# Probabilistic Dimensionality Reduction

Neil D. Lawrence

Amazon Research Cambridge and University of Sheffield, U.K.

Probabilistic Scientific Computing Workshop
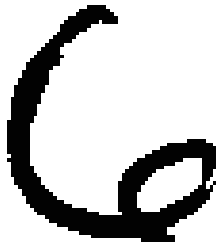ICERM at Brown

6th June 2017

# Outline

**USPS Data Set Handwritten Digit**

- 3648 Dimensions
  - 64 rows by 57 columns

**USPS Data Set Handwritten Digit**

- 3648 Dimensions
  - 64 rows by 57 columns
  - Space contains more than just this digit.

**USPS Data Set Handwritten Digit**

- 3648 Dimensions
    - 64 rows by 57 columns
    - Space contains more than just this digit.
    - Even if we sample every nanosecond from now until the end of the universe, you won't see the original six!

**USPS Data Set Handwritten Digit**

- 3648 Dimensions
    - 64 rows by 57 columns
    - Space contains more than just this digit.
    - Even if we sample every nanosecond from now until the end of the universe, you won't see the original six!

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

**Rotate a 'Prototype'**

# MATLAB Demo

```
demDigitsManifold([1 2], 'all')
```

# MATLAB Demo

```
demDigitsManifold([1 2], 'all')
```

# MATLAB Demo

```
demDigitsManifold([1 2], 'sixnine')
```

# Low Dimensional Manifolds

**Pure Rotation is too Simple**

- In practice the data may undergo several distortions.
  - *e.g.* digits undergo 'thinning', translation and rotation.
- For data with 'structure':
  - we expect fewer distortions than dimensions;
  - we therefore expect the data to live on a lower dimensional manifold.
- Conclusion: deal with high dimensional data by looking for lower dimensional non-linear embedding.

# Existing Methods

**Spectral Approaches**

- ▶ Classical Multidimensional Scaling (MDS) (Mardia et al., 1979).
  - ▶ Uses eigenvectors of similarity matrix.
    - ▶ Isomap (Tenenbaum et al., 2000) is MDS with a particular proximity measure.
  - ▶ Kernel PCA (Schölkopf et al., 1998)
    - ▶ Provides a representation and a mapping — dimensional expansion.
    - ▶ Mapping is implied throught he use of a kernel function as a similarity matrix.
  - ▶ Locally Linear Embedding (Roweis and Saul, 2000).
    - ▶ Looks to preserve locally linear relationships in a low dimensional space.

**Iterative Methods**

- Multidimensional Scaling (MDS)
    - Iterative optimisation of a stress function (Kruskal, 1964).
    - Sammon Mappings (Sammon, 1969).
        - Strictly speaking not a mapping — similar to iterative MDS.
- NeuroScale (Lowe and Tipping, 1997)
    - Augmentation of iterative MDS methods with a mapping.

**Probabilistic Approaches**

- Probabilistic PCA (Tipping and Bishop, 1999; Roweis, 1998)
  - A linear method.

**Probabilistic Approaches**

- ▶ Probabilistic PCA (Tipping and Bishop, 1999; Roweis, 1998)
  - ▸ A linear method.
- ▶ Density Networks (MacKay, 1995)
  - ▸ Use importance sampling and a multi-layer perceptron.

**Probabilistic Approaches**

- Probabilistic PCA (Tipping and Bishop, 1999; Roweis, 1998)
    - A linear method.
- Density Networks (MacKay, 1995)
    - Use importance sampling and a multi-layer perceptron.
- Generative Topographic Mapping (GTM) (Bishop et al., 1998)
    - Uses a grid based sample and an RBF network.

**Probabilistic Approaches**

- ▶ Probabilistic PCA (Tipping and Bishop, 1999; Roweis, 1998)
  - ▶ A linear method.
- ▶ Density Networks (MacKay, 1995)
  - ▶ Use importance sampling and a multi-layer perceptron.
- ▶ Generative Topographic Mapping (GTM) (Bishop et al., 1998)
  - ▶ Uses a grid based sample and an RBF network.

**Difficulty for Probabilistic Approaches**

- ▶ Propagate a probability distribution through a non-linear mapping.

# The New Model

**A Probabilistic Non-linear PCA**

- ▶ PCA has a probabilistic interpretation (Tipping and Bishop, 1999; Roweis, 1998).
- ▶ It is difficult to 'non-linearise'.

**Dual Probabilistic PCA**

- ▶ We present a new probabilistic interpretation of PCA (Lawrence, 2005).
- ▶ This interpretation can be made non-linear.
- ▶ The result is non-linear probabilistic PCA.

## Notation

$q$— dimension of latent/embedded space
$p$— dimension of data space
$n$— number of data points

centred data, $\mathbf{Y} = [\mathbf{y}_{1,:}, \ldots, \mathbf{y}_{n,:}]^\top = \left[\mathbf{y}_{:,1}, \ldots, \mathbf{y}_{:,p}\right] \in \mathfrak{R}^{n \times p}$
latent variables, $\mathbf{X} = [\mathbf{x}_{1,:}, \ldots, \mathbf{x}_{n,:}]^\top = \left[\mathbf{x}_{:,1}, \ldots, \mathbf{x}_{:,q}\right] \in \mathfrak{R}^{n \times q}$
mapping matrix, $\mathbf{W} \in \mathfrak{R}^{p \times q}$

$\mathbf{a}_{i,:}$ is a vector from the $i$th row of a given matrix $\mathbf{A}$
$\mathbf{a}_{:,j}$ is a vector from the $j$th row of a given matrix $\mathbf{A}$

**X and Y are** *design matrices*

- Covariance given by $n^{-1}\mathbf{Y}^\top\mathbf{Y}$.
- Inner product matrix given by $\mathbf{Y}\mathbf{Y}^\top$.

**Linear Latent Variable Model**

- Represent data, $\mathbf{Y}$, with a lower dimensional set of latent variables $\mathbf{X}$.
- Assume a linear relationship of the form

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\epsilon}_{i,:},$$

where

$$\boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2\mathbf{I}\right).$$

# Linear Latent Variable Model

**Probabilistic PCA**

▶ Define *linear-Gaussian relationship* between latent variables and data.



$$p\left(\mathbf{Y}|\mathbf{X},\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2\mathbf{I}\right)$$

# Linear Latent Variable Model

**Probabilistic PCA**

- Define *linear-Gaussian relationship* between latent variables and data.
- **Standard** Latent variable approach:



$$p\left(\mathbf{Y}|\mathbf{X},\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I}\right)$$

# Linear Latent Variable Model

**Probabilistic PCA**

- Define *linear-Gaussian relationship* between latent variables and data.
- **Standard** Latent variable approach:
  - Define Gaussian prior over *latent space*, **X**.



$$p(\mathbf{Y}|\mathbf{X}, \mathbf{W}) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2\mathbf{I}\right)$$

$$p(\mathbf{X}) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{x}_{i,:}|\mathbf{0}, \mathbf{I}\right)$$

# Linear Latent Variable Model

**Probabilistic PCA**

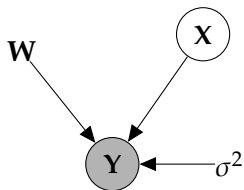- Define *linear-Gaussian relationship* between latent variables and data.

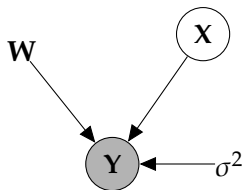- **Standard** Latent variable approach:
  - Define Gaussian prior over *latent space*, $\mathbf{X}$.
  - Integrate out *latent variables*.



$$p\left(\mathbf{Y}|\mathbf{X},\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2\mathbf{I}\right)$$

$$p\left(\mathbf{X}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{x}_{i,:}|\mathbf{0}, \mathbf{I}\right)$$

$$p\left(\mathbf{Y}|\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{0}, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}\right)$$

# Computation of the Marginal Likelihood

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\epsilon}_{i,:}, \quad \mathbf{x}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}\right), \quad \boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2 \mathbf{I}\right)$$

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\epsilon}_{i,:}, \quad \mathbf{x}_{i,:} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2 \mathbf{I}\right)$$

$$\mathbf{W}\mathbf{x}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{W}\mathbf{W}^\top\right),$$

# Computation of the Marginal Likelihood

$$\mathbf{y}_{i,:} = \mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\epsilon}_{i,:}, \quad \mathbf{x}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}\right), \quad \boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2 \mathbf{I}\right)$$

$$\mathbf{W}\mathbf{x}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{W}\mathbf{W}^\top\right),$$

$$\mathbf{W}\mathbf{x}_{i,:} + \boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}\right)$$

# Linear Latent Variable Model II

**Probabilistic PCA Max. Likelihood Soln** (Tipping and Bishop, 1999)



$$p\left(\mathbf{Y}|\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{0}, \mathbf{W}\mathbf{W}^{\top} + \sigma^2\mathbf{I}\right)$$

# Linear Latent Variable Model II

**Probabilistic PCA Max. Likelihood Soln** (Tipping and Bishop, 1999)

$$p\left(\mathbf{Y}|\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{0}, \mathbf{C}\right), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^{\top} + \sigma^2 \mathbf{I}$$

# Linear Latent Variable Model II

**Probabilistic PCA Max. Likelihood Soln** (Tipping and Bishop, 1999)

$$p\left(\mathbf{Y}|\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{0}, \mathbf{C}\right), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^{\top} + \sigma^2\mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{W}\right) = -\frac{n}{2}\log|\mathbf{C}| - \frac{1}{2}\mathrm{tr}\left(\mathbf{C}^{-1}\mathbf{Y}^{\top}\mathbf{Y}\right) + \mathrm{const.}$$

# Linear Latent Variable Model II

**Probabilistic PCA Max. Likelihood Soln** (Tipping and Bishop, 1999)

$$p\left(\mathbf{Y}|\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{0}, \mathbf{C}\right), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^{\top} + \sigma^2\mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{W}\right) = -\frac{n}{2}\log|\mathbf{C}| - \frac{1}{2}\mathrm{tr}\left(\mathbf{C}^{-1}\mathbf{Y}^{\top}\mathbf{Y}\right) + \text{const.}$$

If $\mathbf{U}_q$ are first $q$ principal eigenvectors of $n^{-1}\mathbf{Y}^{\top}\mathbf{Y}$ and the corresponding eigenvalues are $\mathbf{\Lambda}_q$,

# Linear Latent Variable Model II

**Probabilistic PCA Max. Likelihood Soln** (Tipping and Bishop, 1999)

$$p\left(\mathbf{Y}|\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{0}, \mathbf{C}\right), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^{\top} + \sigma^2\mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{W}\right) = -\frac{n}{2}\log|\mathbf{C}| - \frac{1}{2}\text{tr}\left(\mathbf{C}^{-1}\mathbf{Y}^{\top}\mathbf{Y}\right) + \text{const.}$$

If $\mathbf{U}_q$ are first $q$ principal eigenvectors of $n^{-1}\mathbf{Y}^{\top}\mathbf{Y}$ and the corresponding eigenvalues are $\mathbf{\Lambda}_q$,
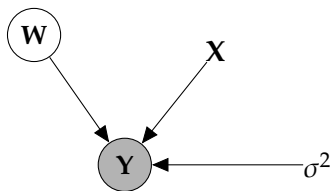
$$\mathbf{W} = \mathbf{U}_q\mathbf{L}\mathbf{R}^{\top}, \quad \mathbf{L} = \left(\mathbf{\Lambda}_q - \sigma^2\mathbf{I}\right)^{\frac{1}{2}}$$

where $\mathbf{R}$ is an arbitrary rotation matrix.

# Linear Latent Variable Model III

**Dual Probabilistic PCA**

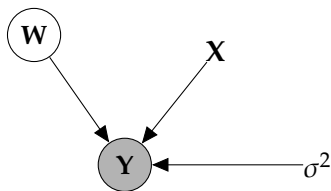- Define *linear-Gaussian relationship* between latent variables and data.



$$p\left(\mathbf{Y}|\mathbf{X}, \mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2 \mathbf{I}\right)$$

# Linear Latent Variable Model III

**Dual Probabilistic PCA**

- Define *linear-Gaussian relationship* between latent variables and data.
- **Novel** Latent variable approach:



$$p\left(\mathbf{Y}|\mathbf{X}, \mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2\mathbf{I}\right)$$

# Linear Latent Variable Model III

**Dual Probabilistic PCA**

- Define *linear-Gaussian relationship* between latent variables and data.

- **Novel** Latent variable approach:
    - Define Gaussian prior over *parameters*, **W**.



$$p\left(\mathbf{Y}|\mathbf{X},\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2\mathbf{I}\right)$$

$$p\left(\mathbf{W}\right) = \prod_{i=1}^{p} \mathcal{N}\left(\mathbf{w}_{i,:}|\mathbf{0},\mathbf{I}\right)$$

# Linear Latent Variable Model III

**Dual Probabilistic PCA**

- Define *linear-Gaussian relationship* between latent variables and data.

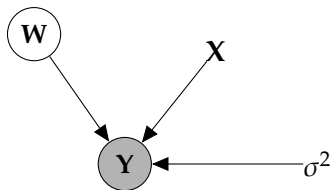- **Novel** Latent variable approach:
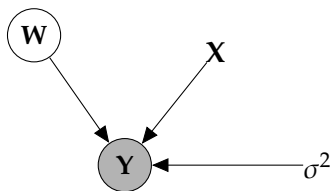  - Define Gaussian prior over *parameters*, $\mathbf{W}$.
  - Integrate out *parameters*.



$$p\left(\mathbf{Y}|\mathbf{X}, \mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{W}\mathbf{x}_{i,:}, \sigma^2\mathbf{I}\right)$$

$$p\left(\mathbf{W}\right) = \prod_{i=1}^{p} \mathcal{N}\left(\mathbf{w}_{i,:}|\mathbf{0}, \mathbf{I}\right)$$

$$p\left(\mathbf{Y}|\mathbf{X}\right) = \prod_{j=1}^{p} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \mathbf{X}\mathbf{X}^{\top} + \sigma^2\mathbf{I}\right)$$

# Computation of the Marginal Likelihood

$$\mathbf{y}_{:,j} = \mathbf{X}\mathbf{w}_{:,j} + \boldsymbol{\epsilon}_{:,j}, \quad \mathbf{w}_{:,j} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}\right), \quad \boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2 \mathbf{I}\right)$$

## Computation of the Marginal Likelihood

$$\mathbf{y}_{:,j} = \mathbf{X}\mathbf{w}_{:,j} + \boldsymbol{\epsilon}_{:,j}, \quad \mathbf{w}_{:,j} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}\right), \quad \boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2\mathbf{I}\right)$$

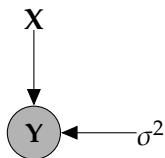$$\mathbf{X}\mathbf{w}_{:,j} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{X}\mathbf{X}^\top\right),$$

$$\mathbf{y}_{:,j} = \mathbf{X}\mathbf{w}_{:,j} + \boldsymbol{\epsilon}_{:,j}, \quad \mathbf{w}_{:,j} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}\right), \quad \boldsymbol{\epsilon}_{i,:} \sim \mathcal{N}\left(\mathbf{0}, \sigma^2 \mathbf{I}\right)$$

$$\mathbf{X}\mathbf{w}_{:,j} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{X}\mathbf{X}^\top\right),$$

$$\mathbf{X}\mathbf{w}_{:,j} + \boldsymbol{\epsilon}_{:,j} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{X}\mathbf{X}^\top + \sigma^2 \mathbf{I}\right)$$

# Linear Latent Variable Model IV

**Dual Probabilistic PCA Max. Likelihood Soln** (Lawrence, 2004, 2005)



$$p\left(\mathbf{Y}|\mathbf{X}\right) = \prod_{j=1}^{p} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \mathbf{X}\mathbf{X}^{\top} + \sigma^2\mathbf{I}\right)$$

# Linear Latent Variable Model IV

**Dual PPCA Max. Likelihood Soln** (Lawrence, 2004, 2005)

$$p\left(\mathbf{Y}|\mathbf{X}\right) = \prod_{j=1}^{p} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \mathbf{K}\right), \quad \mathbf{K} = \mathbf{X}\mathbf{X}^{\top} + \sigma^2\mathbf{I}$$

# Linear Latent Variable Model IV

**PPCA Max. Likelihood Soln** (Tipping and Bishop, 1999)

$$p\left(\mathbf{Y}|\mathbf{X}\right) = \prod_{j=1}^{p} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \mathbf{K}\right), \quad \mathbf{K} = \mathbf{X}\mathbf{X}^{\top} + \sigma^2 \mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{X}\right) = -\frac{p}{2} \log |\mathbf{K}| - \frac{1}{2}\mathrm{tr}\left(\mathbf{K}^{-1}\mathbf{Y}\mathbf{Y}^{\top}\right) + \mathrm{const.}$$

**PPCA Max. Likelihood Soln**

$$p\left(\mathbf{Y}|\mathbf{X}\right) = \prod_{j=1}^{p} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \mathbf{K}\right), \quad \mathbf{K} = \mathbf{X}\mathbf{X}^{\top} + \sigma^2 \mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{X}\right) = -\frac{p}{2}\log|\mathbf{K}| - \frac{1}{2}\mathrm{tr}\left(\mathbf{K}^{-1}\mathbf{Y}\mathbf{Y}^{\top}\right) + \mathrm{const.}$$

If $\mathbf{U}'_q$ are first $q$ principal eigenvectors of $p^{-1}\mathbf{Y}\mathbf{Y}^{\top}$ and the corresponding eigenvalues are $\mathbf{\Lambda}_q$,

# Linear Latent Variable Model IV

**PPCA Max. Likelihood Soln**

$$p\left(\mathbf{Y}|\mathbf{X}\right) = \prod_{j=1}^{p} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0},\mathbf{K}\right), \quad \mathbf{K} = \mathbf{X}\mathbf{X}^{\top} + \sigma^2\mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{X}\right) = -\frac{p}{2}\log|\mathbf{K}| - \frac{1}{2}\text{tr}\left(\mathbf{K}^{-1}\mathbf{Y}\mathbf{Y}^{\top}\right) + \text{const.}$$

If $\mathbf{U}'_q$ are first $q$ principal eigenvectors of $p^{-1}\mathbf{Y}\mathbf{Y}^{\top}$ and the corresponding eigenvalues are $\mathbf{\Lambda}_q$,

$$\mathbf{X} = \mathbf{U}'_q\mathbf{L}\mathbf{R}^{\top}, \quad \mathbf{L} = \left(\mathbf{\Lambda}_q - \sigma^2\mathbf{I}\right)^{\frac{1}{2}}$$

where $\mathbf{R}$ is an arbitrary rotation matrix.

**Dual PPCA Max. Likelihood Soln** (Lawrence, 2004, 2005)

$$p\left(\mathbf{Y}|\mathbf{X}\right) = \prod_{j=1}^{p} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \mathbf{K}\right), \quad \mathbf{K} = \mathbf{X}\mathbf{X}^{\top} + \sigma^2\mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{X}\right) = -\frac{p}{2}\log|\mathbf{K}| - \frac{1}{2}\mathrm{tr}\left(\mathbf{K}^{-1}\mathbf{Y}\mathbf{Y}^{\top}\right) + \mathrm{const.}$$

If $\mathbf{U}_q'$ are first $q$ principal eigenvectors of $p^{-1}\mathbf{Y}\mathbf{Y}^{\top}$ and the corresponding eigenvalues are $\mathbf{\Lambda}_q$,

$$\mathbf{X} = \mathbf{U}_q'\mathbf{L}\mathbf{R}^{\top}, \quad \mathbf{L} = \left(\mathbf{\Lambda}_q - \sigma^2\mathbf{I}\right)^{\frac{1}{2}}$$

where $\mathbf{R}$ is an arbitrary rotation matrix.

## Linear Latent Variable Model IV

**PPCA Max. Likelihood Soln** (Tipping and Bishop, 1999)

$$p\left(\mathbf{Y}|\mathbf{W}\right) = \prod_{i=1}^{n} \mathcal{N}\left(\mathbf{y}_{i,:}|\mathbf{0}, \mathbf{C}\right), \quad \mathbf{C} = \mathbf{W}\mathbf{W}^{\top} + \sigma^2\mathbf{I}$$

$$\log p\left(\mathbf{Y}|\mathbf{W}\right) = -\frac{n}{2}\log|\mathbf{C}| - \frac{1}{2}\text{tr}\left(\mathbf{C}^{-1}\mathbf{Y}^{\top}\mathbf{Y}\right) + \text{const.}$$

If $\mathbf{U}_q$ are first $q$ principal eigenvectors of $n^{-1}\mathbf{Y}^{\top}\mathbf{Y}$ and the corresponding eigenvalues are $\mathbf{\Lambda}_q$,

$$\mathbf{W} = \mathbf{U}_q\mathbf{L}\mathbf{R}^{\top}, \quad \mathbf{L} = \left(\mathbf{\Lambda}_q - \sigma^2\mathbf{I}\right)^{\frac{1}{2}}$$

where $\mathbf{R}$ is an arbitrary rotation matrix.

# Equivalence of Formulations

**The Eigenvalue Problems are equivalent**

- Solution for Probabilistic PCA (solves for the mapping)

$$\mathbf{Y}^\top \mathbf{Y} \mathbf{U}_q = \mathbf{U}_q \mathbf{\Lambda}_q \qquad \mathbf{W} = \mathbf{U}_q \mathbf{L} \mathbf{R}^\top$$
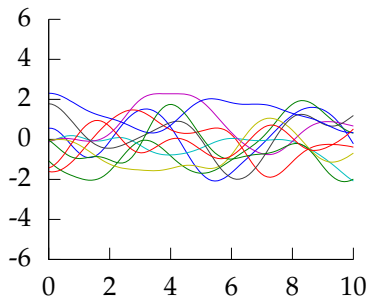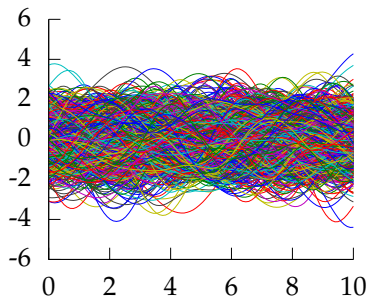
- Solution for Dual Probabilistic PCA (solves for the latent positions)

$$\mathbf{Y} \mathbf{Y}^\top \mathbf{U}_q' = \mathbf{U}_q' \mathbf{\Lambda}_q \qquad \mathbf{X} = \mathbf{U}_q' \mathbf{L} \mathbf{R}^\top$$

- Equivalence is from

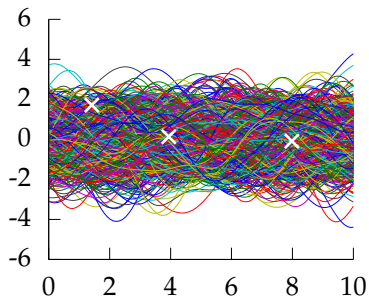$$\mathbf{U}_q = \mathbf{Y}^\top \mathbf{U}_q' \mathbf{\Lambda}_q^{-\frac{1}{2}}$$
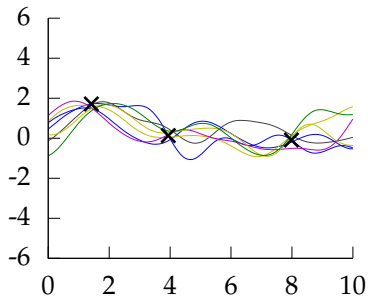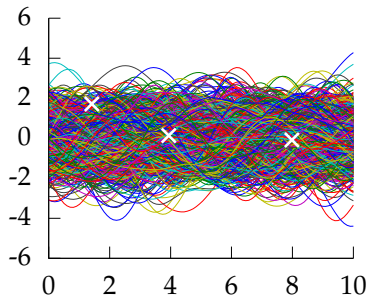
# Gaussian Processes: Extremely Short Overview

# Gaussian Processes: Extremely Short Overview

# Gaussian Processes: Extremely Short Overview

- ▶ http://gpss.cc
- ▶ Next one is in Sheffield in September 2017.
- ▶ Talks and tutorials on line.
- ▶ Jupyter based lab classes.
- ▶ GPy and GPyOpt software available from github.

# Non-Linear Matrix Factorization

- The marginal likelihood of DPPCA is that of a Bayesian linear regression

$$p\left(\mathbf{Y}|\mathbf{X}, \sigma^2, \alpha_x\right) = \prod_{j=1}^{D} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \alpha_w^{-1}\mathbf{X}\mathbf{X}^\top + \sigma^2\mathbf{I}\right).$$

# Non-Linear Matrix Factorization

- The marginal likelihood of DPPCA is that of a Bayesian linear regression

$$p\left(\mathbf{Y}|\mathbf{X}, \sigma^2, \alpha_x\right) = \prod_{j=1}^{D} \mathcal{N}\left(\mathbf{y}_{:,j}|\mathbf{0}, \alpha_w^{-1}\mathbf{K} + \sigma^2\mathbf{I}\right).$$

- Replace inner product matrix with covariance function for non-linear model.

# Missing values

- For the product of GPs marginalizing missing values is straightforward.
- Let $\mathbf{y_i}$ be the observed subset of $\mathbf{y}$.

$$\mathbf{y_i} \sim \mathcal{N}\left(\mu_\mathbf{i}, \Sigma_\mathbf{i,i}\right),$$

- For sparse data

$$p\left(\mathbf{Y}|\mathbf{X}, \sigma^2, \alpha_x\right) = \prod_{j=1}^{D} \mathcal{N}\left(\mathbf{y}_{\mathbf{i}_j, j}|\mathbf{0}, \mathbf{K}_{\mathbf{i}_j, \mathbf{i}_j}\right).$$

# Example: Latent Doodle Space

(Baxter and Anjyo, 2006)

# Example: Latent Doodle Space

<div align="right">(Baxter and Anjyo, 2006)</div>
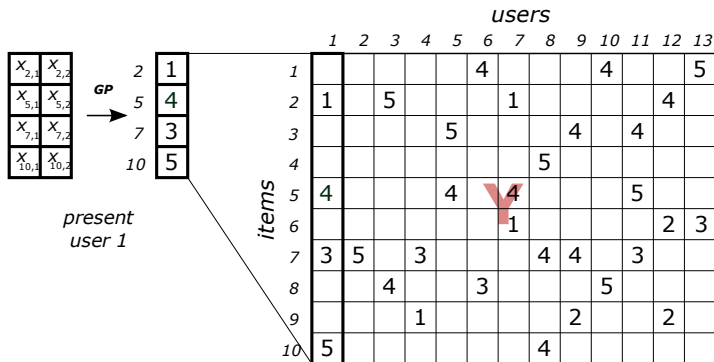
**Generalization with much less Data than Dimensions**

- Powerful uncertainly handling of GPs leads to surprising properties.
- Non-linear models can be used where there are fewer data points than dimensions *without overfitting*.

# Stochastic Gradient Descent



Present data a column at a time.

# Stochastic Gradient Descent



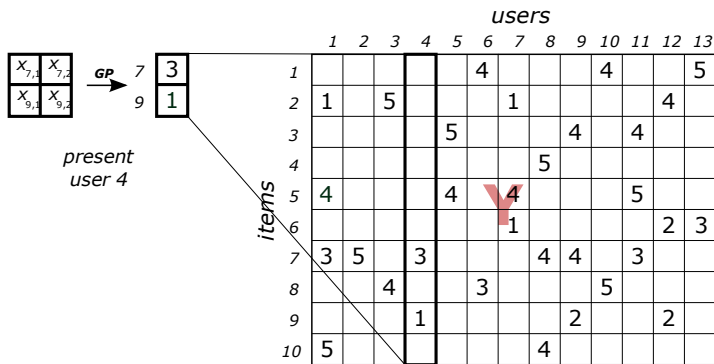Each step updates $\mathbf{X}_{\mathbf{i}_{j},:}$.

Complexity of GP cubic in $N_j$ not $N$.

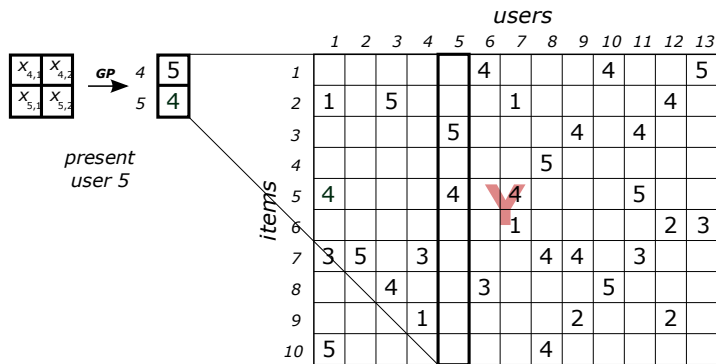# Stochastic Gradient Descent



No Sparse GP approximations required.

# Stochastic Gradient Descent



No Sparse GP approximations required.

# Stochastic Gradient Descent



No Sparse GP approximations required.

# Stochastic Gradient Descent



No Sparse GP approximations required.

# Probabilistic Matrix Factorization for Automated Machine Learning

**Nicoló Fusi**
Microsoft Research
Cambridge, MA, USA
fusi@microsoft.com

**Huseyn Melih Elibol**
Microsoft Research
Cambridge, MA, USA
v-huelib@microsoft.com

## Abstract

In order to achieve state-of-the-art performance, modern machine learning techniques require careful data pre-processing and hyperparameter tuning. Moreover,

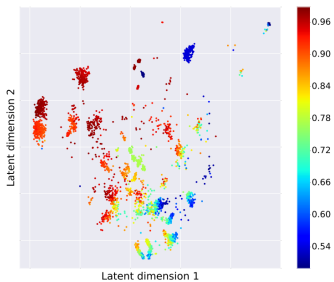Figure 1: Two-dimensional embedding of 5,000 ML pipelines across 576 OpenML datasets. Each point corresponds to a pipeline and is colored by the AUROC obtained by that pipeline in one of the OpenML datasets (OpenML dataset id 943).

# Deep Health

genotype     environment     epigenotype

**G**     **E**     **EG**

$x_1^3$   $x_2^3$   $x_3^3$   $x_4^3$

latent representation
of disease stratification

gene expression

$y_6$

$x_1^2$   $x_2^2$   $x_3^2$   $x_4^2$

$y_1$

survival
analysis

$y_4$   $y_5$   $x_1^1$   $x_2^1$   $x_3^1$   $x_4^1$   $x_5^1$   $y_2$   $y_3$

clinical measurements
and treatment

social
network,
music
data

clinical
notes

**I$_2$**     **I$_1$**

biopsy     X-ray

# Summary

- Many data is usefully summarized with low dimensions.
- Classically pushing probability through non linear functions leads to intractability.
- GP-LVM presents a way around this.
- Recent use case in Automatic Machine Learning

# References I

W. V. Baxter and K.-I. Anjyo. Latent doodle space. In *EUROGRAPHICS*, volume 25, pages 477–485, Vienna, Austria, September 4-8 2006.

C. M. Bishop, M. Svensén, and C. K. I. Williams. GTM: the Generative Topographic Mapping. *Neural Computation*, 10(1):215–234, 1998. [DOI].

C. H. Ek, J. Rihan, P. H. S. Torr, G. Rogez, and N. D. Lawrence. Ambiguity modeling in latent spaces. In A. Popescu-Belis and R. Stiefelhagen, editors, *Machine Learning for Multimodal Interaction (MLMI 2008)*, LNCS, pages 62–73. Springer-Verlag, 28–30 June 2008a. [PDF].

C. H. Ek, P. H. S. Torr, and N. D. Lawrence. Gaussian process latent variable models for human pose estimation. In A. Popescu-Belis, S. Renals, and H. Bourlard, editors, *Machine Learning for Multimodal Interaction (MLMI 2007)*, volume 4892 of *LNCS*, pages 132–143, Brno, Czech Republic, 2008b. Springer-Verlag. [PDF].

K. Grochow, S. L. Martin, A. Hertzmann, and Z. Popovic. Style-based inverse kinematics. In *ACM Transactions on Graphics (SIGGRAPH 2004)*, pages 522–531, 2004.

J. B. Kruskal. Multidimensional scaling by optimizing goodness-of-fit to a nonmetric hypothesis. *Psychometrika*, 29 (1):1–28, 1964. [DOI].

N. D. Lawrence. Gaussian process models for visualisation of high dimensional data. In S. Thrun, L. Saul, and B. Schölkopf, editors, *Advances in Neural Information Processing Systems*, volume 16, pages 329–336, Cambridge, MA, 2004. MIT Press.

N. D. Lawrence. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *Journal of Machine Learning Research*, 6:1783–1816, 11 2005.

D. Lowe and M. E. Tipping. Neuroscale: Novel topographic feature extraction with radial basis function networks. In M. C. Mozer, M. I. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems*, volume 9, pages 543–549, Cambridge, MA, 1997. MIT Press.

D. J. C. MacKay. Bayesian neural networks and density networks. *Nuclear Instruments and Methods in Physics Research, A*, 354(1):73–80, 1995. [DOI].

K. V. Mardia, J. T. Kent, and J. M. Bibby. *Multivariate analysis*. Academic Press, London, 1979. [Google Books] .

V. Priacuriu and I. D. Reid. Nonlinear shape manifolds as shape priors in level set segmentation and tracking. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2011a.

V. Priacuriu and I. D. Reid. Shared shape spaces. In *IEEE International Conference on Computer Vision (ICCV)*, 2011b.

S. T. Roweis. EM algorithms for PCA and SPCA. In M. I. Jordan, M. J. Kearns, and S. A. Solla, editors, *Advances in Neural Information Processing Systems*, volume 10, pages 626–632, Cambridge, MA, 1998. MIT Press.

S. T. Roweis and L. K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500): 2323–2326, 2000. [DOI].

J. W. Sammon. A nonlinear mapping for data structure analysis. *IEEE Transactions on Computers*, C-18(5):401–409, 1969. [DOI].

B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10:1299–1319, 1998. [DOI].

J. B. Tenenbaum, V. de Silva, and J. C. Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000. [DOI].

M. E. Tipping and C. M. Bishop. Probabilistic principal component analysis. *Journal of the Royal Statistical Society, B*, 6 (3):611–622, 1999. [PDF]. [DOI].

R. Urtasun, D. J. Fleet, and P. Fua. 3D people tracking with Gaussian process dynamical models. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 238–245, New York, U.S.A., 17–22 Jun. 2006. IEEE Computer Society Press.

R. Urtasun, D. J. Fleet, A. Hertzmann, and P. Fua. Priors for people tracking from small training sets. In *IEEE International Conference on Computer Vision (ICCV)*, pages 403–410, Bejing, China, 17–21 Oct. 2005. IEEE Computer Society Press.