

DD2434 - Advanced Machine Learning

Hierarchical Models

Carl Henrik Ek
`{chek}@csc.kth.se`

Royal Institute of Technology

December 1st, 2014



Last Lecture

- Representation Learning
 - ▶ Same story as before
 - ▶ Priors even more important
 - ▶ PPCA
 - ▶ GP-LVM
- Quickly: Multidimensional Scaling



Sensory Data

What we are doing

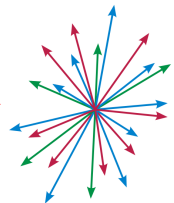
- Sensory representation
 - ▶ Capturing process
 - ▶ Pixels, Waveforms
- Degrees of freedom and dimensionality



Sensory Data

What we are doing

- Sensory representation
 - ▶ Capturing process
 - ▶ Pixels, Waveforms
- Degrees of freedom and dimensionality



Sensory Data

What we are doing

- Sensory representation
 - ▶ Capturing process
 - ▶ Pixels, Waveforms
- Degrees of freedom and dimensionality

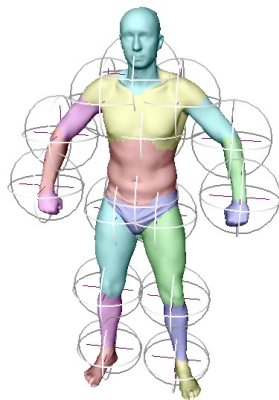


Image data

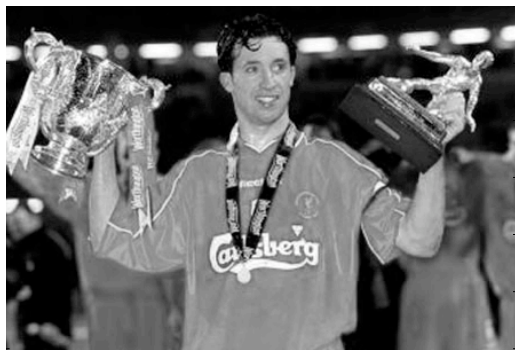


Image data

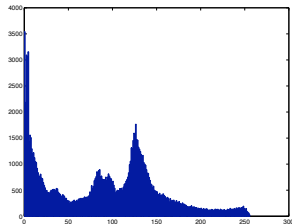
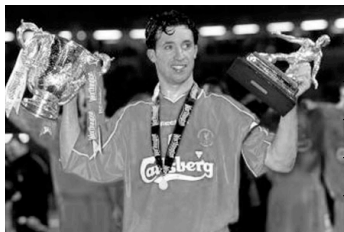


Image data

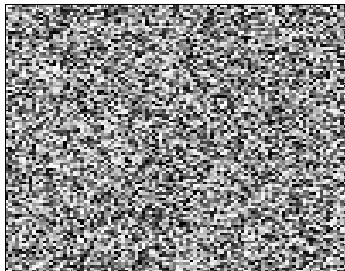
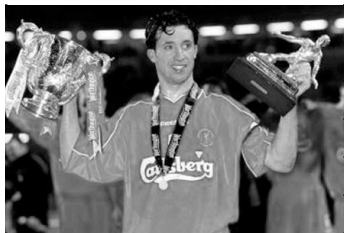


Image data

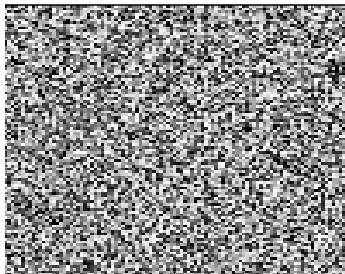
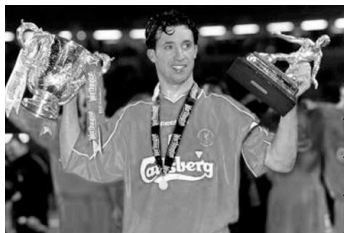


Image data

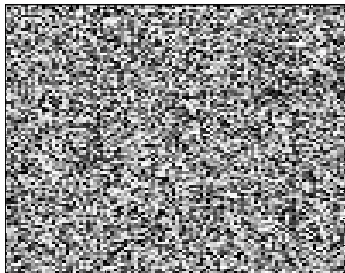
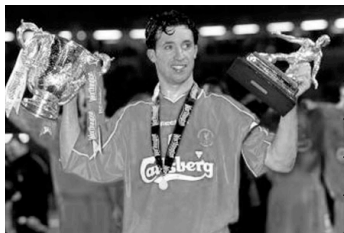
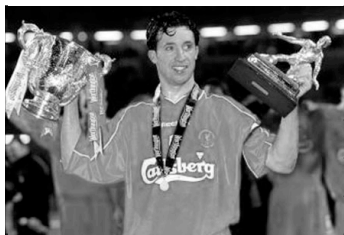


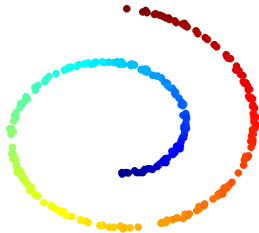
Image data

- Parametrisation
- Degrees of Freedom
- Generating parameters



Motivation

- Want to re-parametrise data
- Computational efficiency
- Discover “data-driven” degrees of freedom
 - ▶ Unravel data-manifold
- Interpretability
- Generalisation



Latent Variable Models¹

$$p(\mathbf{X}) \tag{1}$$

- We have observed some data \mathbf{X}
- Lets assume that $\mathbf{X} \in \mathbb{R}^{N \times d}$ have been generated from $\mathbf{Z} \in \mathbb{R}^{N \times q}$
- \mathbf{Z} - latent variable
- f - generative mapping

¹Murphy 2012, p. 12.

Latent Variable Models¹

$$p(\mathbf{X}|f, \mathbf{Z}) \quad (2)$$

$$\mathbf{f} : \mathbf{Z} \rightarrow \mathbf{X} \quad (3)$$

- We have observed some data \mathbf{X}
- Lets assume that $\mathbf{X} \in \mathbb{R}^{N \times d}$ have been generated from $\mathbf{Z} \in \mathbb{R}^{N \times q}$
- \mathbf{Z} - latent variable
- f - generative mapping

¹Murphy 2012, p. 12.

Latent Variable Models¹

$$p(\mathbf{X}|f, \mathbf{Z}) \tag{4}$$

$$\mathbf{f} : \mathbf{Z} \rightarrow \mathbf{X} \tag{5}$$

- We have observed some data \mathbf{X}
- Lets assume that $\mathbf{X} \in \mathbb{R}^{N \times d}$ have been generated from $\mathbf{Z} \in \mathbb{R}^{N \times q}$
- \mathbf{Z} - latent variable
- f - generative mapping

¹Murphy 2012, p. 12.

Latent Variable Models¹

$$p(\mathbf{X}|f, \mathbf{Z}) \tag{6}$$

$$\mathbf{f} : \mathbf{Z} \rightarrow \mathbf{X} \tag{7}$$

- We have observed some data \mathbf{X}
- Lets assume that $\mathbf{X} \in \mathbb{R}^{N \times d}$ have been generated from $\mathbf{Z} \in \mathbb{R}^{N \times q}$
- \mathbf{Z} - latent variable
- f - generative mapping

¹Murphy 2012, p. 12.

WTF?

The strength of Priors

- Encodes prior belief
- This can also be seen as a preference
 - ▶ Given several perfectly valid solutions which one do i prefer
 - ▶ Regularises solution space
- Latent variable models what do we prefer?

Factor Analysis²

$$\mathbf{x}_i = \mathbf{W}\mathbf{z}_i + \epsilon \quad (8)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \Psi) \quad (9)$$

- Assume the generating mapping to be linear
- For regression we assumed that we knew the inputs \mathbf{Z}
- Now we do not

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$\mathbf{x}_i = \mathbf{W}\mathbf{z}_i + \epsilon \quad (10)$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \Psi) \quad (11)$$

- Assume the generating mapping to be linear
- For regression we assumed that we knew the inputs \mathbf{Z}
- Now we do not

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$\mathbf{x}_i = \mathbf{W}\mathbf{z}_i + \epsilon \quad (12)$$

$$p(\mathbf{X}|\mathbf{Z}, \theta) = \mathcal{N}(\mathbf{W}\mathbf{Z}, \Psi) \quad (13)$$

$$p(\mathbf{Z}) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad (14)$$

- Assume the generating mapping to be linear
- For regression we assumed that we knew the inputs \mathbf{Z}
- Now we do not \Rightarrow specify a prior

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$p(\mathbf{X}|\theta) = \int p(\mathbf{X}|\mathbf{Z}, \theta)p(\mathbf{Z})d\mathbf{Z} = \quad (15)$$

$$= \mathcal{N}(\mathbf{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\boldsymbol{\Sigma}_0\mathbf{W}^T) \quad (16)$$

- **Z** and **W** are related
- Integrate out **Z**
 - ▶ pick $\boldsymbol{\mu}_0 = 0, \boldsymbol{\Sigma}_0 = \mathbf{I}$
- Low dimensional density model of **X**
 - ▶ $\mathcal{O}(QD)$ compared to $\mathcal{O}(D^2)$

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$p(\mathbf{X}|\theta) = \int p(\mathbf{X}|\mathbf{Z}, \theta)p(\mathbf{Z})d\mathbf{Z} = \quad (17)$$

$$= \mathcal{N}(\mathbf{W}\mu_0 + \mu, \Psi + \mathbf{W}\Sigma_0\mathbf{W}^T) \quad (18)$$

$$= \mathcal{N}(\mu, \Psi + \mathbf{W}\mathbf{W}^T) \quad (19)$$

- **Z** and **W** are related
- Integrate out **Z**
 - ▶ pick $\mu_0 = 0, \Sigma_0 = \mathbf{I}$
- Low dimensional density model of **X**
 - ▶ $\mathcal{O}(QD)$ compared to $\mathcal{O}(D^2)$

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$p(\mathbf{X}|\theta) = \int p(\mathbf{X}|\mathbf{Z}, \theta)p(\mathbf{Z})d\mathbf{Z} = \quad (20)$$

$$= \mathcal{N}(\mathbf{W}\mu_0 + \mu, \Psi + \mathbf{W}\Sigma_0\mathbf{W}^T) \quad (21)$$

$$= \mathcal{N}(\mu, \Psi + \mathbf{W}\mathbf{W}^T) \quad (22)$$

- \mathbf{Z} and \mathbf{W} are related
- Integrate out \mathbf{Z}
 - ▶ pick $\mu_0 = 0, \Sigma_0 = \mathbf{I}$
- Low dimensional density model of \mathbf{X}
 - ▶ $\mathcal{O}(QD)$ compared to $\mathcal{O}(D^2)$

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$\tilde{\mathbf{W}} = \mathbf{W}\mathbf{R} \quad (23)$$

$$p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{R}\mathbf{R}^T\mathbf{W}^T) \quad (24)$$

$$= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T) \quad (25)$$

$$(26)$$

Identifiability

- The marginal likelihood is invariant to a rotation
 - ▶ no unique solution
 - ▶ model is the same but interpretation tricky

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$\mathbf{W}_{ML} = \operatorname{argmax}_{\mathbf{W}} p(\mathbf{X}|\boldsymbol{\theta}) \quad (27)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (28)$$

Probabilistic PCA

- Dimensions of \mathbf{X} independent given \mathbf{Z}
 - ▶ \mathbf{W} orthogonal matrix
- Closed form solution Murphy 2012, p. 12.2.2

²Murphy 2012, p. 12.1.1.

Factor Analysis²

$$\mathbf{W}_{ML} = \operatorname{argmax}_{\mathbf{W}} p(\mathbf{X}|\boldsymbol{\theta}) \quad (29)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (30)$$

$$\mathbf{W}_{ML} = \mathbf{U}_q (\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{\frac{1}{2}} \quad (31)$$

$$\mathbf{S} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T \quad (32)$$

Probabilistic PCA

- Dimensions of \mathbf{X} independent given \mathbf{Z}
 - ▶ \mathbf{W} orthogonal matrix
- Closed form solution Murphy 2012, p. 12.2.2

²Murphy 2012, p. 12.1.1.

Factor Analysis²

Summary

- Factor Analysis is a linear continuous latent variable model
- Solution not unique
- PCA is Factor Analysis with two assumptions
 - ▶ factor loadings orthogonal $\mathbf{W}^T \mathbf{W} = \mathbf{I}$
 - ▶ noise free case $\epsilon = \lim_{\sigma^2 \rightarrow 0} \sigma^2 \mathbf{I}$
- PCA is incredibly useful but its important to know what you are assuming, the probabilistic formulation allows you to do just that

²Murphy 2012, p. 12.1.1.

Factor Analysis²

Summary

- Factor Analysis is a linear continuous latent variable model
- Solution not unique
- PCA is Factor Analysis with two assumptions
 - ▶ factor loadings orthogonal $\mathbf{W}^T \mathbf{W} = \mathbf{I}$
 - ▶ noise free case $\epsilon = \lim_{\sigma^2 \rightarrow 0} \sigma^2 \mathbf{I}$
- PCA is incredibly useful but its important to know what you are assuming, the probabilistic formulation allows you to do just that

²Murphy 2012, p. 12.1.1.

Gaussian Process Latent Variable Models

History repeats itself

- In PPCA we assumed no uncertainty in the mapping
- We can use \mathcal{GPs} over mapping
- Gaussian Process Latent Variable Model [Lawrence 2005]

Gaussian Process Latent Variable Models

History repeats itself

- In PPCA we assumed no uncertainty in the mapping
- We can use \mathcal{GPs} over mapping
- Gaussian Process Latent Variable Model [Lawrence 2005]

Gaussian Process Latent Variable Models

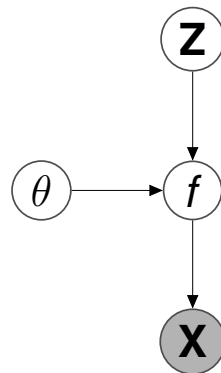
$$p(\mathbf{X}|\mathbf{f}, \mathbf{Z}, \theta) \tag{33}$$

- In PPCA we marginalised out \mathbf{Z} and optimised for \mathbf{W}
- Not possible for a general \mathcal{GP}

Gaussian Process Latent Variable Models

GP-LVM

- General co-variance function (Ex. SE)
- \mathbf{Z} appears non-linearly in relation to \mathbf{X}
- Marginalisation of \mathbf{Z} intractable



Gaussian Process Latent Variable Models

$$\operatorname{argmax}_{\mathbf{Z}, \theta} p(\mathbf{X}|\mathbf{Z}, \theta)p(\mathbf{Z}) \quad (34)$$

$$p(\mathbf{X}|\mathbf{Z}, \theta) = \int p(\mathbf{X}|\mathbf{f})p(\mathbf{f}|\mathbf{Z}, \theta)d\mathbf{f} \quad (35)$$

$$p(\mathbf{Z}) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (36)$$

- **GP**-prior sufficiently regularises objective
- Need to set dimensionality of \mathbf{Z}

Gaussian Process Latent Variable Models

- You can add different priors on latent representations
 - ▶ Topological
 - ▶ Dynamic GP and a GP
 - ▶ Classification
- Any preference you can formulate as a prior

Gaussian Process Latent Variable Models

$$\mathbf{z}_{t+1} = g(\mathbf{z}_t) + \epsilon_z \quad (37)$$

$$g \sim \mathcal{GP}(\mathbf{0}, k(\mathbf{z}_i, \mathbf{z}_j)) \quad (38)$$

- You can add different priors on latent representations
 - ▶ Topological
 - ▶ Dynamic GP and a GP
 - ▶ Classification
- Any preference you can formulate as a prior

Gaussian Process Latent Variable Models

- You can add different priors on latent representations
 - ▶ Topological
 - ▶ Dynamic GP and a GP
 - ▶ Classification
- Any preference you can formulate as a prior

Gaussian Process Latent Variable Models

- You can add different priors on latent representations
 - ▶ Topological
 - ▶ Dynamic GP and a GP
 - ▶ Classification
- Any preference you can formulate as a prior

Multidimensional Scaling

- N entities with proximity relations δ_{ij}
- Must be metric
- Find embedding $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_N]^T$ to minimize

$$E_{MDS} = \|\mathbf{D} - \Delta\|_F$$
$$\begin{cases} \mathbf{D}_{ij} = \|\mathbf{y}_i - \mathbf{y}_j\|_{L2} \\ \Delta_{ij} = \delta_{ij} \end{cases}$$

$$\begin{aligned}
\|\mathbf{A}\|_F &= \sqrt{\text{trace}(\mathbf{A}\mathbf{A}^T)} = \sqrt{\sum_{i=1}^N \lambda_i^2} \\
\|\mathbf{D} - \Delta\|_F &= \left\{ \Delta = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \Rightarrow \Delta = \sum_{i=1}^N \lambda_i \mathbf{v}_i \mathbf{v}_i^T \right\} = \\
&= \left\| \mathbf{D} - \sum_{i=1}^N \lambda_i \mathbf{v}_i \mathbf{v}_i^T \right\|_F = \left\| \sum_{i=1}^d q_i \mathbf{v}_i \mathbf{v}_i^T - \sum_{i=1}^N \lambda_i \mathbf{v}_i \mathbf{v}_i^T \right\|_F = \\
&= \left\| \sum_{i=1}^d (q_i - \lambda_i) \mathbf{v}_i \mathbf{v}_i^T - \sum_{i=d+1}^N \lambda_i \mathbf{v}_i \mathbf{v}_i^T \right\|_F
\end{aligned}$$

$$\text{Choose } \mathbf{D} = \mathbf{A}_{\rightarrow d} \Rightarrow E_{MDS} = \sqrt{\sum_{i=d+1}^N \lambda_i^2}$$

Multidimensional Scaling

Generate geometrical configuration \mathbf{Y} that could generate \mathbf{D}

1. Convert distance matrix D to Gram matrix $\mathbf{G} = \mathbf{Y}\mathbf{Y}^T$

▶ Proof

2. Diagonalise Gram matrix G

$$\begin{aligned}\mathbf{G} &= \mathbf{Y}\mathbf{Y}^T = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T = \left(\mathbf{V}\mathbf{\Lambda}^{\frac{1}{2}}\right) \left(\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{V}^T\right) = \\ &= \left(\mathbf{V}\mathbf{\Lambda}^{\frac{1}{2}}\right) \left(\mathbf{V} \left(\mathbf{\Lambda}^{\frac{1}{2}}\right)^T\right)^T = \left(\mathbf{V}\mathbf{\Lambda}^{\frac{1}{2}}\right) \left(\mathbf{V}\mathbf{\Lambda}^{\frac{1}{2}}\right)^T\end{aligned}$$

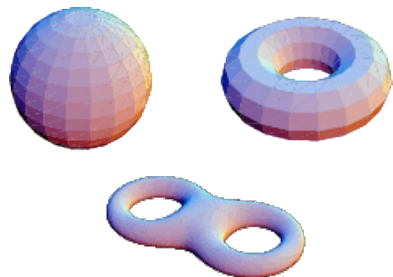
3. Chose $\mathbf{Y} = \mathbf{V}\mathbf{\Lambda}^{\frac{1}{2}}$
4. Dimension of \mathbf{Y} : $\text{rank}(\mathbf{Y}\mathbf{Y}^T) = \text{rank}(\mathbf{G}) = \text{rank}(\mathbf{D}) = d$

▶ PCA Equivalence

Non linearities

Manifold

- Generalisation of low dimensional object embedded in high dimensional space
- Similarity?
- Local similarity
- Extend local similarity to global



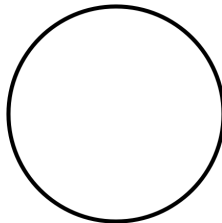
Non linearities

Definition

“In mathematics, a manifold is a topological space that near each point resembles Euclidean space”^a

^a<http://en.wikipedia.org/wiki/Manifold>

Non linearities

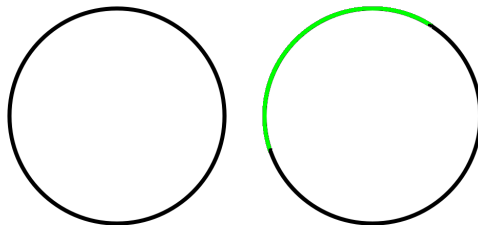


Definition

“In mathematics, a manifold is a topological space that near each point resembles Euclidean space”^a

^a<http://en.wikipedia.org/wiki/Manifold>

Non linearities

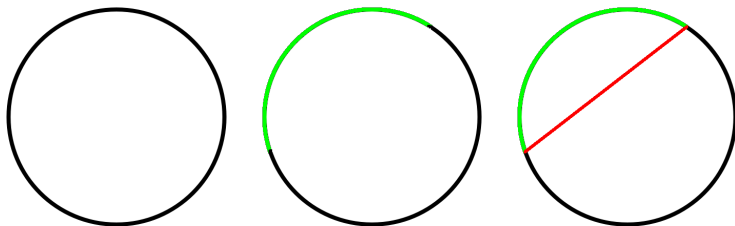


Definition

“In mathematics, a manifold is a topological space that near each point resembles Euclidean space”^a

^a<http://en.wikipedia.org/wiki/Manifold>

Non linearities

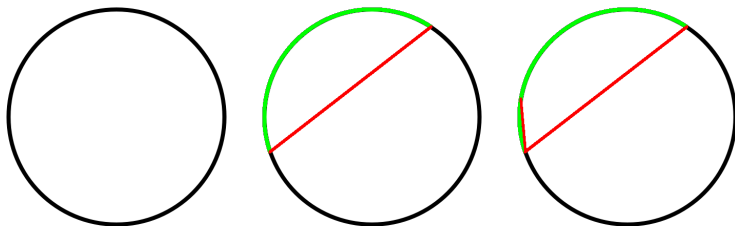


Definition

“In mathematics, a manifold is a topological space that near each point resembles Euclidean space”^a

^a<http://en.wikipedia.org/wiki/Manifold>

Non linearities

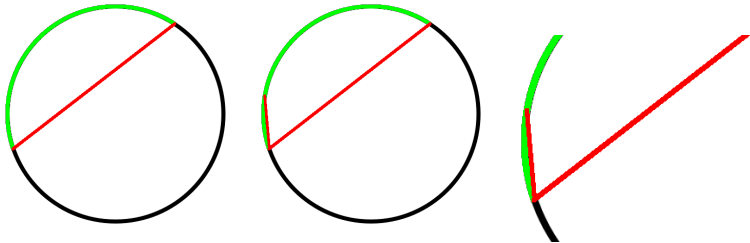


Definition

“In mathematics, a manifold is a topological space that near each point resembles Euclidean space”^a

^a<http://en.wikipedia.org/wiki/Manifold>

Non linearities



Definition

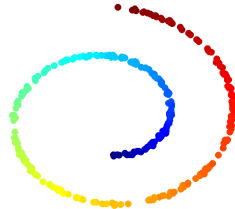
“In mathematics, a manifold is a topological space that near each point resembles Euclidean space”^a

^a<http://en.wikipedia.org/wiki/Manifold>

Non linearities

Manifold

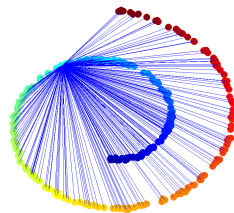
- Generalisation of low dimensional object embedded in high dimensional space
- **Similarity?**
- Local similarity
- Extend local similarity to global



Non linearities

Manifold

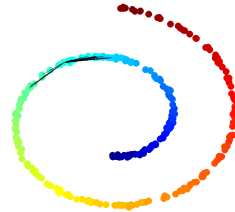
- Generalisation of low dimensional object embedded in high dimensional space
- **Similarity?**
- Local similarity
- Extend local similarity to global



Non linearities

Manifold

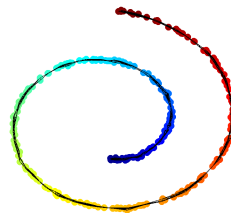
- Generalisation of low dimensional object embedded in high dimensional space
- Similarity?
- Local similarity
- Extend local similarity to global



Non linearities

Manifold

- Generalisation of low dimensional object embedded in high dimensional space
- Similarity?
- Local similarity
- Extend local similarity to global



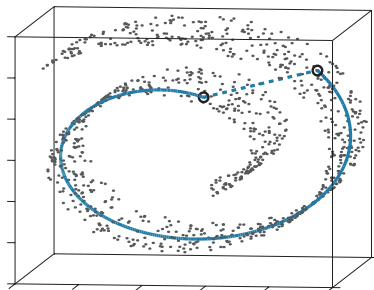
Non linearities



Proximity Graph

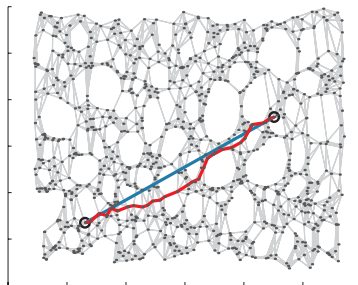
1. Identify neighbors of each data point $\mathbf{x}_i \in N(\mathbf{x}_i)$
2. Build graph $\mathbf{P} = \left\{ \underbrace{\mathbf{X}}_{\text{vertexset}}, \underbrace{\mathbf{W}}_{\text{edgeset}} \right\}$
 - ▶ Put edges between vertices's in neighborhood
 - ▶ Assume \mathbf{P} connected (and in most cases symmetric)
3. **Objective:** Complete \mathbf{P} to make it fully connected
4. Different algorithms have different strategies
 - ▶ What are the edge weights?
 - ▶ How to complete \mathbf{P}

Maximum Variance Unfolding



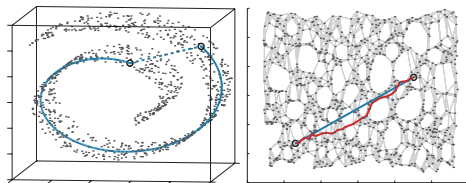
Any “fold” of the manifold between two points will **decrease** the *Euclidean* distance between the points while the *Manifold* distance remains **constant**

Maximum Variance Unfolding



If manifold is **maximally** stretched between two points the *Euclidean* distance will **equal** the *Manifold* distance

Maximum Variance Unfolding

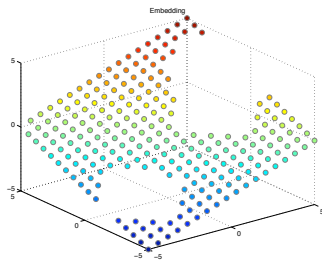


Maximise all pairwise distance outside local neighborhood (upper bound)

$$\max \sum_{i=1}^N \sum_{j=1}^N \|\mathbf{y}_i - \mathbf{y}_j\|_{L2}^2$$
$$\Rightarrow \max(\text{trace}(\mathbf{K}))$$

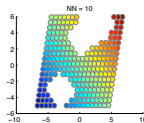
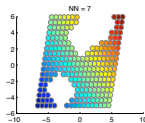
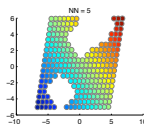
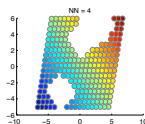
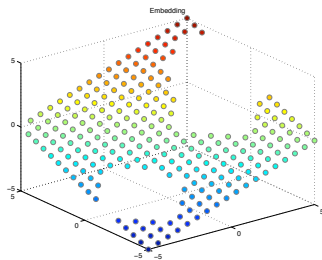
▶ Proof

Maximum Variance Unfolding: Example³



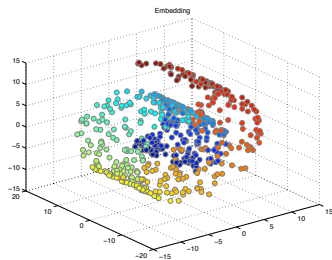
³/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example³



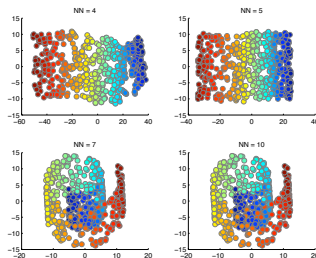
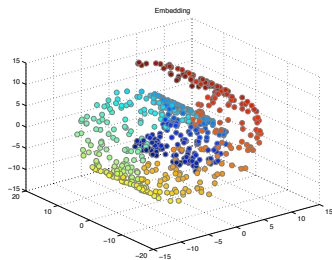
³/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example³



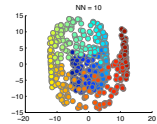
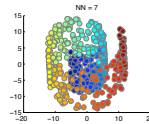
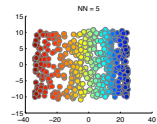
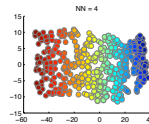
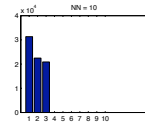
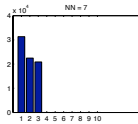
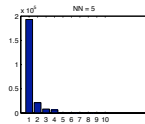
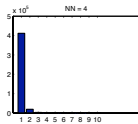
`3/algos/mvu_embed.m`

Maximum Variance Unfolding: Example³



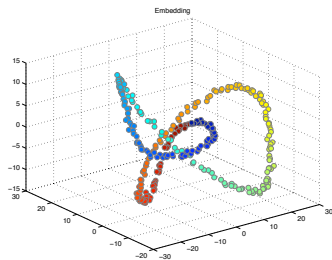
³/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example³



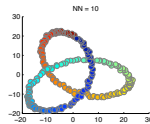
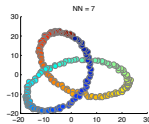
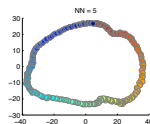
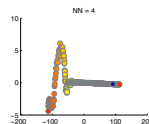
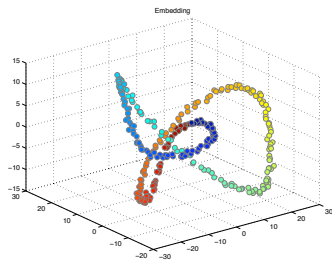
³/algos/mvu_embed.m

Maximum Variance Unfolding: Example³



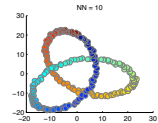
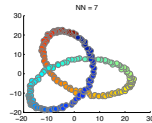
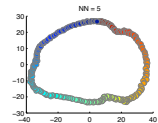
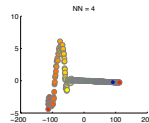
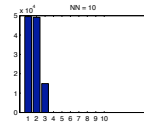
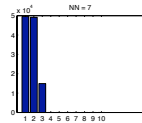
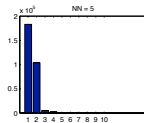
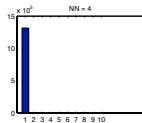
³/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example³



³/algorithms/mvu_embed.m

Maximum Variance Unfolding: Example³



³/algos/mvu_embed.m

Introduction

Recap

Hierarchical Models

Summary

Outline

- Hierarchical Models
 - ▶ motivation
 - ▶ history
 - ▶ neural networks
 - ▶ deep models
 - ▶ Why is this exciting?
- Summary of my part



$$f : \mathbf{X} \rightarrow \mathbf{Y} \quad (39)$$

Problem set-up

- Some data \mathbf{X} (input)
- Some task \mathbf{Y} (output)
- Estimate mapping from data
- Using a hierarchy

$$f : \mathbf{X} \rightarrow \mathbf{Y} \quad (40)$$

$$\mathbf{X} \rightarrow \mathbf{H}_1 \rightarrow \mathbf{H}_2 \rightarrow \dots \rightarrow \mathbf{Y} \quad (41)$$

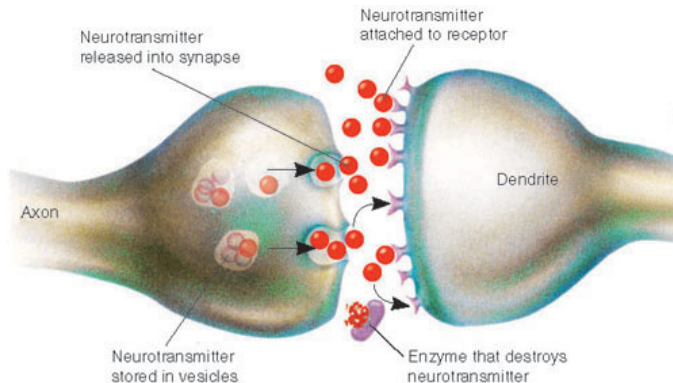
Problem set-up

- Some data \mathbf{X} (input)
- Some task \mathbf{Y} (output)
- Estimate mapping from data
- Using a hierarchy

Standing on the shoulders of giants

Deep Learning and Neural Networks

Hierarchical Models



Hierarchical Models

History 1940-1990

- Artificial Neuron McCulloch and Pitts 1943 Rosenblatt 1958
- Only linear functions Minsky and Papert 1969
- Multi-layered Perceptron Rumelhart *et al.* 1986
- Back-propagation

Hierarchical Models

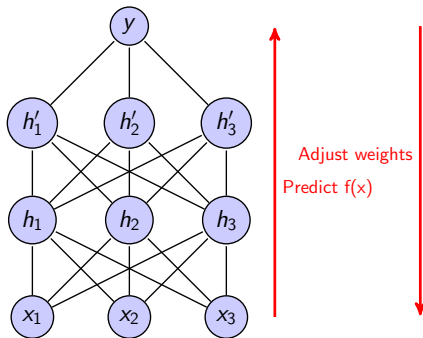
$$y_i = \rho \left(\sum_{j=0}^N w_{ij} x_j \right) \quad (42)$$

$$\rho(t) = \frac{1}{1 + e^{-t}} \quad (43)$$

Artificial Neuron

- x_j signal j into neuron i
- w_{ij} weight of signal from j
- ρ activation function

Hierarchical Models



Hierarchical Models

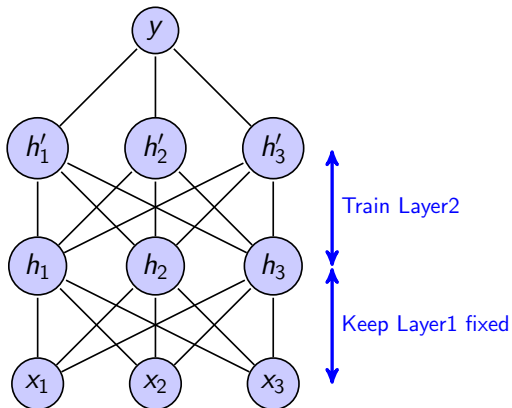


Hierarchical Models

History 2004-2010

- Vanishing Gradients
- Restricted Boltzman Machine
- Layer-wise training Hinton *et al.* 2006
 - ▶ *“If you want to do Computer Vision first learn Computer Graphics”*
- Allows for unlabeled data

Hierarchical Models



Hierarchical Models

History 2010-

- Heuristic structures
 - ▶ Convolutional Neural Networks
- Big-Data
- Infrastructural changes
 - ▶ GPUs
 - ▶ Distributed computations

Hierarchical Models



Human: "A group of men playing Frisbee in the park."

Computer model: "A group of young people playing a game of Frisbee."

1 of 6



Hierarchical Models



Human: "A young hockey player playing in the ice rink."

Computer model: "Two hockey players are fighting over the puck."

2 of 6



Hierarchical Models



Human: "Three different types of pizza on top of a stove."
Computer model: "A pizza sitting on top of a pan on top of a stove."

5 of 6



Hierarchical Models

26 June 2012 Last updated at 16:03 GMT



Google computer works out how to spot cats

A Google research team has trained a network of 1,000 computers wired up like a brain to recognise cats.

The team built a neural network, which mimics the working of a biological brain, that worked out how to spot pictures of cats in just three days.

The cat-spotting computer was created as part of a larger project to investigate machine learning.



Millions of images were used to train the neural network

How to proceed

- **Very active field of research**
- **Very impressive results**
 - ▶ on some tasks
- Some science and lots of engineering
- I'll try to give you a flavour of the field
- ... and my opinions

How to proceed

- Very active field of research
- Very impressive results
 - ▶ on some tasks
- Some science and lots of engineering
- I'll try to give you a flavour of the field
- ... and my opinions

How to proceed

- Very active field of research
- Very impressive results
 - ▶ on some tasks
- Some science and lots of engineering
- I'll try to give you a flavour of the field
- ... and my opinions

How to proceed

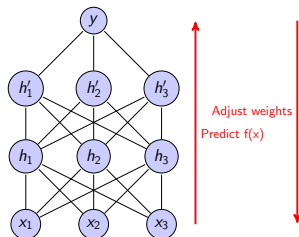
- Very active field of research
- Very impressive results
 - ▶ on some tasks
- Some science and lots of engineering
- I'll try to give you a flavour of the field
- ... and my opinions

How to proceed

- Very active field of research
- Very impressive results
 - ▶ on some tasks
- Some science and lots of engineering
- I'll try to give you a flavour of the field
- ... and my opinions

Revival of NN

- Back-prop does not handle depth
- Depth requires more data
- Restricted Boltzmann Machine
- Layer-wise training



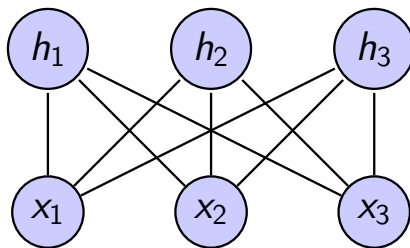
Restricted Boltzmann Machine⁴

$$p(\mathbf{x}, \mathbf{h}|\theta) = \frac{1}{Z(\theta)} \prod_r^R \prod_k^K \psi_{rk}(x_r, h_k) \quad (44)$$

- Product of Experts vs. Mixtures of Experts
 - ▶ Allows for “sharp” distributions
- $Z(\theta)$ forces normalisation
- Hidden units binary

⁴Murphy 2012, p. 27.7.

Restricted Boltzmann Machine⁴



⁴Murphy 2012, p. 27.7.

Restricted Boltzmann Machine⁴

$$p(\mathbf{h}|\mathbf{x}, \theta) = \prod_k p(h_k|\mathbf{x}, \theta) \quad (45)$$

$$p(\mathbf{x}|\mathbf{h}, \theta) = \prod_r p(x_r|\mathbf{h}, \theta) \quad (46)$$

- Variables are conditionally independent
- Learn θ using gradient based means

⁴Murphy 2012, p. 27.7.

Restricted Boltzmann Machine⁴

Binary RBM

$$p(\mathbf{x}, \mathbf{h}|\theta) = \frac{1}{Z(\theta)} e^{-E(\mathbf{x}, \mathbf{h}; \theta)} \quad (47)$$

$$E(\mathbf{x}, \mathbf{h}; \theta) = - \sum_r^R \sum_k^K x_r h_k \tilde{W}_{rk} - \sum_r^R x_r b_r - \sum_k^K h_k c_k \quad (48)$$

$$p(\mathbf{h}|\mathbf{x}, \theta) = \prod_k^K p(h_k|\mathbf{x}, \theta) = \prod_k^K \text{Ber}(h_k|\text{sigm}(\mathbf{w}_{:,k}\mathbf{x})) \quad (49)$$

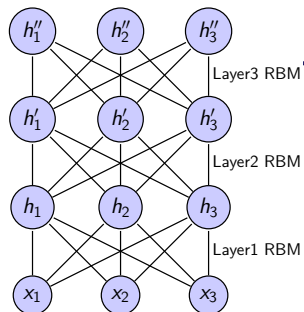
$$\mathbb{E}[\mathbf{h}|\mathbf{x}, \theta] = \text{sigm}(\mathbf{W}^T \mathbf{x}) \quad (50)$$

$$\mathbb{E}[\mathbf{x}|\mathbf{h}, \theta] = \text{sigm}(\mathbf{W} \mathbf{h}) \quad (51)$$

⁴Murphy 2012, p. 27.7.

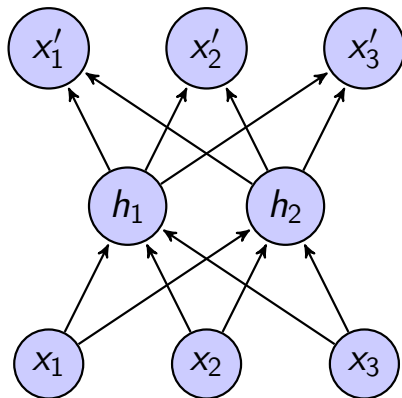
Deep Belief Networks⁵

- Stack several RBMs
- Layer-wise independence
- Each RBM works as a prior for the next level
- *“If you want to do Computer Vision first learn Computer Graphics”*



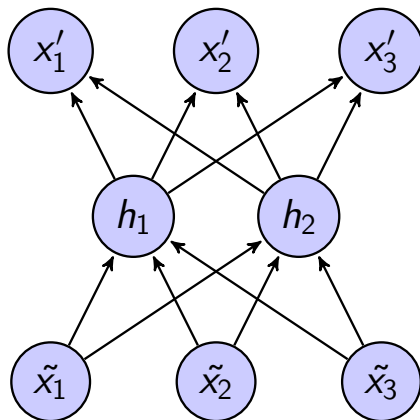
⁵Murphy 2012, p. 28.2.3.

Auto-encoders⁶



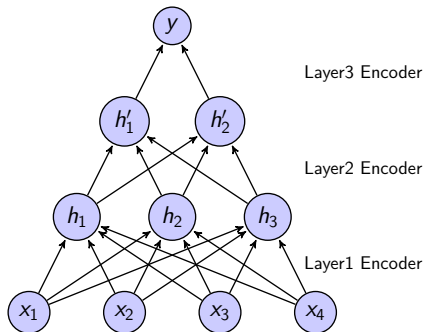
⁶Vincent *et al.* 2010.

Auto-encoders⁶



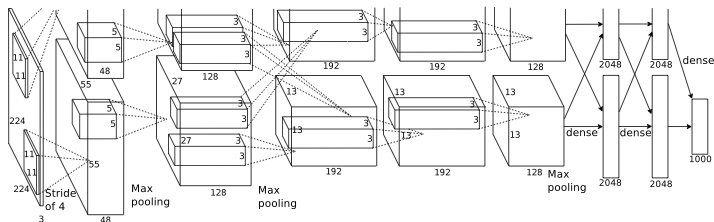
⁶Vincent *et al.* 2010.

Auto-encoders⁶



⁶Vincent *et al.* 2010.

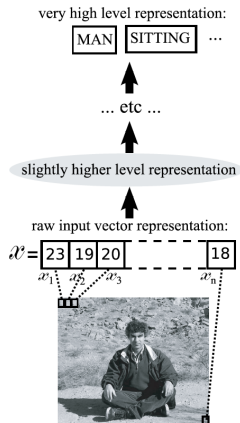
Convolutional Neural Networks⁷



Very structured architecture allows for non-layerwise training

⁷Berkely Caffe

Why⁸



⁸Bengio *et al.* 2013.

Why⁸



⁸Bengio *et al.* 2013.

Why⁸

“It’s true there’s been a lot of work on trying to apply statistical models to various linguistic problems. I think there have been some successes, but a lot of failures. There is a notion of success which I think is novel in the history of science. It interprets success as approximating unanalyzed data.”

[Noam Chomsky]

⁸Bengio *et al.* 2013.

Why⁸

Carls Rant

- These things clearly works
- The science is not to make them work but Why they work
- Quickest short-term progress is often not reached by principles
- We run the risk of disapointing a lot of people by getting lost

⁸Bengio *et al.* 2013.



Deep Gaussian Processes⁹

- Why does a probabilistic model work?
- A good model has sensible priors
- Samples from priors tells us what we prefer to model
- What are hierarchical priors?

⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

- Why does a probabilistic model work?
- A good model has sensible priors
 - Samples from priors tells us what we prefer to model
 - What are hierarchical priors?

⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

- Why does a probabilistic model work?
- A good model has sensible priors
- Samples from priors tells us what we prefer to model
- What are hierarchical priors?

⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

- Why does a probabilistic model work?
- A good model has sensible priors
- Samples from priors tells us what we prefer to model
- What are hierarchical priors?

⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

$$f(\mathbf{x}) = \frac{1}{K} \sum_i^K w_i h_i(\mathbf{x}) = \mathbf{w}^T \mathbf{h}(\mathbf{x}) \quad (52)$$

⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

$$f(\mathbf{x}) = \frac{1}{K} \sum_i^K w_i h_i(\mathbf{x}) = \mathbf{w}^T \mathbf{h}(\mathbf{x}) \quad (53)$$

$$= \mathbf{w}^T \mathbf{h}^{(2)}(\mathbf{h}^{(1)}(\mathbf{x})) \quad (54)$$

$$k_1(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{h}(\mathbf{x}_i)^T \mathbf{h}(\mathbf{x}_j) \quad (55)$$

$$k_2(\mathbf{x}_i, \mathbf{x}_j) = [\mathbf{h}^{(2)}(\mathbf{h}^{(1)}(\mathbf{x}_i))]^T \mathbf{h}^{(2)}(\mathbf{h}^{(1)}(\mathbf{x}_j)) \quad (56)$$

⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

$$f(\mathbf{x}) = \frac{1}{K} \sum_i^K w_i h_i(\mathbf{x}) = \mathbf{w}^T \mathbf{h}(\mathbf{x}) \quad (57)$$

$$= \mathbf{w}^T \mathbf{h}^{(2)}(\mathbf{h}^{(1)}(\mathbf{x})) \quad (58)$$

$$k_1(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{h}(\mathbf{x}_i)^T \mathbf{h}(\mathbf{x}_j) \quad (59)$$

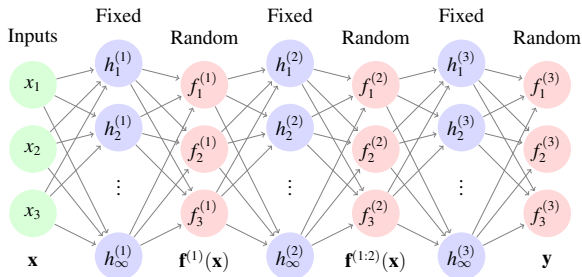
$$k_2(\mathbf{x}_i, \mathbf{x}_j) = [\mathbf{h}^{(2)}(\mathbf{h}^{(1)}(\mathbf{x}_i))]^T \mathbf{h}^{(2)}(\mathbf{h}^{(1)}(\mathbf{x}_j)) \quad (60)$$

$k(\mathbf{x}_i, \mathbf{x}_j)$ has closed form for SE kernel

$$k_{L+1}(\mathbf{x}_i, \mathbf{x}_j) = e^{k_L(\mathbf{x}_i, \mathbf{x}_j) - 1} \quad (61)$$

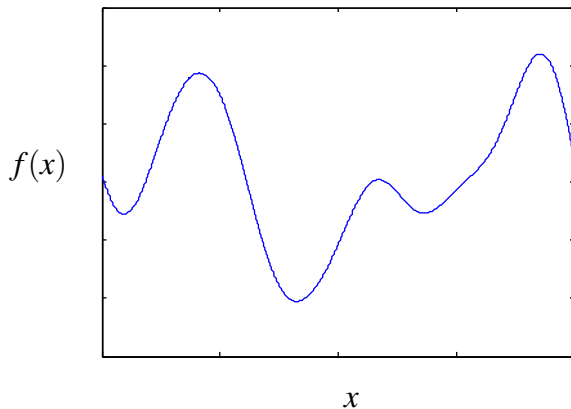
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



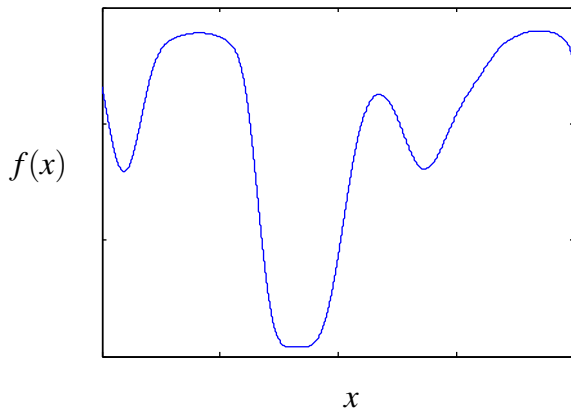
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



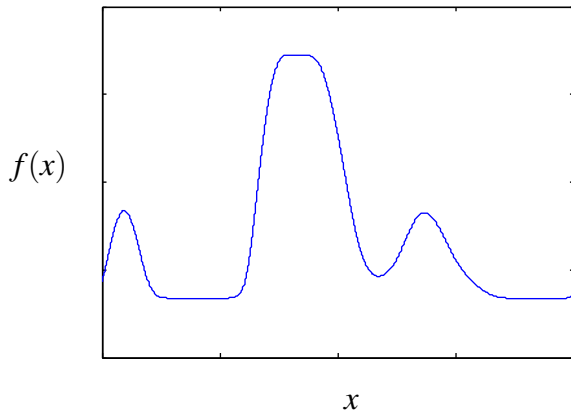
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



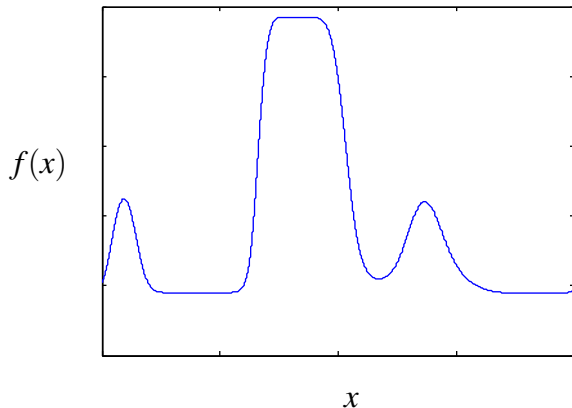
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



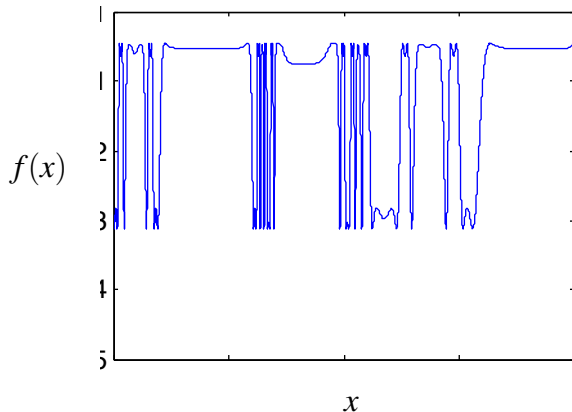
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



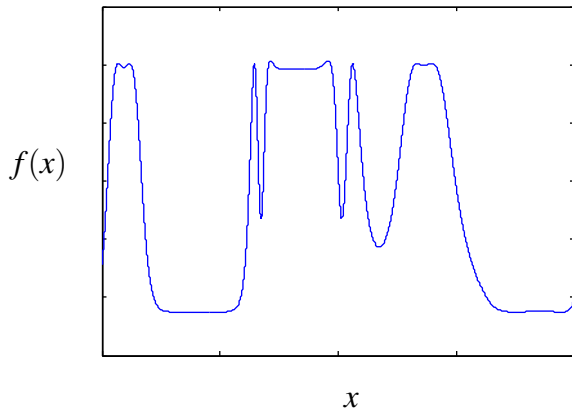
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



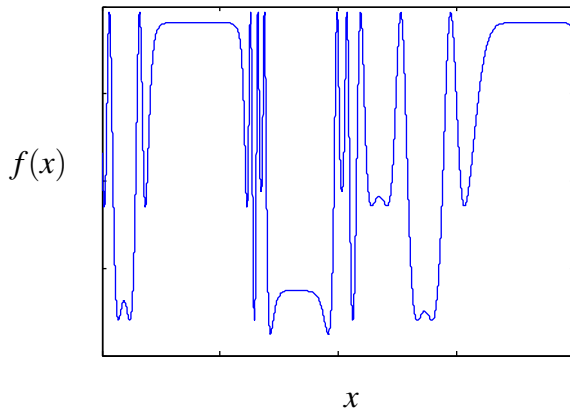
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



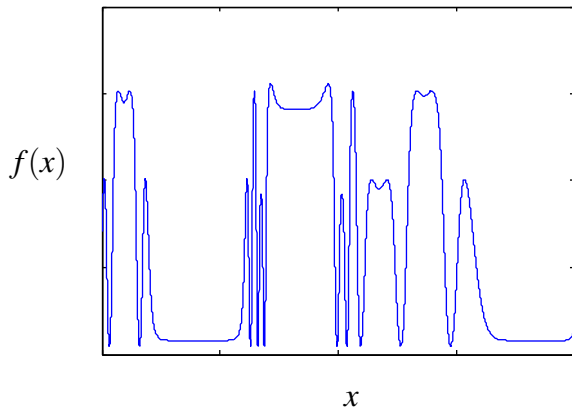
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



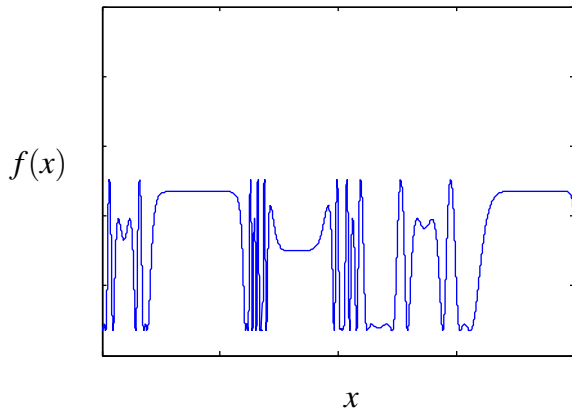
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



⁹Duvenaud *et al.* 2014.

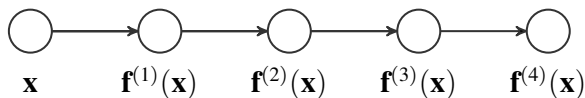
Deep Gaussian Processes⁹



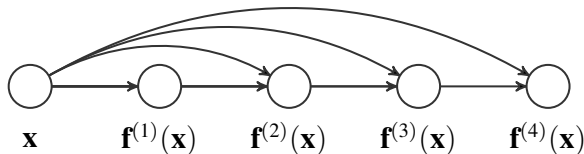
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

Standard architecture:

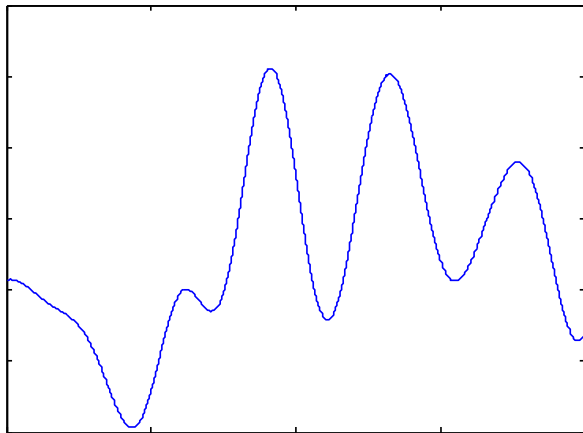


Input-connected architecture:



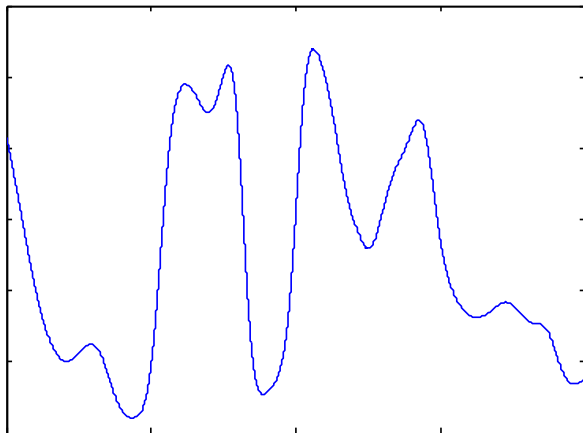
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



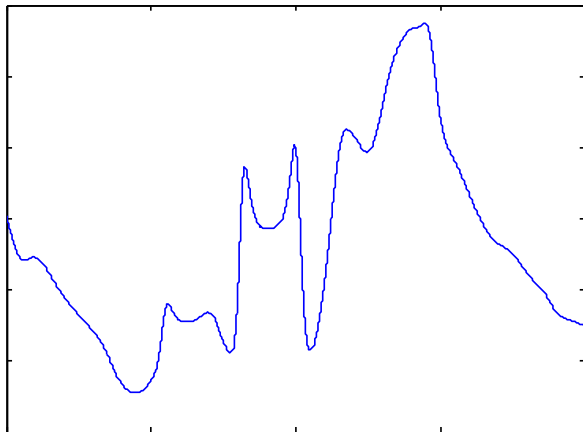
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



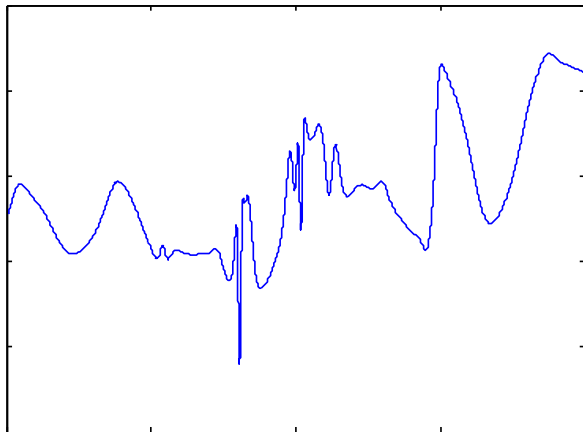
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



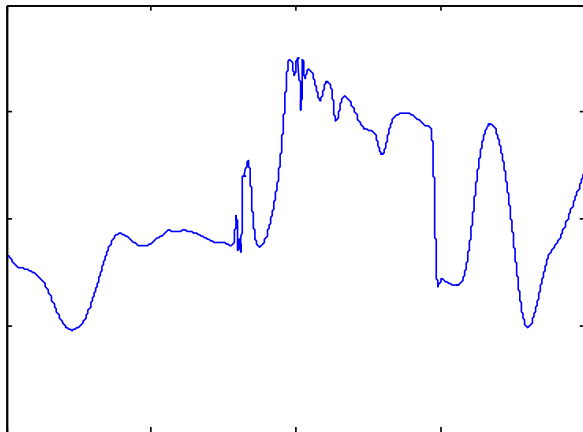
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



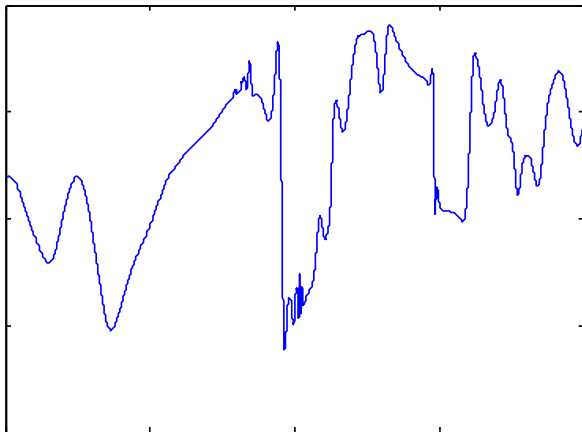
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



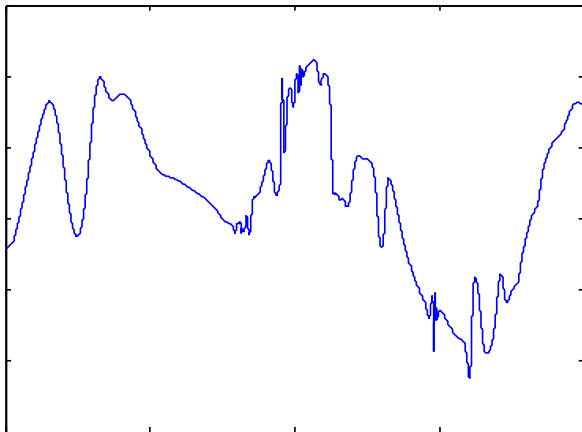
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



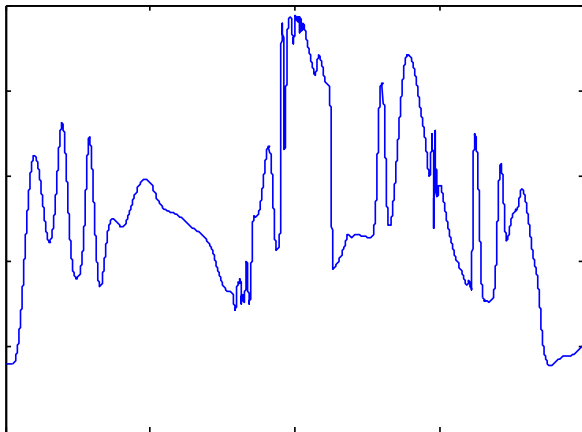
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



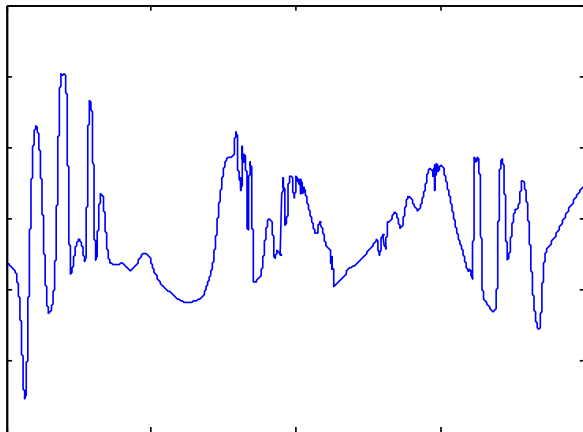
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



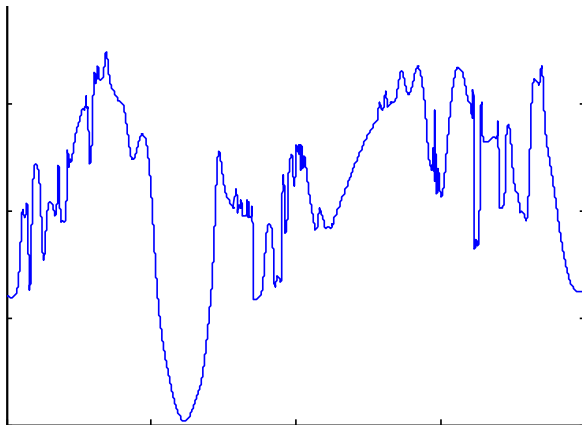
⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹



⁹Duvenaud *et al.* 2014.

Deep Gaussian Processes⁹

- Priors allows us to analyse design before seeing data
- Deep GPs shows what depth provides
 - ▶ non-stationary functions
- Allows for deep models on small data-sets
- Shed light on some current design heuristics

⁹Duvenaud *et al.* 2014.

Future

- If we have enough data we do not need priors (Laplace)
- which interesting problems do we have that for?
- no priors (or not formulated priors) makes us headless chickens
- when we need a lot of data to solve a simple problem you should be worried

Future

- If we have enough data we do not need priors (Laplace)
- which interesting problems do we have that for?
- no priors (or not formulated priors) makes us headless chickens
- when we need a lot of data to solve a simple problem you should be worried

Future

- If we have enough data we do not need priors (Laplace)
- which interesting problems do we have that for?
- no priors (or not formulated priors) makes us headless chickens
- when we need a lot of data to solve a simple problem you should be worried

Future

- If we have enough data we do not need priors (Laplace)
- which interesting problems do we have that for?
- no priors (or not formulated priors) makes us headless chickens
- when we need a lot of data to solve a simple problem you should be worried

Introduction

Recap

Hierarchical Models

Summary

End of Part 2

- Bayesian modelling
 - ▶ specify likelihood and prior
 - ▶ inference through posterior
- Strength of priors
- Sensible assumptions and approximations (MAP, ML, Variational)
- We have been very abstract on purpose to focus on understanding learning [Chomsky]

What do you need to do?

- Translate to your own problems/data
- How have you solved problems before, think of the assumptions you made
- What are sensible priors/likelihoods/structures
- What assumptions do I need to make?
- Don't be afraid of being abstract, when you get too close to the problem you often make assumptions that you are not aware of
- Get your hands dirty, i.e. develop your own priors for developing models

What do you need to do?

- Translate to your own problems/data
- How have you solved problems before, think of the assumptions you made
- What are sensible priors/likelihoods/structures
- What assumptions do I need to make?
- Don't be afraid of being abstract, when you get too close to the problem you often make assumptions that you are not aware of
- Get your hands dirty, i.e. develop your own priors for developing models

Take home message

- Machine learning is really simple, it should be as even Carl have learnt quite a few things in life
- Formulating learning so that it can be externalised might be very hard and really involved but that is just labour
- Make assumptions, lots of them, that is the basis of learning, but be aware of them

Take home message

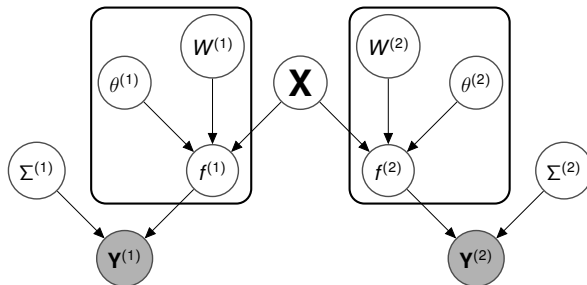
- Machine learning is really simple, it should be as even Carl have learnt quite a few things in life
- Formulating learning so that it can be externalised might be very hard and really involved but that is just labour
- Make assumptions, lots of them, that is the basis of learning, but be aware of them

Take home message



- Machine learning is really simple, it should be as even Carl have learnt quite a few things in life
- Formulating learning so that it can be externalised might be very hard and really involved but that is just labour
- Make assumptions, lots of them, that is the basis of learning, but be aware of them

e.o.f.



My Research



References I

-  **Kevin P Murphy.** *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012. ISBN: 0262018020, 9780262018029.
-  **Neil D Lawrence.** “Probabilistic non-linear principal component analysis with Gaussian process latent variable models”. In: *The Journal of Machine Learning Research* 6 (2005), pp. 1783–1816. URL: <http://dl.acm.org/citation.cfm?id=1194904>.

References II

-  **Warren S McCulloch and Walter Pitts.** “A logical calculus of the ideas immanent in nervous activity”. *English*. In: *The Bulletin of Mathematical Biophysics* 5.4 (Dec. 1943), pp. 115–133. DOI: [10.1007/BF02478259](https://doi.org/10.1007/BF02478259). URL: <http://link.springer.com/10.1007/BF02478259>.
-  **F Rosenblatt.** “The perceptron: a probabilistic model for information storage and organization in the brain”. In: *Psychology Review* (Nov. 1958), pp. 386–408. URL: <http://www.ncbi.nlm.nih.gov/pubmed/13602029>.

References III





Marvin Minsky and Seymour Papert. “Perceptrons. An Introduction to Computational Geometry.” English. In: *Science* 165.3895 (Aug. 1969), pp. 780–782. DOI: [10.1126/science.165.3895.780](https://doi.org/10.1126/science.165.3895.780). URL: <http://www.sciencemag.org/cgi/doi/10.1126/science.165.3895.780>.





D E Rumelhart *et al.* “Learning representations by back-propagating errors”. In: *Nature* 323.9 (Oct. 1986), pp. 533–536. URL: http://www.iro.umontreal.ca/~pift6266/A06/refs/backprop_old.pdf.

References IV

-  **Geoffrey E Hinton *et al.*** “A Fast Learning Algorithm for Deep Belief Nets”. English. In: *Neural Computation* 18.7 (July 2006), pp. 1527–1554. DOI: [10.1162/jmlr.2003.4.7-8.1235](https://doi.org/10.1162/jmlr.2003.4.7-8.1235). URL: <http://www.mitpressjournals.org/doi/abs/10.1162/neco.2006.18.7.1527>.
-  **Pascal Vincent *et al.*** “Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion”. In: *The Journal of Machine Learning Research* 11 (Mar. 2010), pp. 3371–3408. URL: <http://dl.acm.org/citation.cfm?id=1756006.1953039>.

References V

-  **Yoshua Bengio *et al.*** “Representation learning: A review and new perspectives”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35 (Aug. 2013), pp. 1798–1828. ISSN: 0162-8828. URL: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=6472238.
-  **David Duvenaud *et al.*** *Avoiding pathologies in very deep networks*. 2014. URL: <http://jmlr.org/proceedings/papers/v33/duvenaud14.pdf>.

Appendix

Similar Matrices: Self-Similarity

$$\mathbf{A} = \mathbf{I} \mathbf{A} \mathbf{I}^{-1} = \mathbf{I}^{-1} \mathbf{A} \mathbf{I}$$

◀ Return

Similar Matrices: Symmetry

$$\begin{aligned} \mathbf{A} \sim \mathbf{B} &\Rightarrow \mathbf{B} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P} \\ \det \mathbf{B} &= \det(\mathbf{P}^{-1}\mathbf{A}\mathbf{P}) = \det(\mathbf{P}^{-1})\det(\mathbf{A})\det(\mathbf{P}) = \\ &= \det(\mathbf{A})\det(\mathbf{P}^{-1})\det(\mathbf{P}) = \det(\mathbf{A})\frac{1}{\det(\mathbf{P})}\det(\mathbf{P}) = \\ &\det(\mathbf{B}) \end{aligned}$$

[◀ Return](#)

Similar Matrices: Trace

$$\begin{aligned} \mathbf{A} \quad \sim \quad \mathbf{B} &\Rightarrow \mathbf{B} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P} \\ \text{trace}(\mathbf{B}) &= \text{trace}(\mathbf{P}^{-1}\mathbf{A}\mathbf{P}) = \{\text{trace}(\mathbf{A}\mathbf{P}) = \text{trace}(\mathbf{A}\mathbf{P})\} = \\ &= \text{trace} \left(\left(\mathbf{P}\mathbf{P}^{-1} \right) \mathbf{A} \right) = \text{trace}(\mathbf{A}) \end{aligned}$$

[◀ Return](#)

Similar Matrices: Power

$$\begin{aligned}
 \mathbf{A} &\sim \mathbf{B} \Rightarrow \mathbf{B} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P} \\
 \mathbf{B}^2 &= (\mathbf{P}^{-1}\mathbf{A}\mathbf{P})^2 = (\mathbf{P}^{-1}\mathbf{A}\mathbf{P})(\mathbf{P}^{-1}\mathbf{A}\mathbf{P}) = \\
 &= (\mathbf{P}^{-1}\mathbf{A}) \left(\underbrace{\mathbf{P}\mathbf{P}^{-1}}_{=\mathbf{I}} \right) (\mathbf{A}\mathbf{P}) = \\
 &= \mathbf{P}^{-1}\mathbf{A}\mathbf{A}\mathbf{P} = \mathbf{P}^{-1}\mathbf{A}^2\mathbf{P}
 \end{aligned}$$

Prove further powers by induction over exponent

[◀ Return](#)

Similar Matrices: Invertability

$$\begin{aligned}\mathbf{A} &\sim \mathbf{B} \Rightarrow \mathbf{B} = \mathbf{P}^{-1}\mathbf{A}\mathbf{P} \\ &\Rightarrow \det(\mathbf{A}) = \det(\mathbf{B})\end{aligned}$$

\mathbf{A}^{-1} Exists if $\det(\mathbf{A}) \neq 0$

$$\det(\mathbf{B}) \neq 0 \iff \det(\mathbf{A}) \neq 0$$

[◀ Return](#)

$$\begin{aligned} \mathbf{A}_{ij} &= \sum_{k=1}^N \mathbf{v}_{ik} \mathbf{D}_{kk} (\mathbf{v}^T)_{kj} = \sum_{k=1}^N (\mathbf{v}_k)_i \lambda_k (\mathbf{v}_k)_j \\ &= \sum_{k=1}^N \left(\lambda_k \mathbf{v}_k \mathbf{v}_k^T \right)_{ij} \end{aligned}$$

[◀ Return](#)

Rank Approximation

$$\begin{aligned}
 \|\mathbf{A} - \mathbf{B}\|_F &= \left\| \sum_{i=1}^N \lambda_i \mathbf{v}_i \mathbf{v}_i^T - \sum_{i=1}^N q_i \mathbf{v}_i \mathbf{v}_i^T \right\|_F = \\
 &= \left\| \sum_{i=1}^N (\lambda_i - q_i) \mathbf{v}_i \mathbf{v}_i^T \right\| = \\
 &= \left\{ \left((\lambda_i - q_i) \mathbf{v}_i \underbrace{\mathbf{v}_i^T \mathbf{v}_i}_{=1} \right) \mathbf{v}_i = (\lambda_i - q_i) \mathbf{v}_i \right\} = \\
 &= \sqrt{\sum_{i=1}^N (\lambda_i - q_i)^2} \quad \leftarrow \text{Return}
 \end{aligned}$$

Multidimensional Scaling

Define:

$$d_{ij}^2 = \sum_{k=1}^d (x_{ki} - x_{kj})^2 = \mathbf{x}_i^T \mathbf{x}_i + \mathbf{x}_j^T \mathbf{x}_j - 2\mathbf{x}_i \mathbf{x}_j$$

$$g_{ij} = \sum_{k=1}^d x_{ki} x_{kj} = \mathbf{x}_i^T \mathbf{x}_j$$

$$\Rightarrow d_{ij}^2 = g_{ii} + g_{jj} - 2g_{ij}$$

Centering:
$$\sum_{i=1}^N g_{ij} = \sum_{i=1}^N \mathbf{x}_i^T \mathbf{x}_j = \underbrace{\left(\sum_{i=1}^N \mathbf{x}_i^T \right)}_{=0} \mathbf{x}_j = 0$$

Multidimensional Scaling

Want to Express **G** in terms of **D**

$$g_{ij} = \frac{1}{2}(g_{ii} + g_{jj} - d_{ij}^2)$$

$$\frac{1}{N} \sum_{i=1}^N d_{ij}^2 = g_{jj} + \frac{1}{N} \sum_{i=1}^N g_{ii}$$

$$\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = \frac{2}{N} \sum_{i=1}^N g_{ii}$$

$$\Rightarrow g_{ij} = \frac{1}{2} \left(\frac{1}{N} \left(\sum_{k=1}^N d_{kj}^2 + \sum_{k=1}^N d_{ik}^2 - \frac{1}{N} \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 \right) - d_{ij}^2 \right)$$

◀ Return: MDS

◀ Return: MVU

PCA MDS Equivalence

$$\begin{aligned}
 \mathbf{G} &= \mathbf{X}\mathbf{X}^T = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T \\
 \Rightarrow & (\mathbf{X}\mathbf{X}^T)\mathbf{v}_i = \lambda_i\mathbf{v}_i \\
 \Rightarrow & \frac{1}{N-1}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)\mathbf{v}_i = \lambda_i\frac{1}{N-1}\mathbf{X}^T\mathbf{v}_i \\
 \Rightarrow & \underbrace{\frac{1}{N-1}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)}_{\mathbf{S}}\mathbf{v}_i = \lambda_i\frac{1}{N-1}\mathbf{X}^T\mathbf{v}_i \\
 \Rightarrow & \mathbf{S} \underbrace{(\mathbf{X}^T\mathbf{v}_i)}_{\text{eigenvectors?}} = \underbrace{\frac{\lambda_i}{N-1}}_{\text{eigenvalue?}} \underbrace{(\mathbf{X}^T\mathbf{v}_i)}_{\text{eigenvector?}}
 \end{aligned}$$

PCA MDS Equivalence

Enforce orthogonality

$$\begin{aligned}(\mathbf{X}^T \mathbf{v}_i)^T (\mathbf{X}^T \mathbf{v}_i) &= \mathbf{v}_i^T \mathbf{X} \mathbf{X}^T \mathbf{v}_i = \lambda_i \\ \Rightarrow \frac{1}{\sqrt{\lambda_i}} \mathbf{v}_i^T \mathbf{X} \mathbf{X}^T \mathbf{v}_i \frac{1}{\sqrt{\lambda_i}} &= \left(\frac{1}{\sqrt{\lambda_i}} \right)^2 \lambda_i = 1 \\ (\mathbf{X}^T \mathbf{v}_i) \frac{1}{\sqrt{\lambda_i}})^T (\mathbf{X}^T \mathbf{v}_i \frac{1}{\sqrt{\lambda_i}}) &= 1\end{aligned}$$

PCA MDS Equivalence

$$\begin{aligned}
 \text{Define: } \mathbf{v}_i^{\text{PCA}} &= \mathbf{X}^T \mathbf{v}_i \frac{1}{\sqrt{\lambda_i}} \\
 \mathbf{y}_i^{\text{PCA}} &= \mathbf{X} \mathbf{v}_i^{\text{PCA}} = \mathbf{X} \mathbf{X}^T \mathbf{v}_i \frac{1}{\sqrt{\lambda_i}} = \\
 &= \lambda_i \mathbf{v}_i \frac{1}{\sqrt{\lambda_i}} = \sqrt{\lambda_i} \mathbf{v}_i \\
 \mathbf{y}_i^{\text{MDS}} &= \mathbf{v}_i \sqrt{\lambda_i} = \sqrt{\lambda_i} \mathbf{v}_i \\
 \Rightarrow \mathbf{y}_i^{\text{PCA}} &= \mathbf{y}_i^{\text{MDS}}
 \end{aligned}$$

← PCA

Maximum Variance Unfolding: Objective

$$\begin{aligned}
 \sum_{i=1}^N g_{ii} &= \sum_{i=1}^N \frac{1}{2} \left(\frac{1}{N} \left(\sum_{k=1}^N d_{kj}^2 + \sum_{k=1}^N d_{ik}^2 - \frac{1}{N} \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 \right) - d_{ii}^2 \right) = \\
 &= \underbrace{\frac{1}{2N} \sum_{i=1}^N \sum_{k=1}^N d_{ki}^2 + \frac{1}{2N} \sum_{i=1}^N \sum_{k=1}^N d_{ik}^2}_{\text{symmetry} = \frac{1}{2N} 2 \sum_{i=1}^N \sum_{k=1}^N d_{ki}^2} - \\
 &- \frac{1}{2N^2} N \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 - \frac{1}{2} \sum_i \underbrace{d_{ii}^2}_{=0} =
 \end{aligned}$$

Maximum Variance Unfolding: Objective

$$\begin{aligned}
 &= \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^N d_{ki}^2 - \frac{1}{2N} \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 = \\
 &= \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 \\
 \text{trace}(\mathbf{G}) &= \sum_{i=1}^N g_{ii} = \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = \\
 &= \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N \|\mathbf{y}_i - \mathbf{y}_j\|_{L_2}^2
 \end{aligned}$$

◀ Return

Maximum Variance Unfolding: Centering

$$\begin{aligned}
 \sum_{i=1}^N \sum_{j=1}^N g_{ij} &= \sum_{i=1}^N \sum_{j=1}^N \frac{1}{2} \left(\frac{1}{N} \left(\sum_{k=1}^N d_{kj}^2 + \sum_{k=1}^N d_{ik}^2 - \right. \right. \\
 &\quad \left. \left. - \frac{1}{N} \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 \right) - d_{ij}^2 \right) = \\
 &= \frac{1}{2N} \underbrace{\sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N d_{kj}^2}_{=N \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2} + \frac{1}{2N} \underbrace{\sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N d_{ik}^2}_{=N \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2} -
 \end{aligned}$$

Maximum Variance Unfolding: Centering

$$\begin{aligned}
 & - \frac{1}{2N^2} \underbrace{\sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2}_{=N^2 \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = \\
 & = \underbrace{\left(\frac{1}{2} + \frac{1}{2} - \frac{1}{2} - \frac{1}{2} \right)}_{=0} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = 0 \\
 \left\| \sum_{i=1}^N \mathbf{y}_i \right\|_{L_2}^2 & \Rightarrow \sum_{i=1}^N \sum_{j=1}^N \mathbf{K}_{ij} = 0
 \end{aligned}$$

[← Return](#)

Spectral Theorem

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \quad \mathbf{A} = \mathbf{V} \Delta \mathbf{V}^T, \quad \|\mathbf{x}\|_{L2} = 1$$

$$\mathbf{x} = \mathbf{1} \sum_{i=1}^N \alpha_i \mathbf{v}_i$$

$$\|\alpha\| = 1$$

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \left(\sum_{i=1}^N \alpha_i \mathbf{v}_i \right)^T \mathbf{A} \left(\sum_{i=1}^N \alpha_i \mathbf{v}_i \right) =$$

$$= \left(\sum_{i=1}^N \alpha_i \mathbf{v}_i \right)^T \left(\sum_{i=1}^N \lambda_i \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{i=1}^N \alpha_i \mathbf{v}_i \right) =$$

Spectral Theorem

$$\begin{aligned}
 &= \left(\sum_{i=1}^N \alpha_i \mathbf{v}_i \right)^T \left(\sum_{i=1}^N \lambda_i \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{i=1}^N \alpha_i \mathbf{v}_i \right) = \\
 &= \left\{ \mathbf{v}_i^T \mathbf{v}_j = \begin{cases} 1 & i=j \\ 0 & \text{otherwise} \end{cases} \right\} = \\
 &= \sum_{i=1}^N \alpha_i^2 \lambda_i \underbrace{\mathbf{v}_i^T \mathbf{v}_i}_{=1} \underbrace{\mathbf{v}_i^T \mathbf{v}_i}_{=1} = \\
 &= \sum_{i=1}^N \alpha_i^2 \lambda_i \begin{cases} \max & : \mathbf{x}^T \mathbf{A} \mathbf{x} = \lambda_1 & \mathbf{x} = \mathbf{v}_1 \\ \min & : \mathbf{x}^T \mathbf{A} \mathbf{x} = \lambda_N & \mathbf{x} = \mathbf{v}_N \end{cases}
 \end{aligned}$$

[← Return LLE](#)
[← Return Laplacian](#)

Maximum Variance Unfolding: Objective

$$\begin{aligned}
 \sum_{i=1}^N g_{ii} &= \sum_{i=1}^N \frac{1}{2} \left(\frac{1}{N} \left(\sum_{k=1}^N d_{kj}^2 + \sum_{k=1}^N d_{ik}^2 - \frac{1}{N} \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 \right) - d_{ii}^2 \right) = \\
 &= \underbrace{\frac{1}{2N} \sum_{i=1}^N \sum_{k=1}^N d_{ki}^2 + \frac{1}{2N} \sum_{i=1}^N \sum_{k=1}^N d_{ik}^2}_{\text{symmetry} = \frac{1}{2N} 2 \sum_{i=1}^N \sum_{k=1}^N d_{ki}^2} - \\
 &- \frac{1}{2N^2} N \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 - \frac{1}{2} \sum_i \underbrace{d_{ii}^2}_{=0} =
 \end{aligned}$$

Maximum Variance Unfolding: Objective

$$\begin{aligned}
 &= \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^N d_{ki}^2 - \frac{1}{2N} \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 = \\
 &= \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 \\
 \text{trace}(\mathbf{G}) &= \sum_{i=1}^N g_{ii} = \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = \\
 &= \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N \|\mathbf{y}_i - \mathbf{y}_j\|_{L_2}^2
 \end{aligned}$$

◀ Return

Maximum Variance Unfolding: Centering

$$\begin{aligned}
 \sum_{i=1}^N \sum_{j=1}^N g_{ij} &= \sum_{i=1}^N \sum_{j=1}^N \frac{1}{2} \left(\frac{1}{N} \left(\sum_{k=1}^N d_{kj}^2 + \sum_{k=1}^N d_{ik}^2 - \right. \right. \\
 &\quad \left. \left. - \frac{1}{N} \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2 \right) - d_{ij}^2 \right) = \\
 &= \frac{1}{2N} \underbrace{\sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N d_{kj}^2}_{=N \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2} + \frac{1}{2N} \underbrace{\sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N d_{ik}^2}_{=N \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2} -
 \end{aligned}$$

Maximum Variance Unfolding: Centering

$$\begin{aligned}
 & - \frac{1}{2N^2} \underbrace{\sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N \sum_{p=1}^N d_{kp}^2}_{=N^2 \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2} - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = \\
 & = \underbrace{\left(\frac{1}{2} + \frac{1}{2} - \frac{1}{2} - \frac{1}{2} \right)}_{=0} \sum_{i=1}^N \sum_{j=1}^N d_{ij}^2 = 0 \\
 \left\| \sum_{i=1}^N \mathbf{y}_i \right\|_{L2}^2 & \Rightarrow \sum_{i=1}^N \sum_{j=1}^N \mathbf{K}_{ij} = 0
 \end{aligned}$$

◀ Return