes

$$\frac{1}{N} \sum_{n=1}^N \{ \mathbf{u}_1^T \mathbf{x}_n - \mathbf{u}_1^T \bar{\mathbf{x}} \}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$



PCA Derivation

Maximise the projected variance $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$ with respect to \mathbf{u}_1

* problem: $\|\mathbf{u}_1\| \rightarrow \infty$ remember:

\mathbf{u}_1 is a unit vector, i.e. $\mathbf{u}_1^T \mathbf{u}_1 = 1$

Introduce Lagrange multiplier λ_1

$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1)$$

Derivative equal to zero: \mathbf{u}_1 Must be an eigen vector of \mathbf{S}

Maximum variance when \mathbf{u}_1 is the eigenvector with the largest eigenvalue λ_1



PCA Theorem

Theorem

Each \mathbf{x}_n can be written as:
$$\mathbf{x}_n = \bar{\mathbf{x}} + \sum_{i=1}^M z_{ni} \mathbf{u}_i$$

Where \mathbf{u}_i are the M eigenvectors of \mathbf{S} with non-zero eigenvalues.

The scalars z_{ni} are the coordinates of \mathbf{x}_n in the principal component space.



PCA to compress data

Have we achieved compression of the data?

If we consider $M = D$ then there is no compression,

This case results in a realignment of the data w.r.t. principal components.

However! Correlated data will result in z_{ni} being zero or close to zero.

i.e. $\{\mathbf{x}_n\}_{n=1,\dots,N}$ lies in a lower-dimensional subspace.



PCA to compress data

Achieve compression: $M < D$

1. Sort the eigenvectors \mathbf{u}_i according to their associated eigenvalue.

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M$$

2. Discard smaller eigen values, i.e. assuming $\lambda_i \approx 0 \quad \forall i > k$

Then

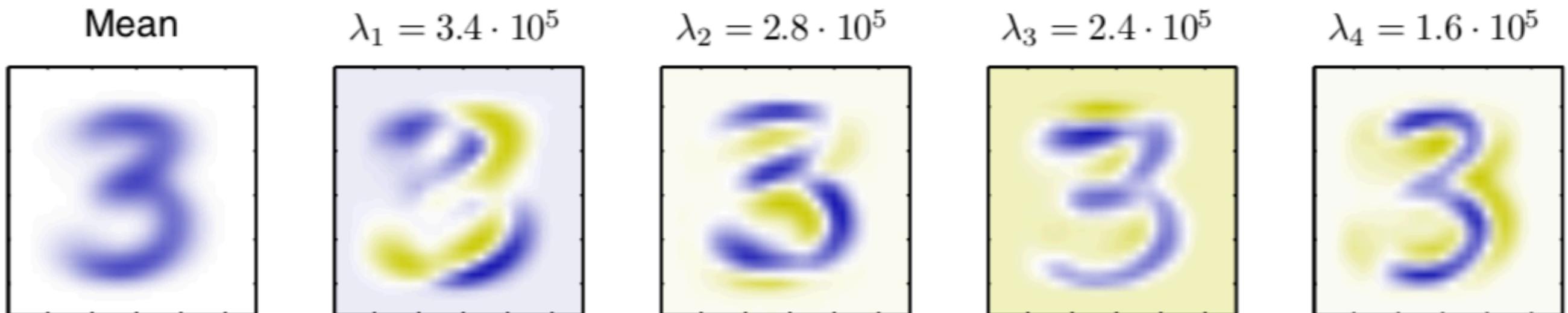
$$\mathbf{x}_n \approx \bar{\mathbf{x}} + \sum_{i=1}^K z_{ni} \mathbf{u}_i$$



THE UNIVERSITY OF
SYDNEY

Reconstruction of digits

An eigen vector of the covariance matrix is a vector in the original D dimensional space.





Reconstruction of digits

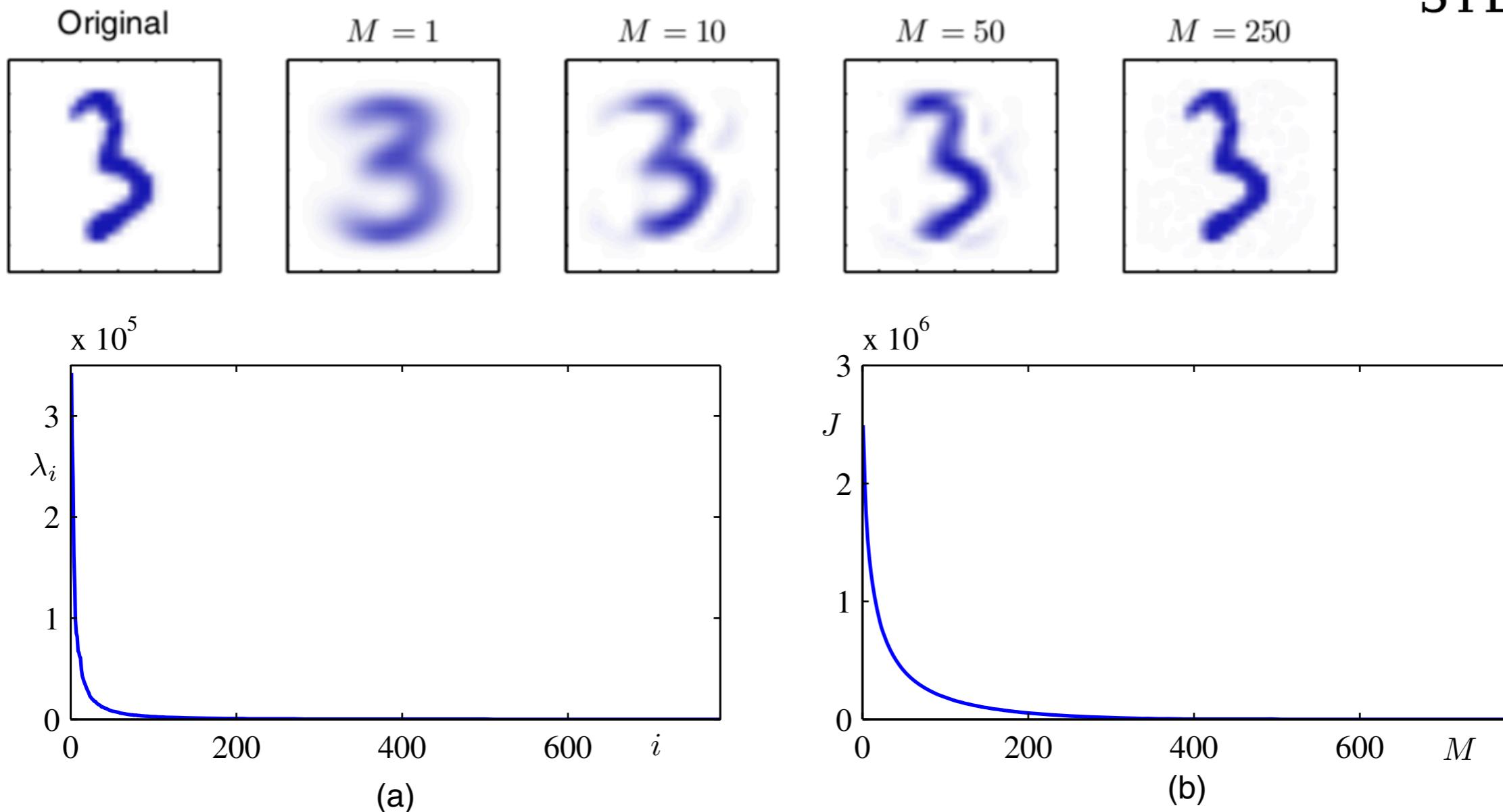


Figure 12.4 (a) Plot of the eigenvalue spectrum for the off-line digits data set. (b) Plot of the sum of the discarded eigenvalues, which represents the sum-of-squares distortion J introduced by projecting the data onto a principal component subspace of dimensionality M .



PCA Algorithm

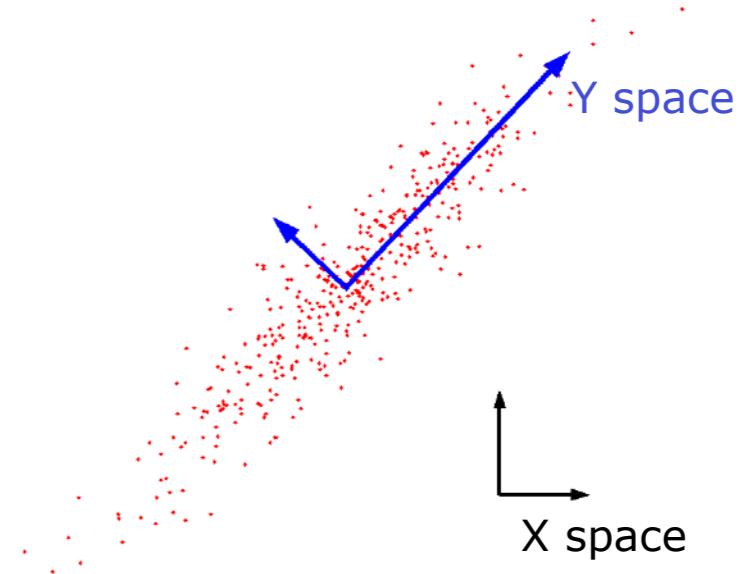
Goal:

To project linearly while keeping the variance of the data

Covariance matrix \mathbf{S} of the data

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

Low-dimensional space (Y)



Eigenvectors and eigenvalues of \mathbf{S}

\mathbf{u}_i = main directions

λ_1 = variance along each direction

Projection & Reconstruction

$$\mathbf{x}_n \approx \bar{\mathbf{x}} + \sum_{i=1}^K z_{ni} \mathbf{u}_i$$

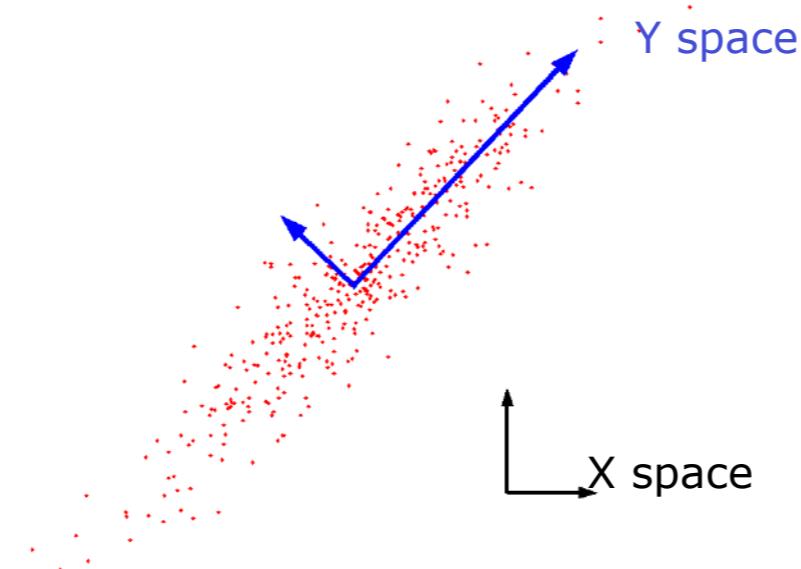
High-dimensional space (X)



THE UNIVERSITY OF
SYDNEY

Metric multidimensional scaling

Is it necessary to know the X-axes in order to find the Y-axes ones?



Of course not!

It should be possible to find the principal directions (the Y axes) only from the **relative positions** between data.

That is metric MDS...



THE UNIVERSITY OF
SYDNEY

Metric multidimensional scaling

Suppose all we were given were distances or symmetric “dissimilarities” Δ_{ij} .

$$\Delta = \begin{bmatrix} 0 & \Delta_{12} & \Delta_{13} & \Delta_{14} \\ \Delta_{12} & 0 & \Delta_{23} & \Delta_{24} \\ \Delta_{13} & \Delta_{23} & 0 & \Delta_{34} \\ \Delta_{14} & \Delta_{24} & \Delta_{34} & 0 \end{bmatrix}$$

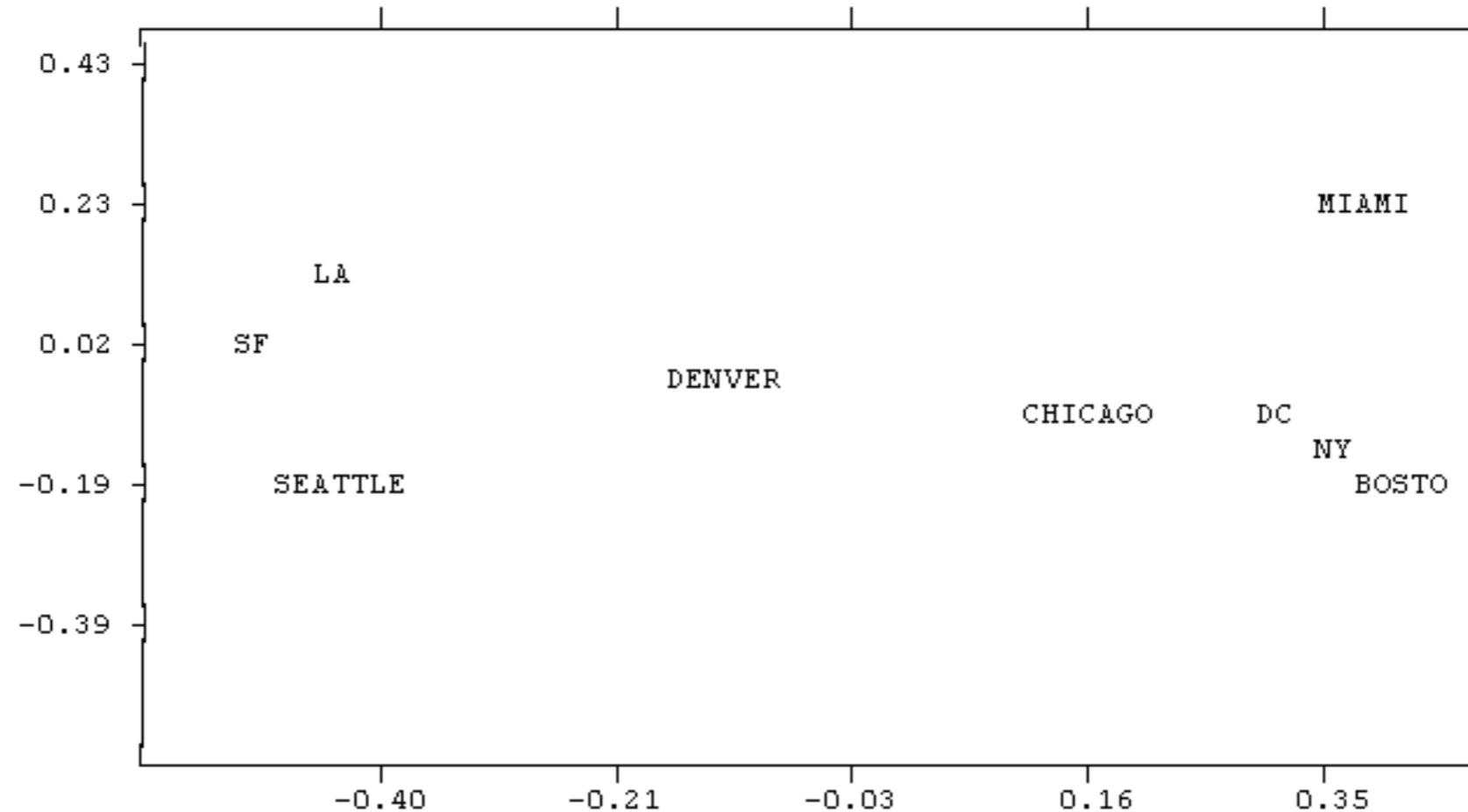
Goal: Find vectors \mathbf{y}_i such that $\|\mathbf{y}_i - \mathbf{y}_j\| \approx \Delta_{ij}$.

This is called **Multidimensional Scaling (MDS)**.



Example: From distances to a map

	1	2	3	4	5	6	7	8	9	
	BOST	NY	DC	MIAM	CHIC	SEAT	SF	LA	DENV	
	-----	-----	-----	-----	-----	-----	-----	-----	-----	
1	BOSTON	0	206	429	1504	963	2976	3095	2979	1949
2	NY	206	0	233	1308	802	2815	2934	2786	1771
3	DC	429	233	0	1075	671	2684	2799	2631	1616
4	MIAMI	1504	1308	1075	0	1329	3273	3053	2687	2037
5	CHICAGO	963	802	671	1329	0	2013	2142	2054	996
6	SEATTLE	2976	2815	2684	3273	2013	0	808	1131	1307
7	SF	3095	2934	2799	3053	2142	808	0	379	1235
8	LA	2979	2786	2631	2687	2054	1131	379	0	1059
9	DENVER	1949	1771	1616	2037	996	1307	1235	1059	0





Metric MDS - Algorithm

Goal:

To project linearly while keeping the $(N-1) \times N/2$ pairwise distances

Computation:

1. Matrix D of the squared distances

$$D = [d_{i,j}^2]$$

2. Eigenvectors and eigenvalues of D

V_i = coordinates along the main directions

λ_i = variance along each direction

3. Projection

$$Y = \sqrt{\text{diag}(\lambda_{1 \leq i \leq p})} V_{1 \leq i \leq p}^T$$

Result of PCA = result of metric MDS !

Only distances are needed i.e. more independent from representation.



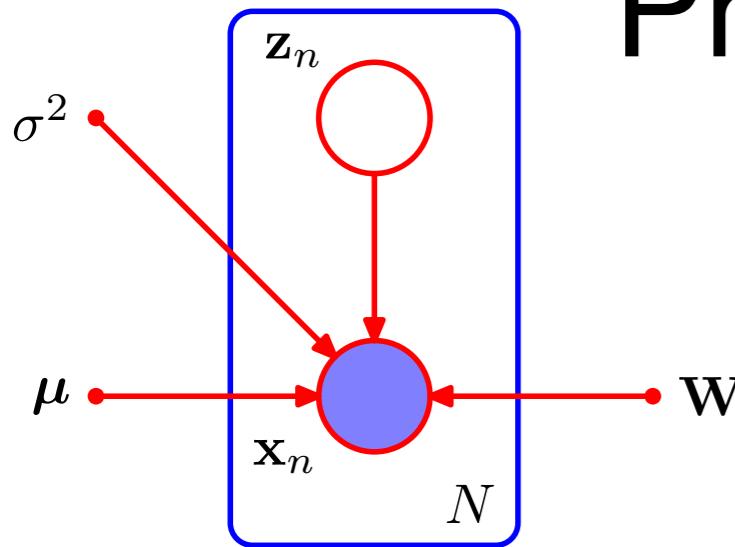
THE UNIVERSITY OF
SYDNEY

Probabilistic PCA

- Probabilistic version assuming Gaussian data
- Solution can be found as a maximum likelihood optimisation
- Computationally efficient
- Generative model: it can be sampled from to generate examples in the low dimensional space.



Probabilistic PCA



Prior for z (latent variable)

$$p(z) = \mathcal{N}(z|0, I).$$

Conditional distribution for x

$$p(x|z) = \mathcal{N}(x|Wz + \mu, \sigma^2 I)$$

x and z are linearly related

$$x = Wz + \mu + \epsilon$$

D-dimensional noise variable

Observed x → D-dimensional

$$\epsilon \sim \mathcal{N}(0, \sigma^2 I)$$

Latent z → M-dimensional



Generative view of PPCA

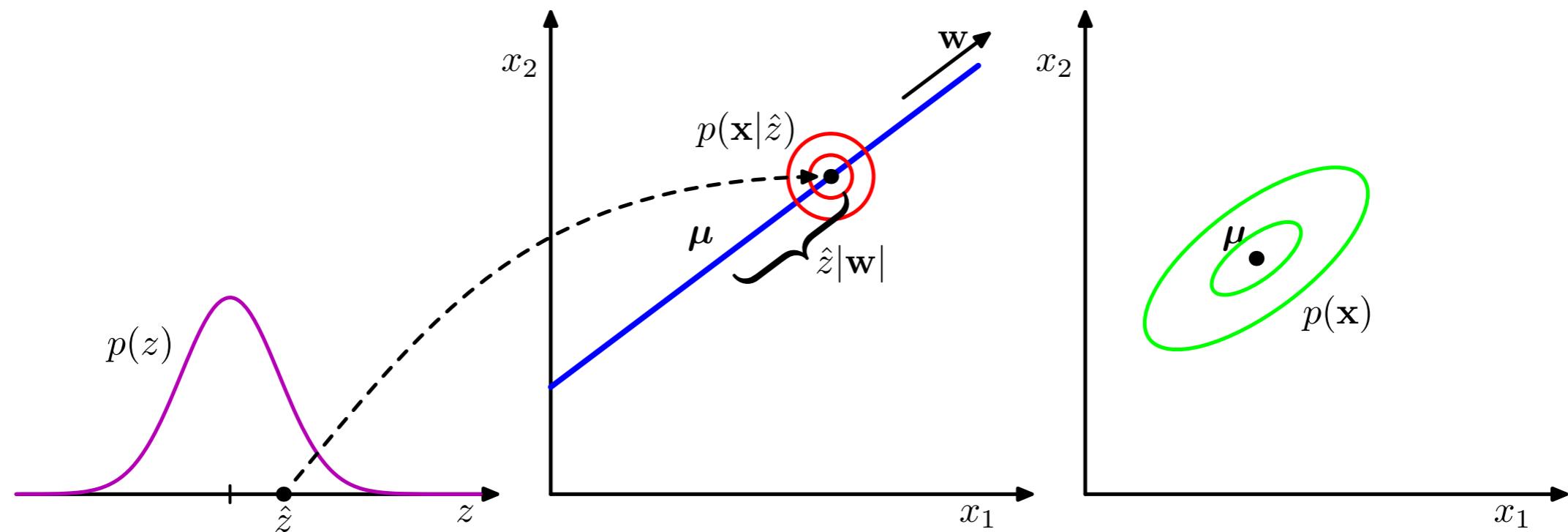


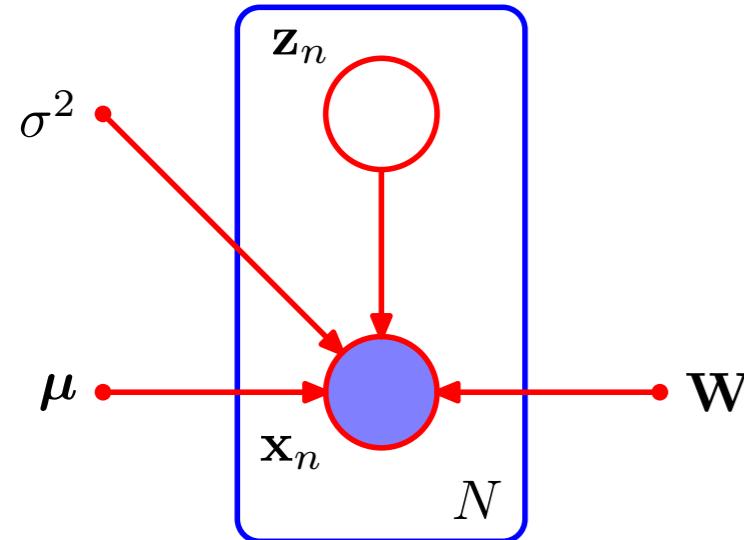
Figure 12.9 An illustration of the generative view of the probabilistic PCA model for a two-dimensional data space and a one-dimensional latent space. An observed data point \mathbf{x} is generated by first drawing a value \hat{z} for the latent variable from its prior distribution $p(z)$ and then drawing a value for \mathbf{x} from an isotropic Gaussian distribution (illustrated by the red circles) having mean $w\hat{z} + \mu$ and covariance $\sigma^2\mathbf{I}$. The green ellipses show the density contours for the marginal distribution $p(\mathbf{x})$.

Bishop's book, page 572



Predictions with PPCA

We want to compute



$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

From Bayes' rule:

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) d\mathbf{z}.$$

This has closed form: $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$

where $\mathbf{C} = \mathbf{WW}^T + \sigma^2 \mathbf{I}$.

Finally,

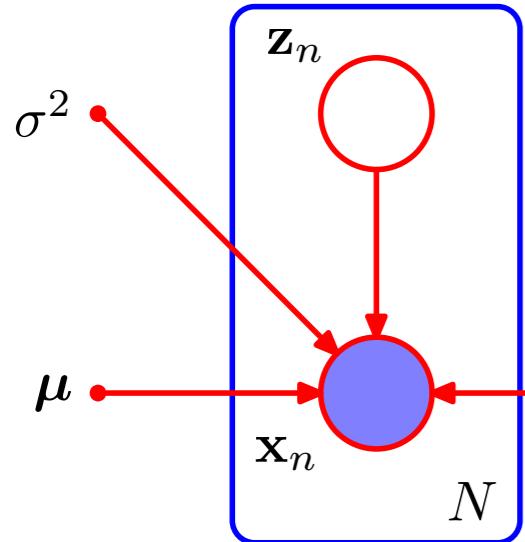
$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2}\mathbf{M})$$

where $\mathbf{M} = \mathbf{W}^T\mathbf{W} + \sigma^2 \mathbf{I}$.



Estimating parameters: Maximum likelihood solution

Given a dataset



$$\mathbf{X} = \{\mathbf{x}_n\}$$

The log-likelihood is

$$\begin{aligned}\ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})\end{aligned}$$

Setting the derivatives to zero, the parameters can be computed in closed form

$$\boldsymbol{\mu} = \bar{\mathbf{x}}$$

$$\mathbf{W}_{\text{ML}} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$$

$$\sigma_{\text{ML}}^2 = \frac{1}{D-M} \sum_{i=M+1}^D \lambda_i$$

Check Bishop's book for the definitions



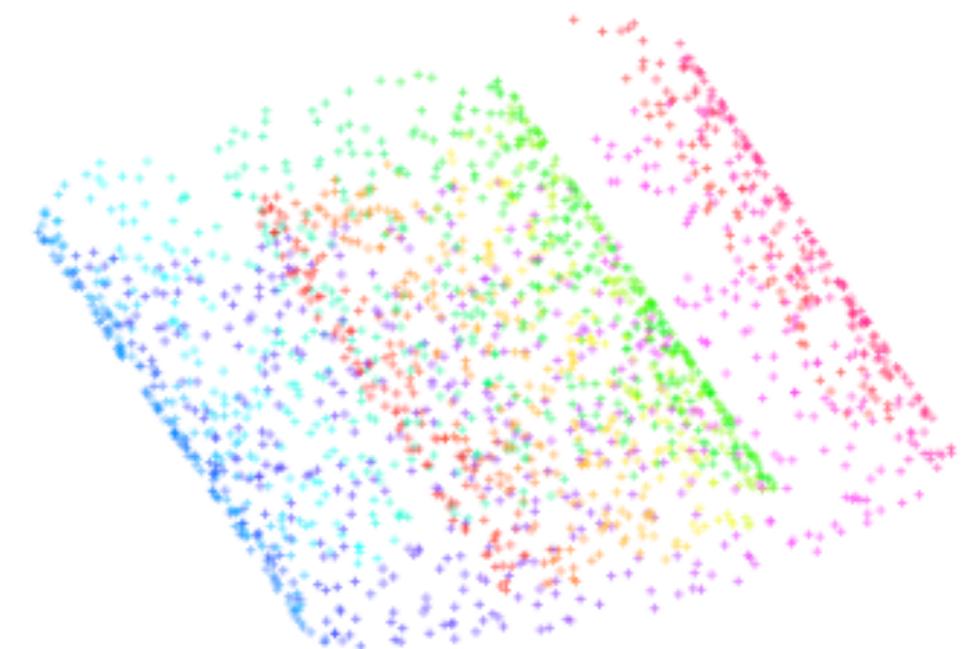
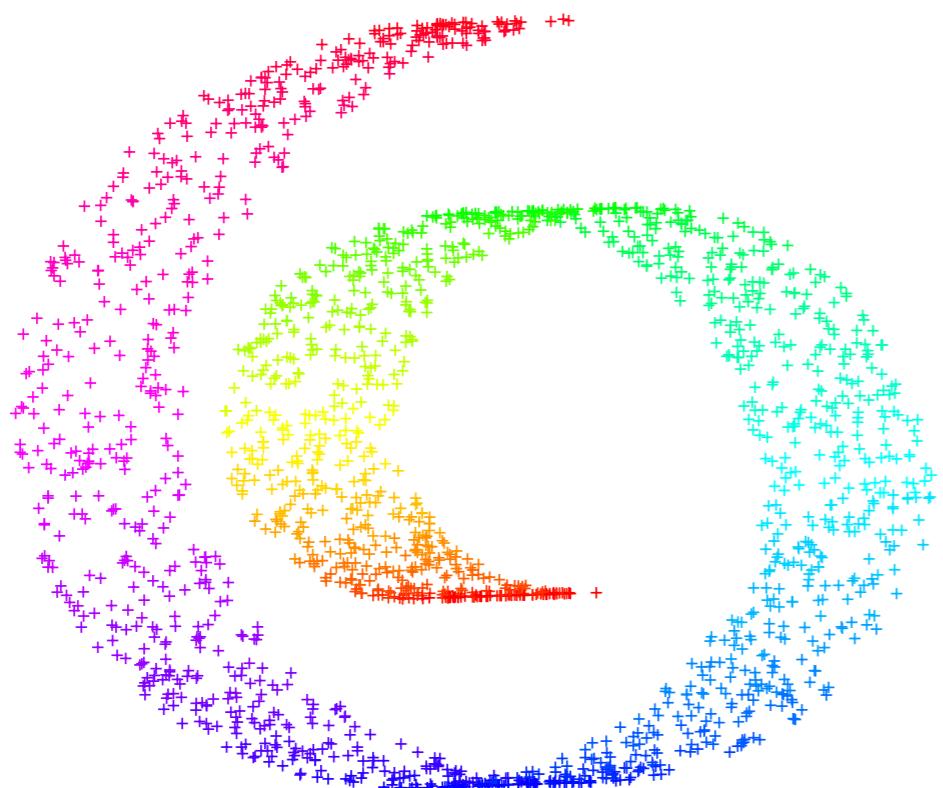
THE UNIVERSITY OF
SYDNEY

Research Topics



THE UNIVERSITY OF
SYDNEY

Problems with Linear Mappings



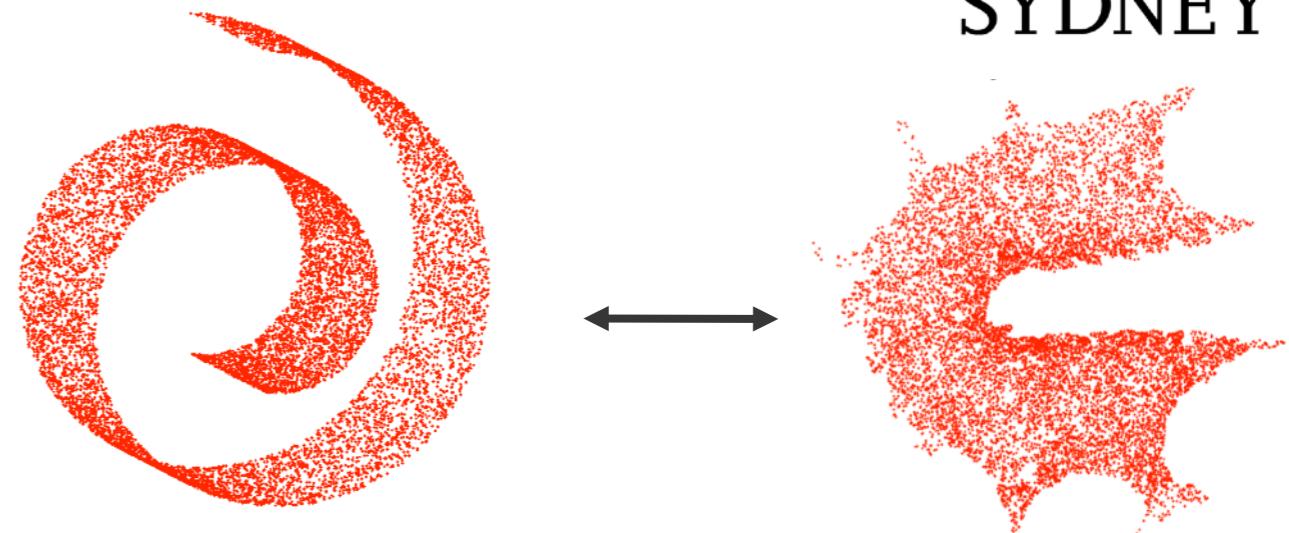
PCA



THE UNIVERSITY OF
SYDNEY

A perfect method

1. A bijective mapping ?
2. A “nice” mapping ?
3. A mapping that preserves distances ?
4. A mapping that preserves topology (neighbors) ?



Importance (and difficulty) to **evaluate** projections



Graph distances

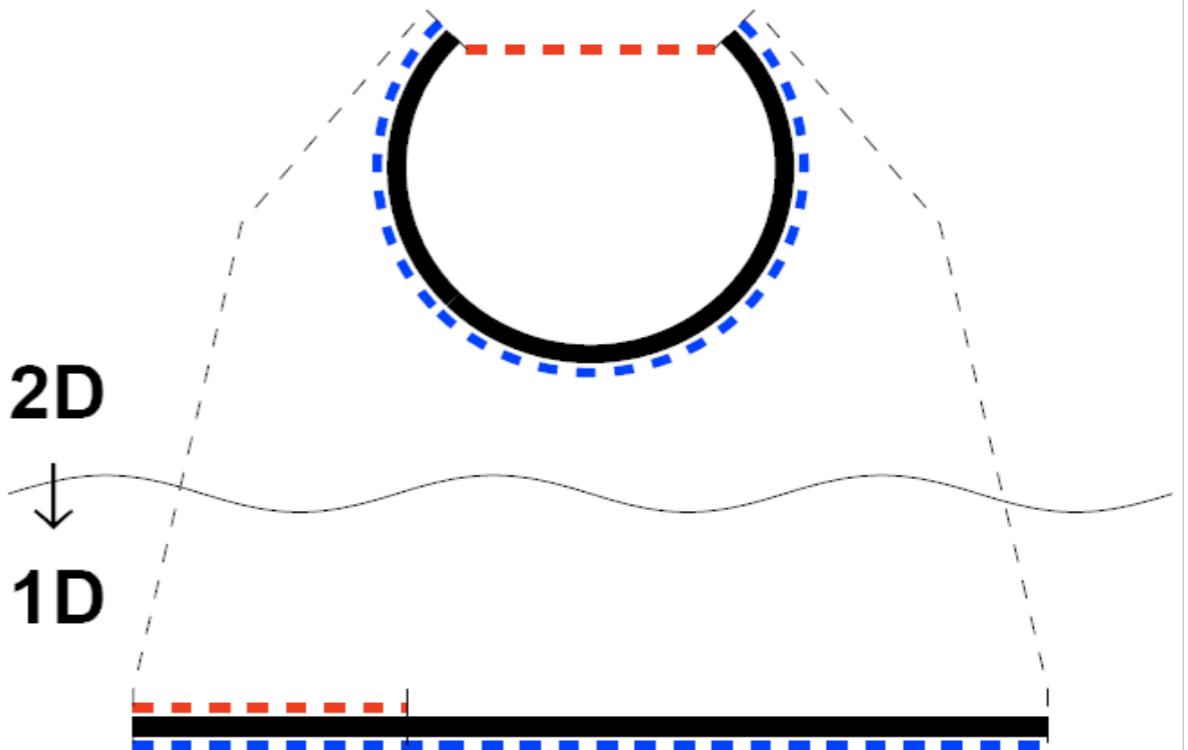
Euclidean distances must always be increased to be projected

Solution: to measure distances along the manifold

Less discrepancy between d_x and d_y , so easier optimization, etc.

So-called geodesic distance

- Manifold
- Geodesic distance
- Euclidean distance





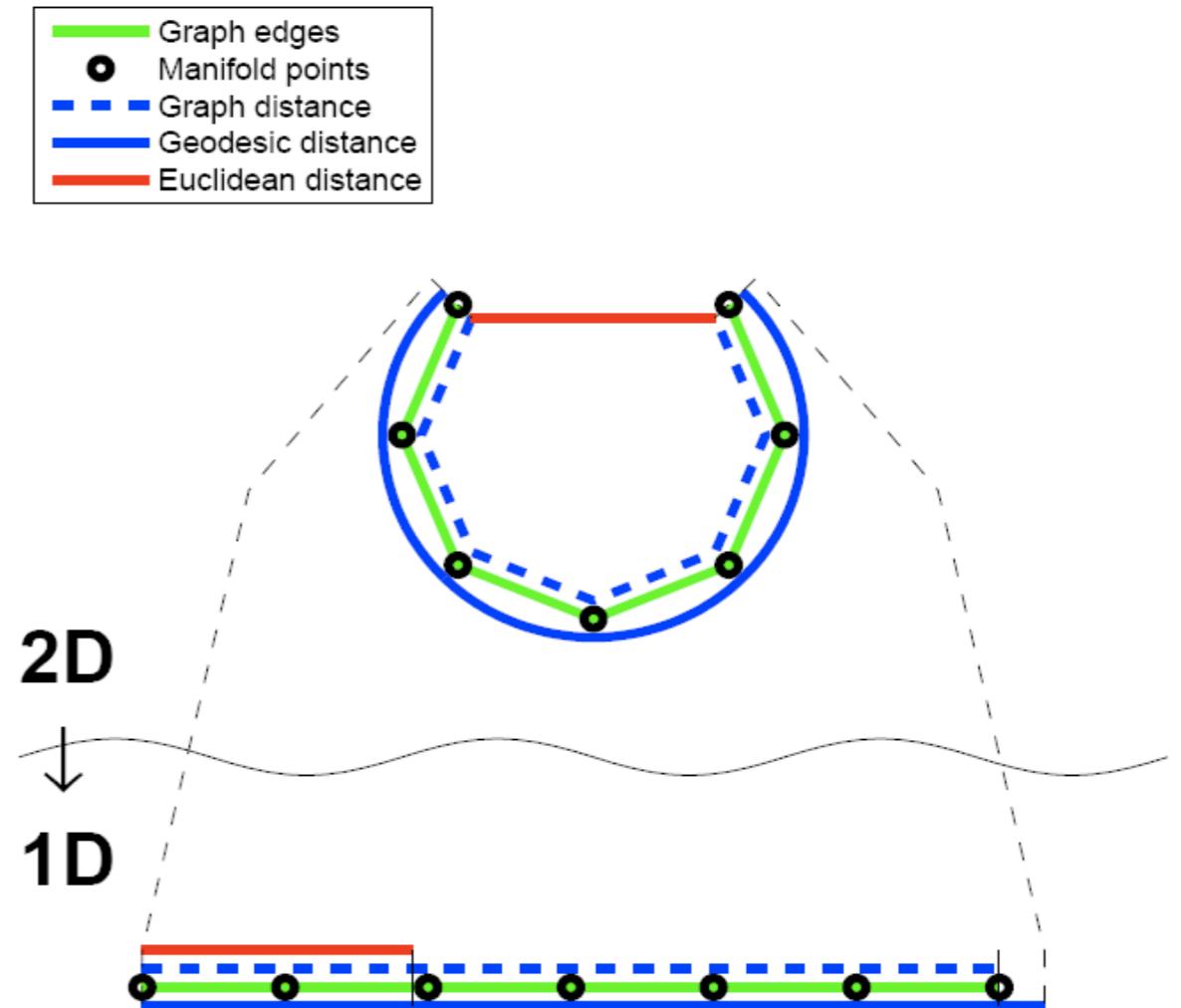
THE UNIVERSITY OF
SYDNEY

Approximating graph distances

Problem: in $d > 1$, geodesic distances may be measured along an **infinity of paths**
→ functional optimisation, **intractable** in practice

Idea: approximate manifolds with graphs

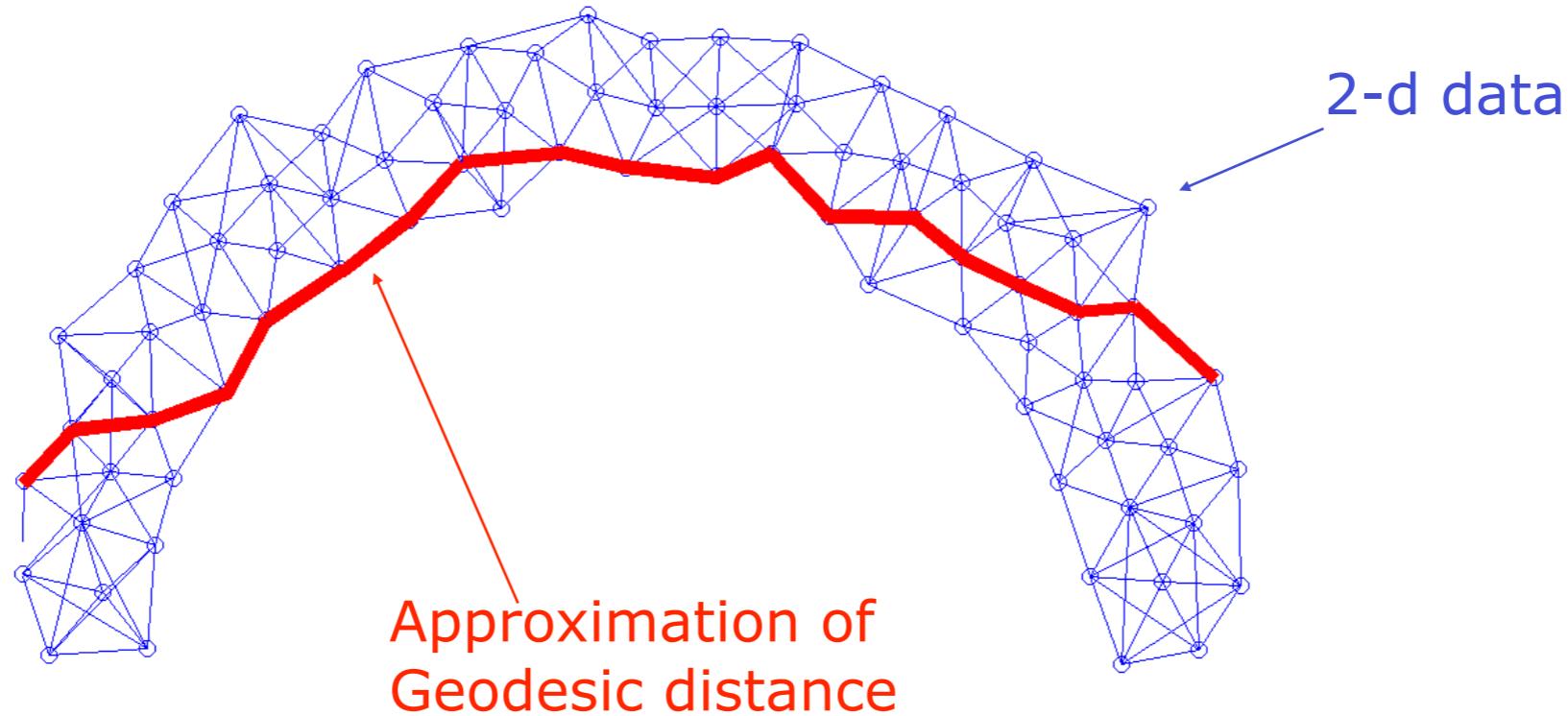
Graph distance
= sum of edge lengths
 \approx geodesic distance





THE UNIVERSITY OF
SYDNEY

Geodesic distances



How to build the graph from the data?

Connect each data to its k nearest neighbours, or

Connect each data to all other ones in a ε -ball

Ensure connectivity of the graph



THE UNIVERSITY OF
SYDNEY

Isomap

Build distance matrix

Compute nearest neighbour graph

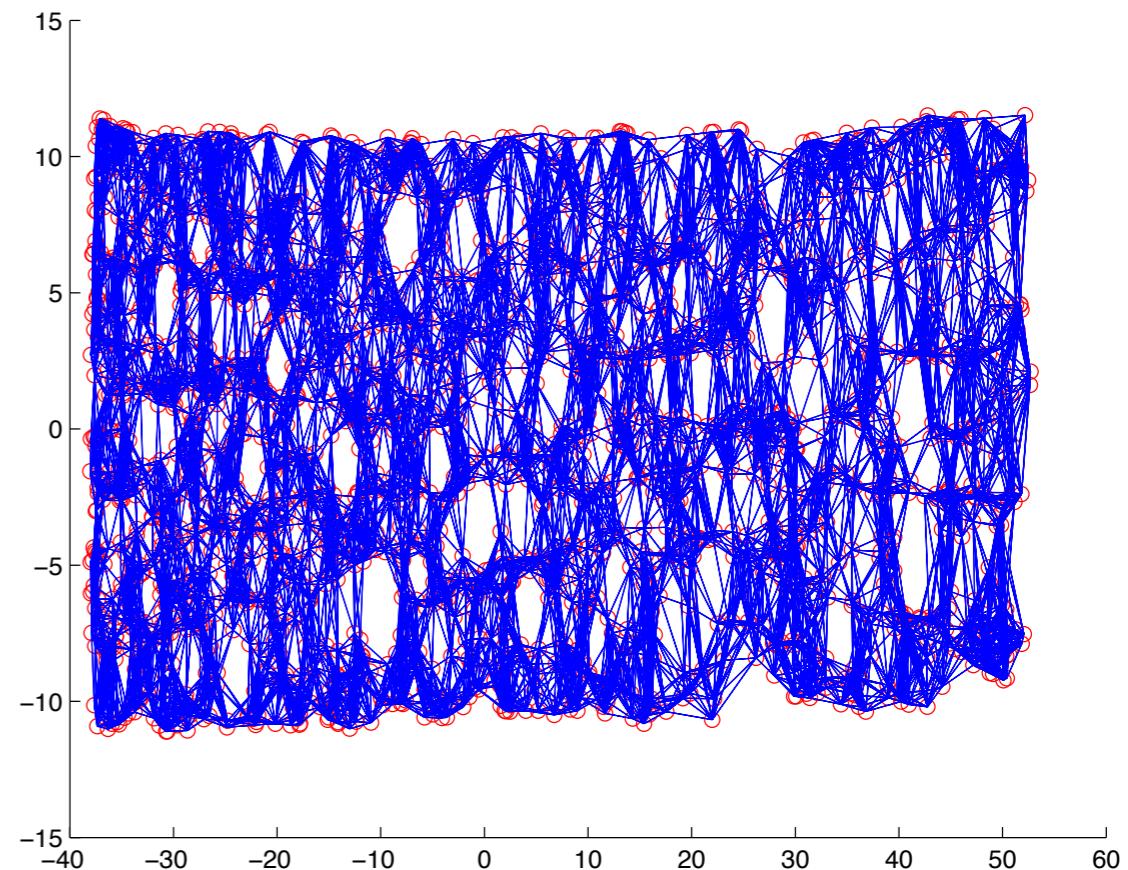
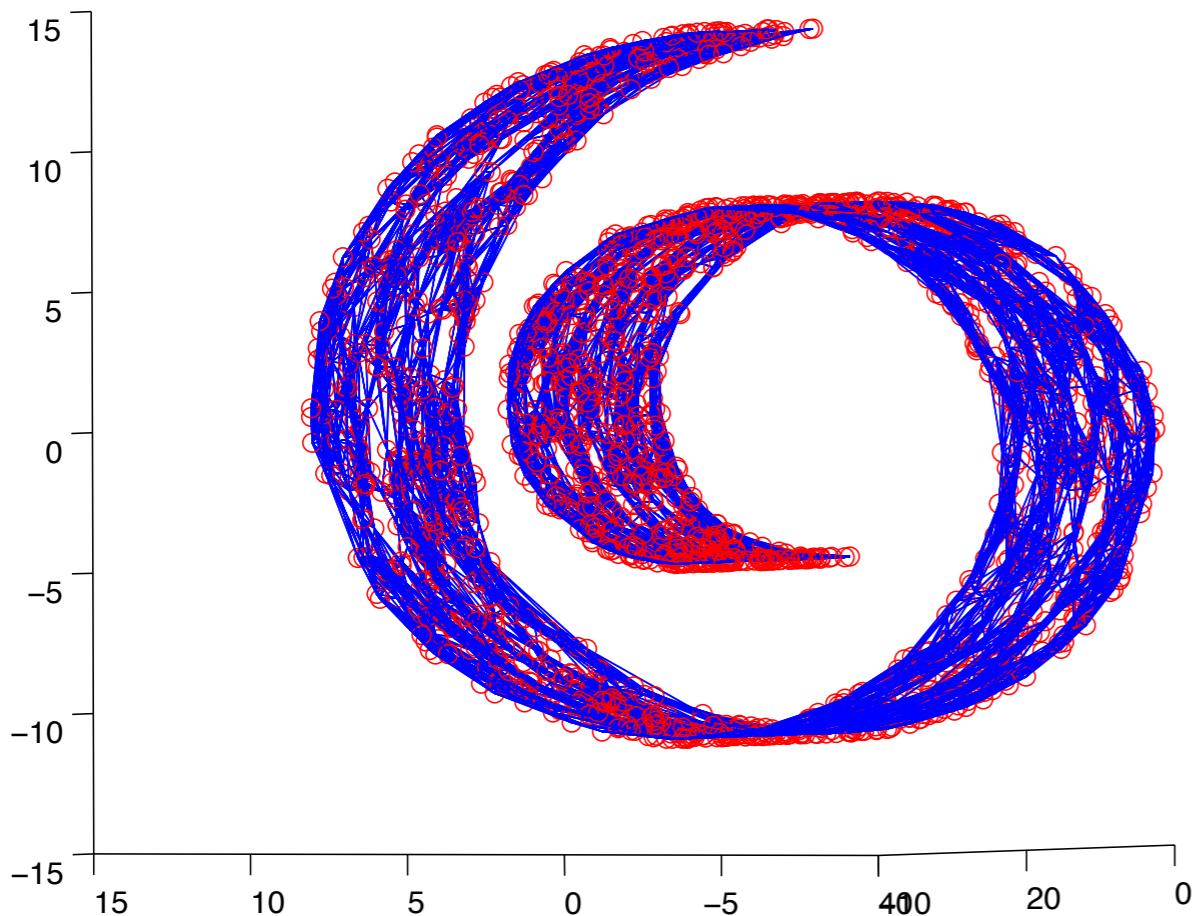
Approximate geodesic distance by shortest path

Metric MDS with geodesic (graph) distances



Isomap

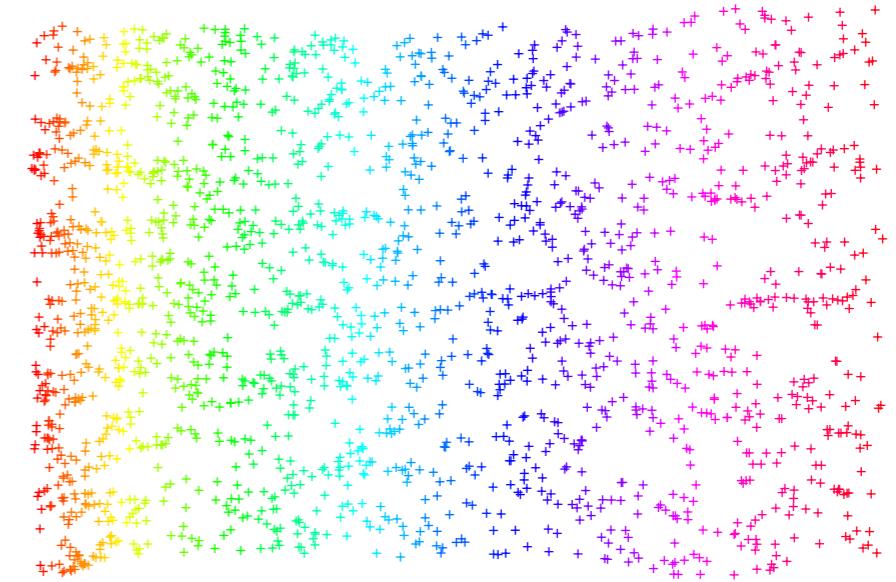
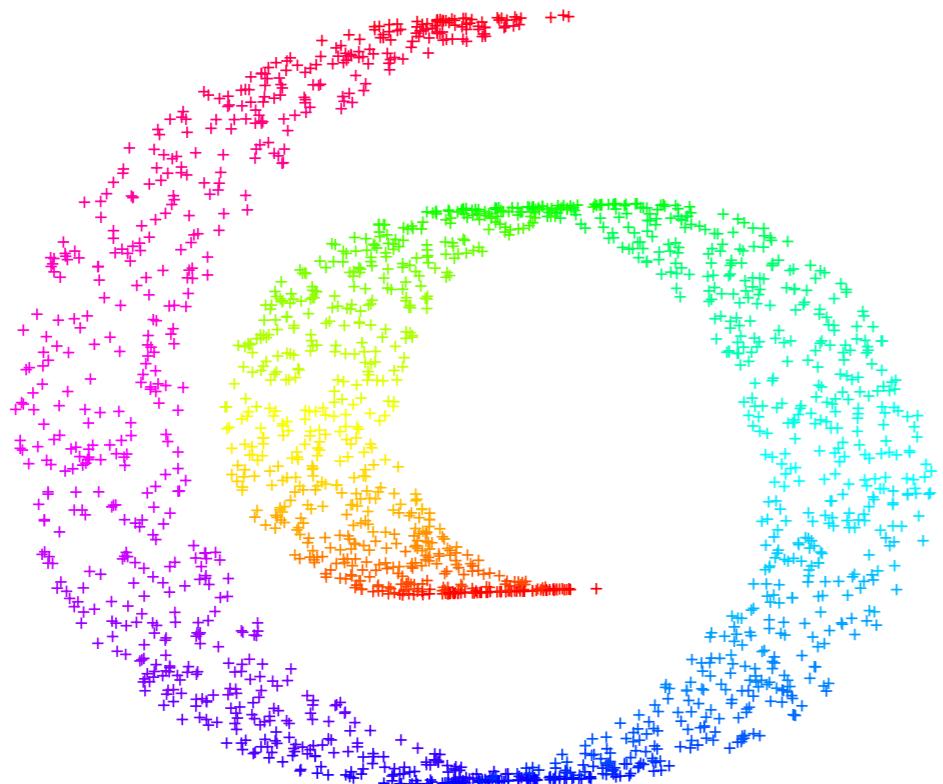
THE UNIVERSITY OF
SYDNEY





Isomap

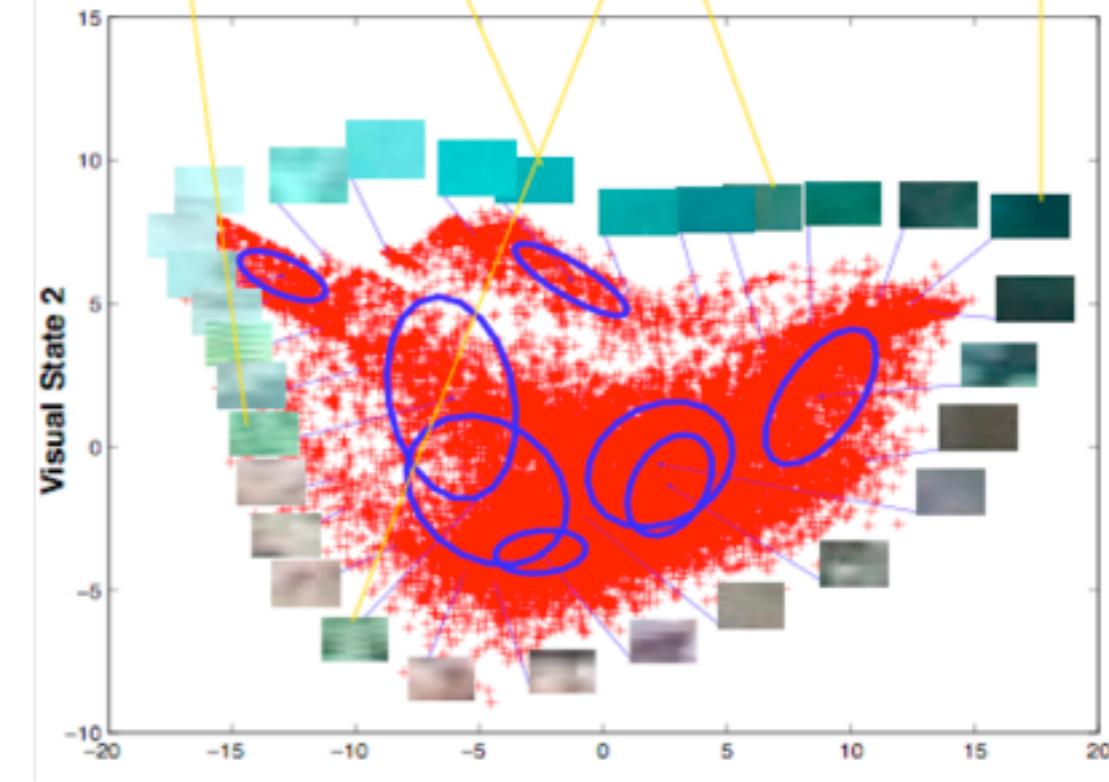
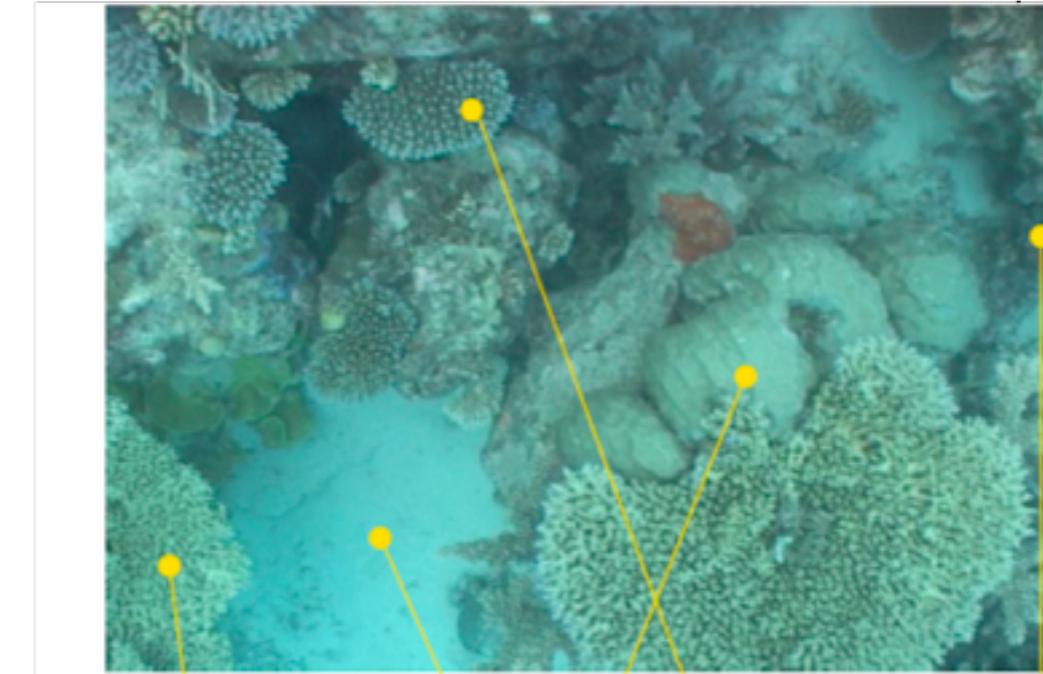
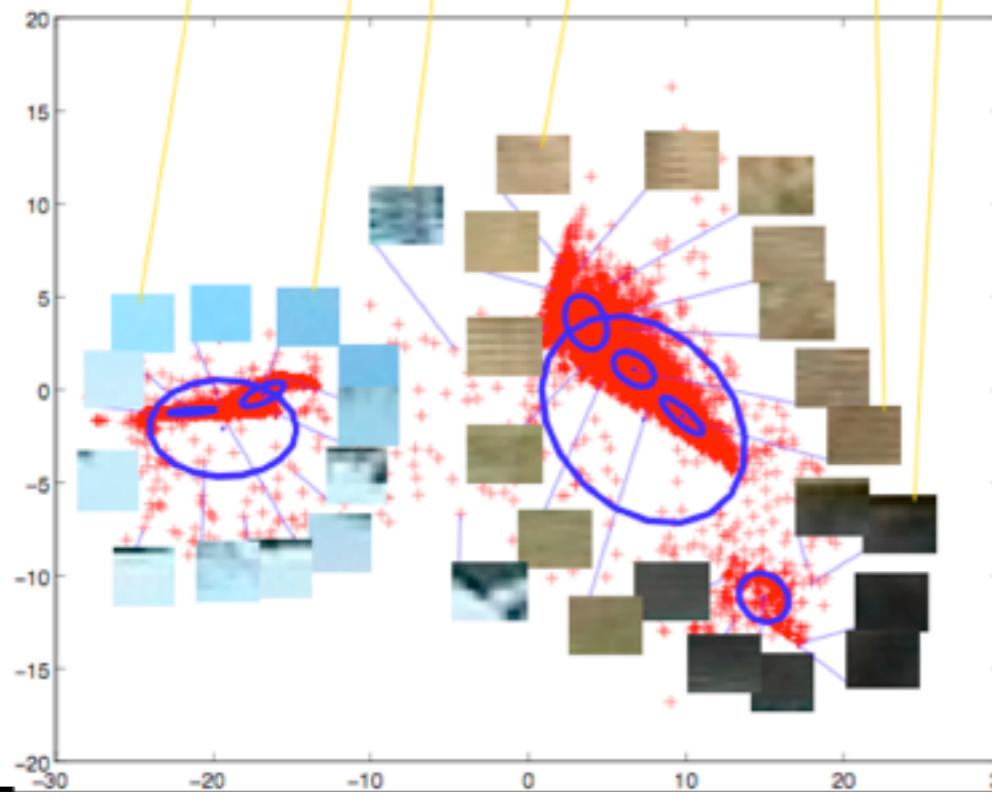
THE UNIVERSITY OF
SYDNEY





Unstructured Environments

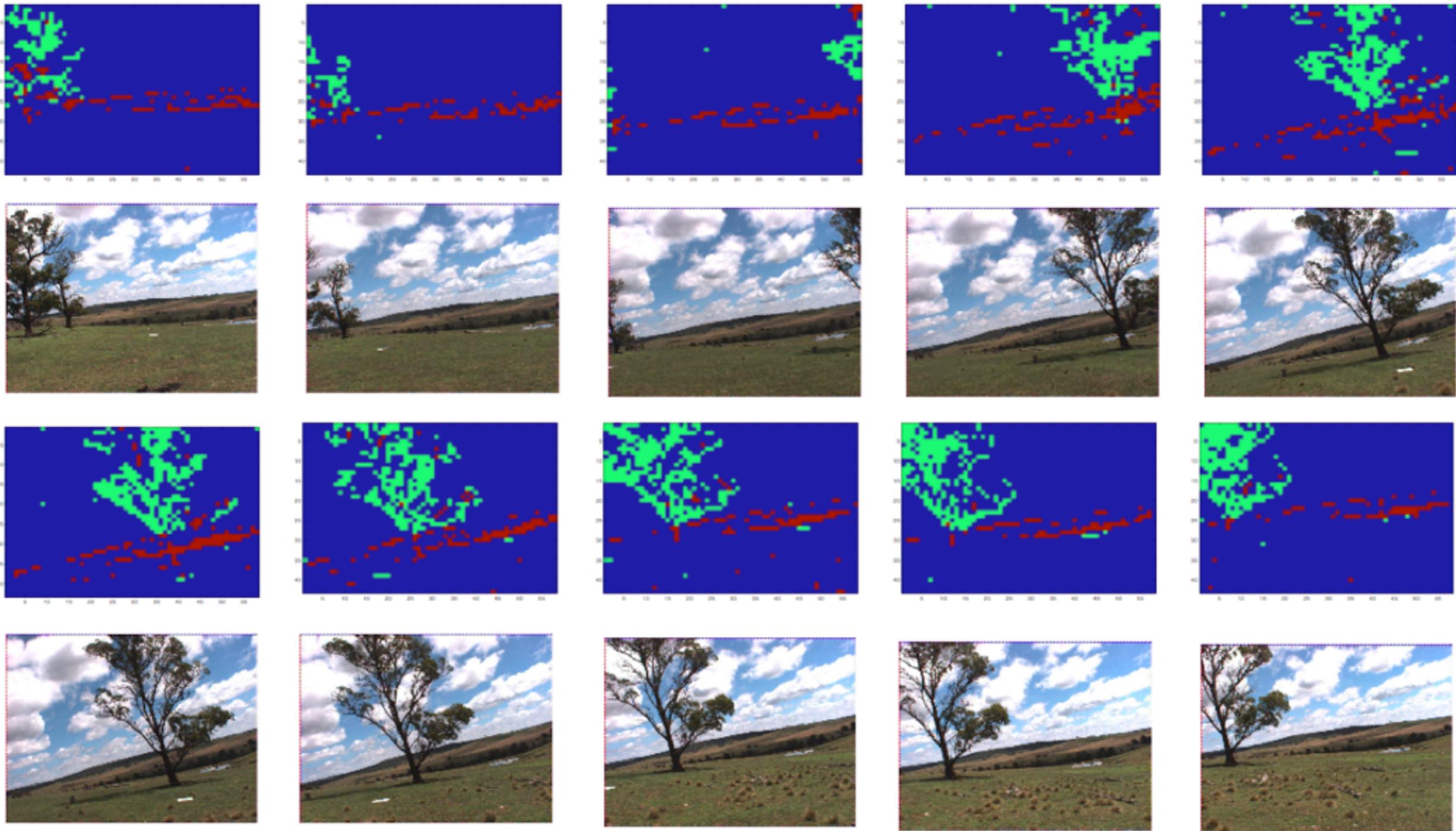
THE UNIVERSITY OF
SYDNEY





Segmentation

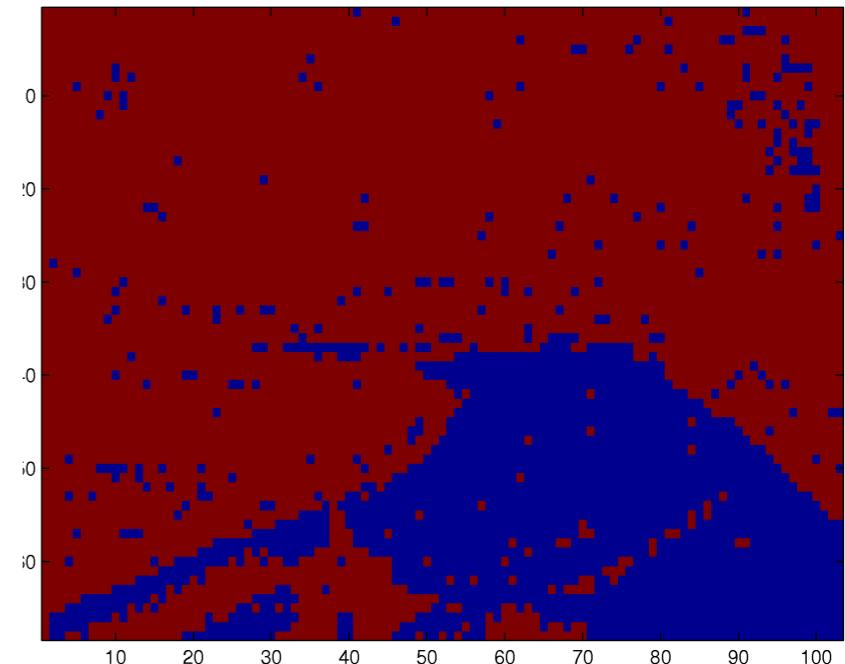
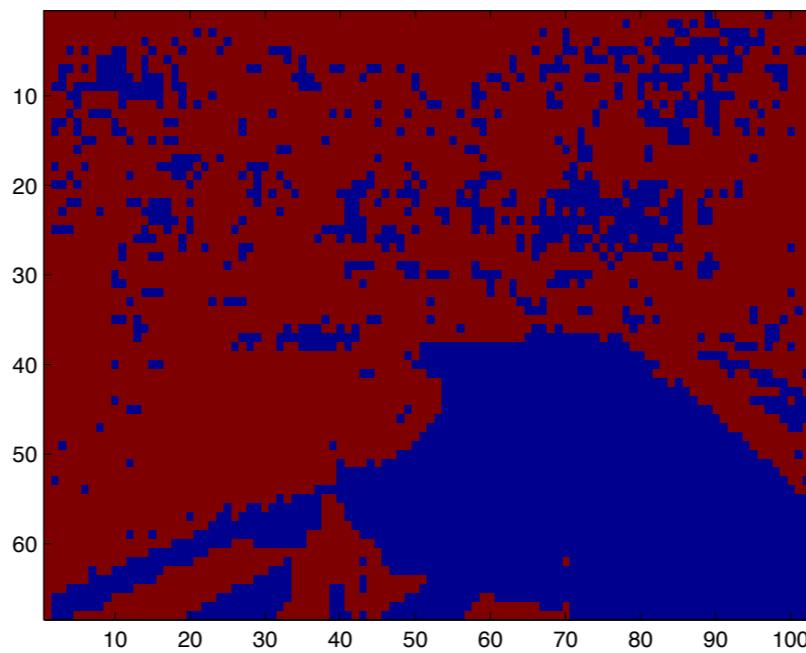
THE UNIVERSITY OF
SYDNEY





THE UNIVERSITY OF
SYDNEY

Road Segmentation





THE UNIVERSITY OF
SYDNEY

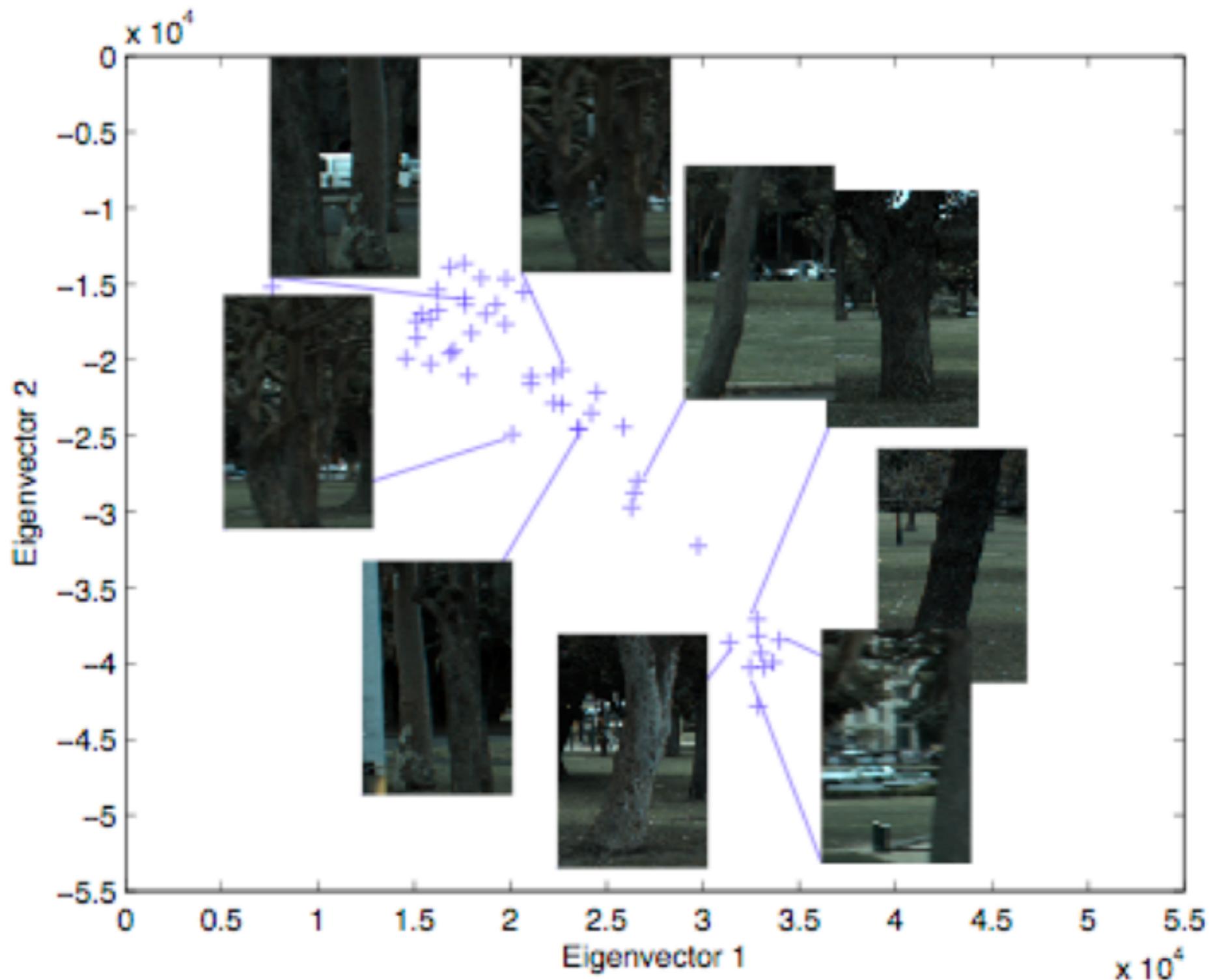
Multi-View Representations





THE UNIVERSITY OF
EY

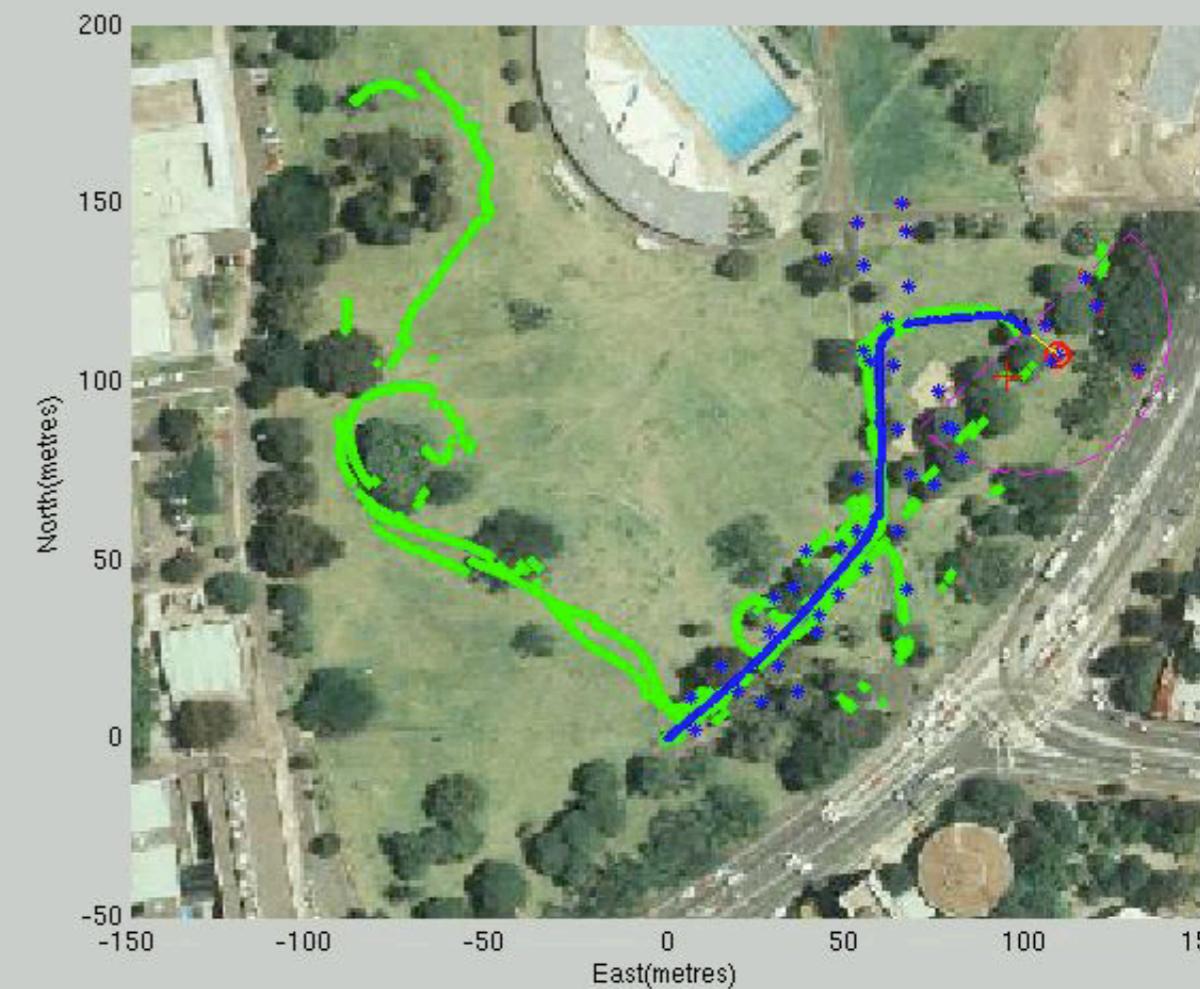
Trees Manifold





THE UNIVERSITY OF
SYDNEY

Closing the Loop in SLAM



time=214.14 secs cputime=2768.69 secs
pos:[100.941 114.661 -46.951]