

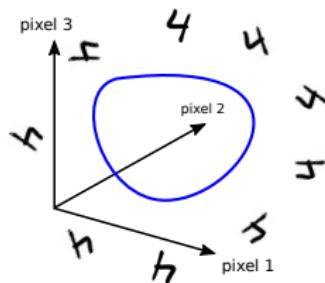
# Non-linear Dimensionality Reduction

Rich Turner

# Applications and reasons to study dimensionality reduction

## modelling data on/near manifolds

e.g. objects + transformations  
= non-linear manifolds



## simple building blocks

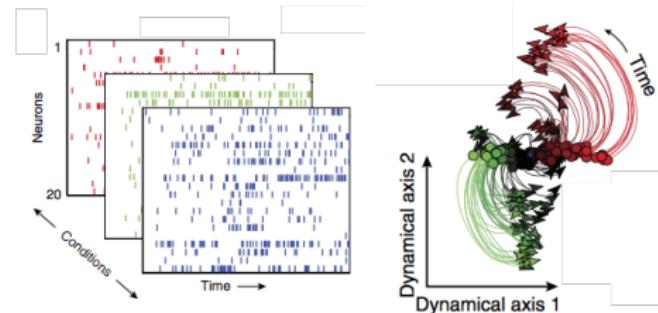
preprocessing/feature learning:  
reducing computational complexity  
improving statistical efficiency

## compose into complex models:

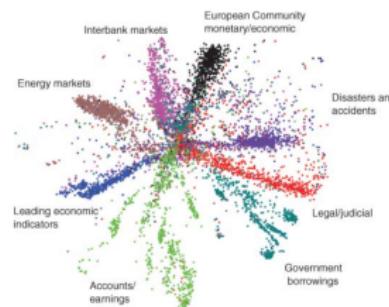
FA -> LGSSMs  
GPLVM -> GPSSMs

## visualisation

understanding structure  
in high-dimensional data

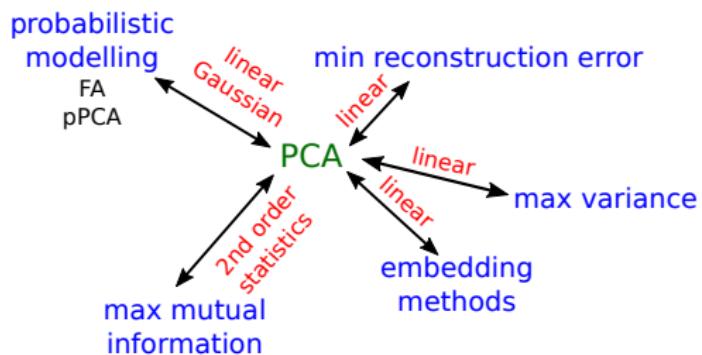


Cunningham and Yu,  
Nature Neuro, 2014

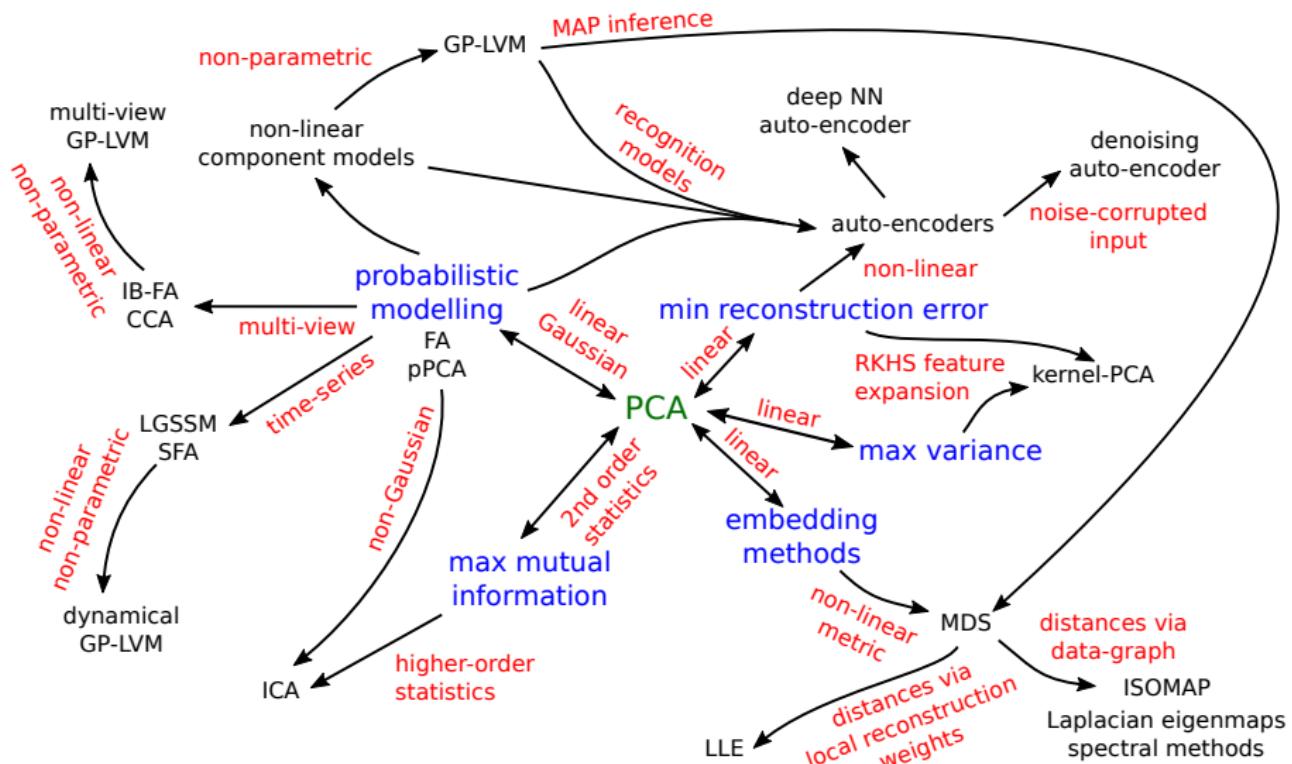


Hinton and Salakhutdinov  
Science, 2006

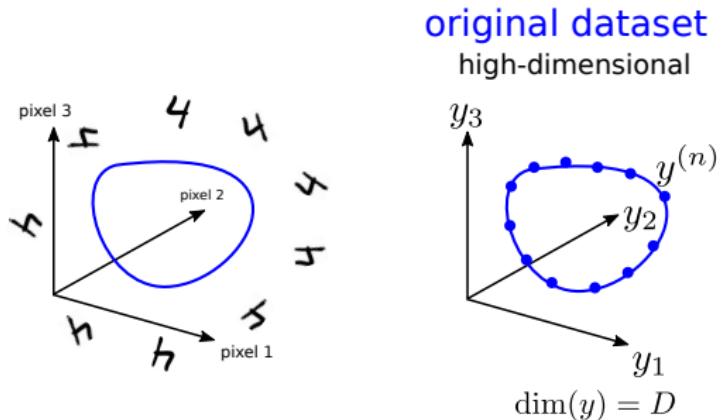
## Dimensionality reduction: conceptual space



# Dimensionality reduction: conceptual space

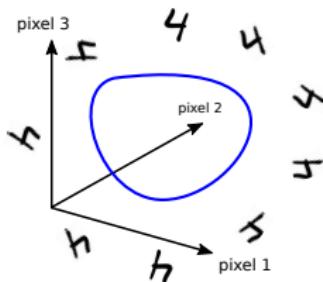


## Dimensionality reduction via distance preserving embeddings



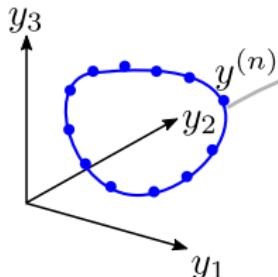
$$d_{nm}^{(y)} = ||y^{(n)} - y^{(m)}||$$

## Dimensionality reduction via distance preserving embeddings



original dataset

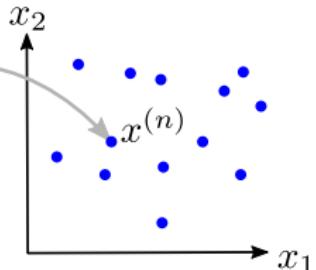
high-dimensional



$$\dim(y) = D$$

new dataset

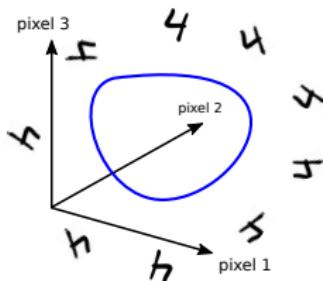
low-dimensional



$$\dim(x) = K$$

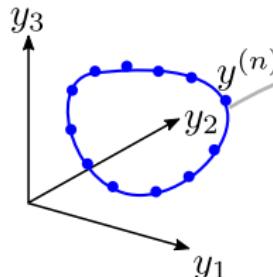
$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

## Dimensionality reduction via distance preserving embeddings



original dataset

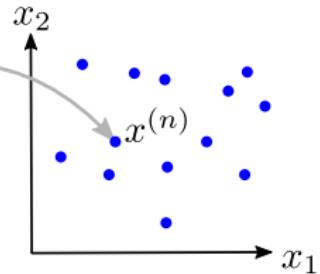
high-dimensional



$$\dim(y) = D$$

new dataset

low-dimensional



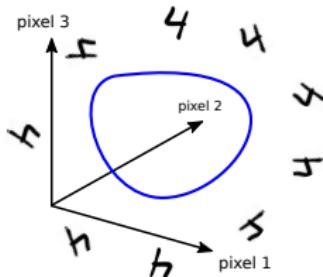
$$\dim(x) = K$$

$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

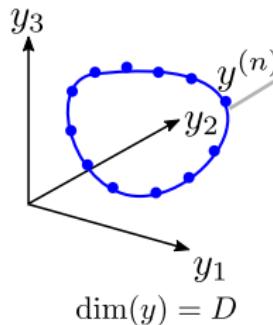
$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

## Dimensionality reduction via distance preserving embeddings



original dataset

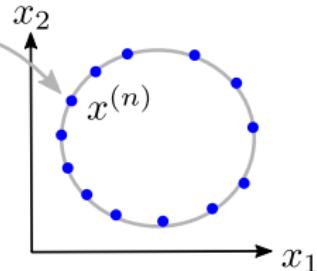
high-dimensional



$$\dim(y) = D$$

new dataset

low-dimensional



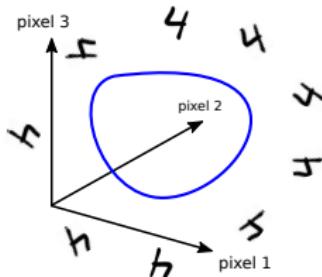
$$\dim(x) = K$$

$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

## Dimensionality reduction via distance preserving embeddings



original dataset

high-dimensional

$y_3$

$y_2$

$$\dim(y) = D$$

new dataset

low-dimensional

$x_2$

$x^{(n)}$

$x_1$

$$\dim(x) = K$$

$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

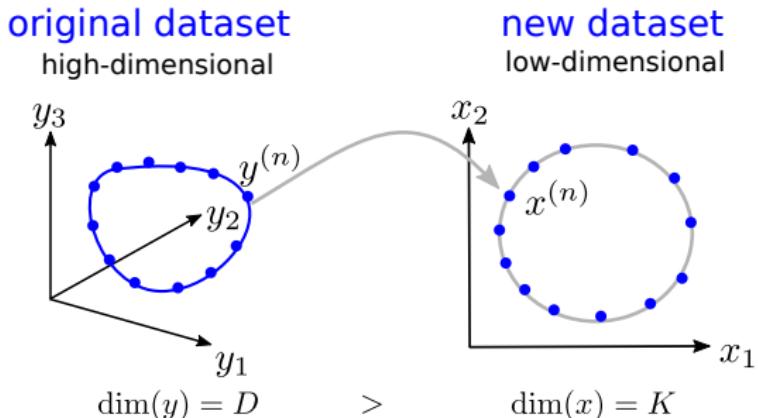
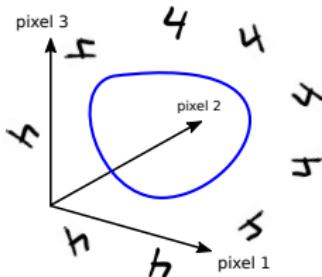
$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

PCA: Euclidean metric (squared error)

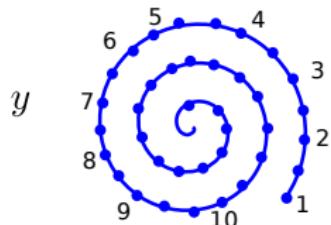
MDS: general distance metric e.g. ISOMAP

Tenenbaum et al  
Science, 2000

# Dimensionality reduction via distance preserving embeddings



**Question:**  
**Will MDS with Euclidean distances make sensible 1D embeddings of:**



$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

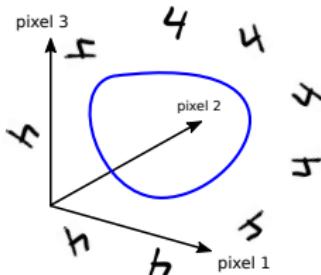
$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

PCA: Euclidean metric (squared error)

MDS: general distance metric e.g. ISOMAP

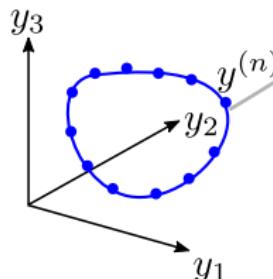
Tenenbaum et al  
Science, 2000

# Dimensionality reduction via distance preserving embeddings



original dataset

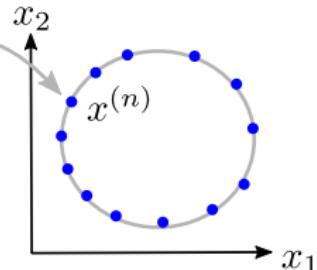
high-dimensional



$$\dim(y) = D$$

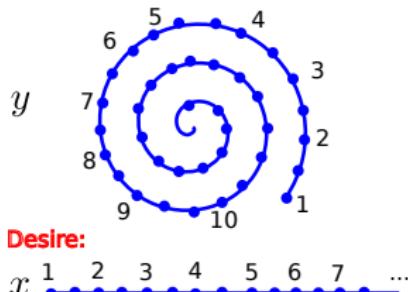
new dataset

low-dimensional



$$\dim(x) = K$$

**Question:**  
Will MDS with Euclidean distances  
make sensible 1D embeddings of:



**Desire:**

$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

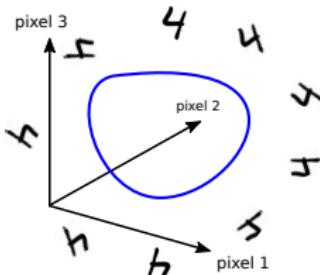
$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

PCA: Euclidean metric (squared error)

MDS: general distance metric e.g. ISOMAP

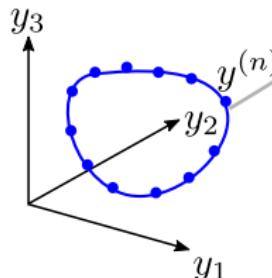
Tenenbaum et al  
Science, 2000

# Dimensionality reduction via distance preserving embeddings



original dataset

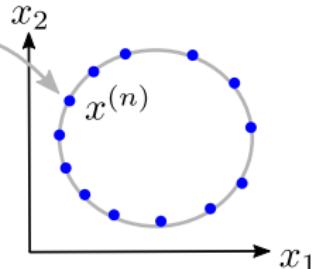
high-dimensional



$$\dim(y) = D$$

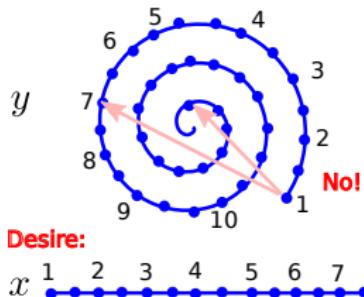
new dataset

low-dimensional



$$\dim(x) = K$$

**Question:**  
Will MDS with Euclidean distances  
make sensible 1D embeddings of:



**Desire:**

$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

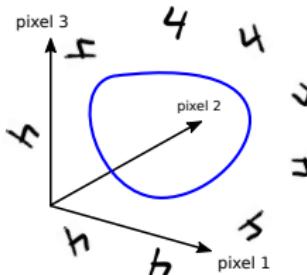
$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

PCA: Euclidean metric (squared error)

MDS: general distance metric e.g. ISOMAP

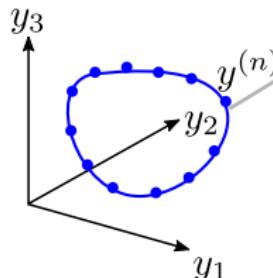
Tenenbaum et al  
Science, 2000

# Dimensionality reduction via distance preserving embeddings

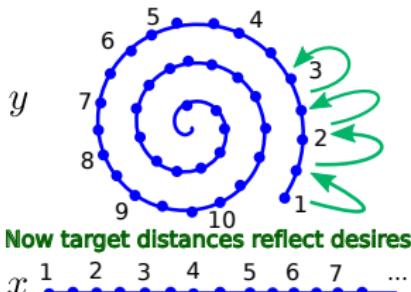


original dataset

high-dimensional

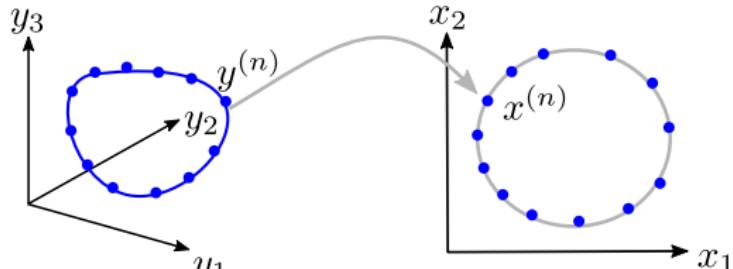


**ISOMAP:**  
Geodesic distance via  
neighbourhood graph



new dataset

low-dimensional



$$\dim(y) = D$$

>

$$\dim(x) = K$$

$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

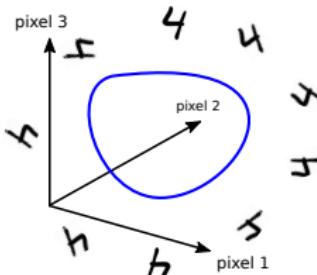
$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

PCA: Euclidean metric (squared error)

MDS: general distance metric e.g. ISOMAP

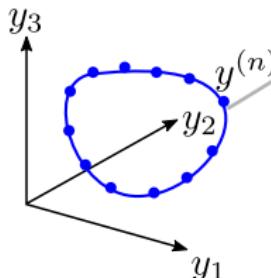
Tenenbaum et al  
Science, 2000

# Dimensionality reduction via distance preserving embeddings

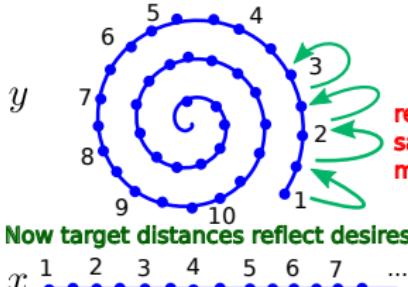


original dataset

high-dimensional



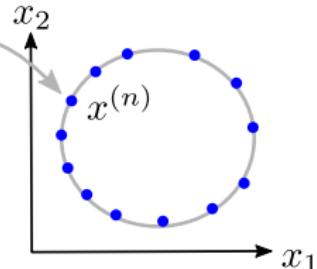
**ISOMAP:**  
Geodesic distance via  
neighbourhood graph



Now target distances reflect desires

new dataset

low-dimensional



$$\dim(y) = D >$$

$$\dim(x) = K$$

$$d_{nm}^{(y)} = \|y^{(n)} - y^{(m)}\| \approx d_{nm}^{(x)} = \|x^{(n)} - x^{(m)}\|$$

embed by optimising new datapoints to match distances:

requires dense sampling of manifold

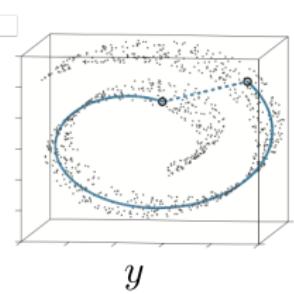
$$\arg \min_{\{x^{(n)}\}_{n=1}^N} \sum_{n < m} \|d_{nm}^{(y)} - d_{nm}^{(x)}\|$$

PCA: Euclidean metric (squared error)

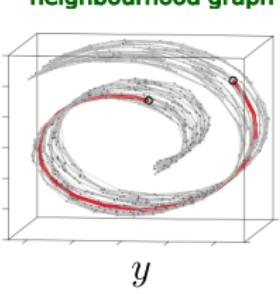
MDS: general distance metric e.g. ISOMAP

# Dimensionality reduction via distance preserving embeddings

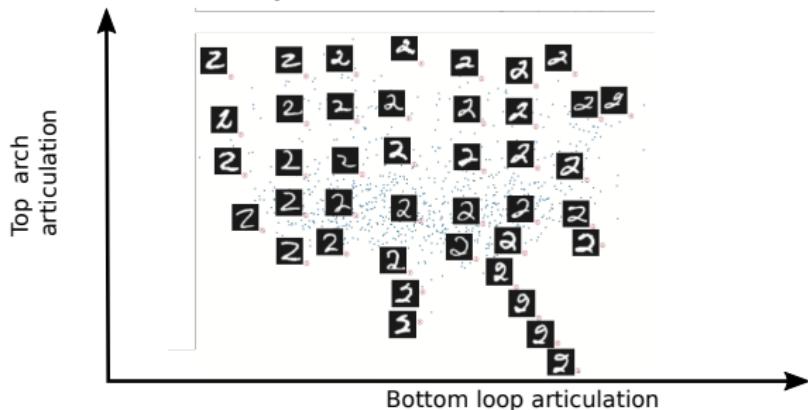
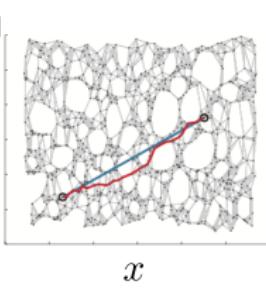
original dataset



geodesic distance via  
neighbourhood graph



new dataset



ISOMAP, Tenenbaum et al Science 2000  
LLE, Roweis et al. Science 2000  
tSNE, Hinton and van der Maaten JMLR 2008

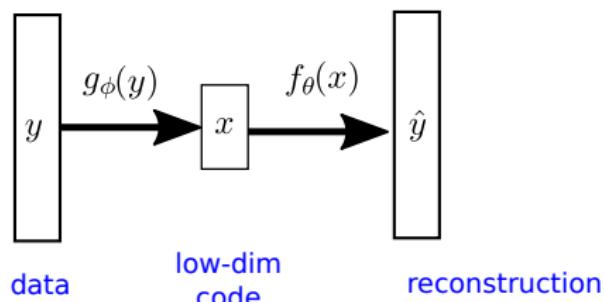
## Dimensionality reduction via distance preserving embeddings

### Limitations (also strengths?)

- ▶ non-linear embedding-based methods require optimisation of new representation  $x$  ( $N \times K$  parameters)
- ▶ works well for low-dimensional embeddings  $K = 2$  or  $3$ , but slow for higher dimensions
- ▶ does not provide quick way to map new data-points into new representation ( $y^{(new)} \rightarrow x^{(new)}$  involves optimisation)
- ▶ does not provide a way for mapping out of reduced space ( $x^{(new)} \rightarrow y^{(new)}$ ?)

**Next: auto-encoders fix above using supervised learning algorithm (non-linear regression) for unsupervised dimensionality reduction**

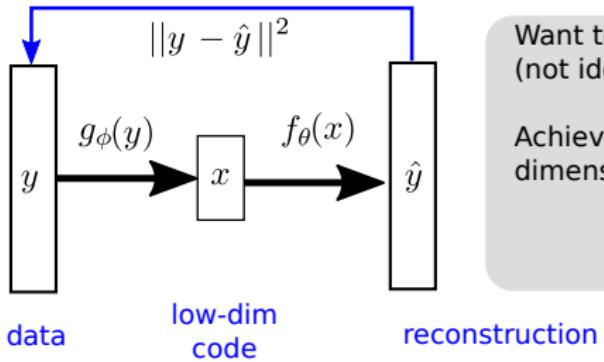
## Dimensionality reduction using auto-encoders



data      low-dim code      reconstruction

$$x(y) = g_\phi(y) \quad \hat{y}(x) = f_\theta(x)$$

## Dimensionality reduction using auto-encoders



Want to learn interesting embedding  
(not identity mapping)

Achieve via constraint:  
dimensionality

data                  low-dim code                  reconstruction

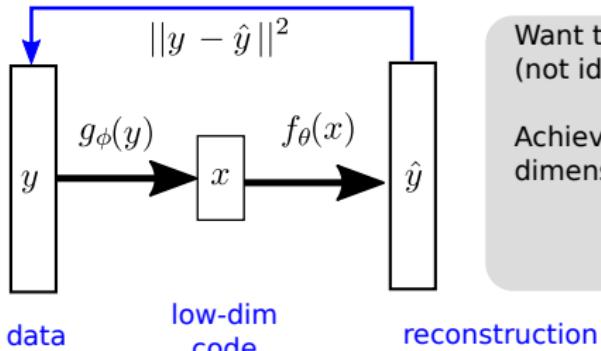
$$x(y) = g_\phi(y)$$

$$\hat{y}(x) = f_\theta(x)$$

cost function

$$\arg \min_{\theta, \phi} \sum_{n=1}^N ||y^{(n)} - \hat{y}^{(n)}||^2$$

## Dimensionality reduction using auto-encoders



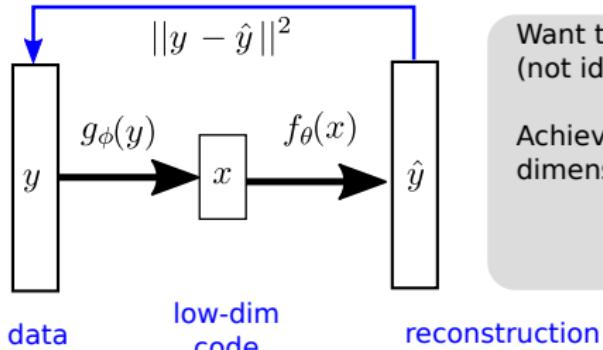
Want to learn interesting embedding  
(not identity mapping)

Achieve via constraint:  
dimensionality, sparsity, function complexity

cost function

$$\arg \min_{\theta, \phi} \sum_{n=1}^N \|y^{(n)} - \hat{y}^{(n)}\|^2 + \text{constraints}$$

## Dimensionality reduction using auto-encoders



Want to learn interesting embedding  
(not identity mapping)

Achieve via constraint:  
dimensionality, sparsity, function complexity

data                  low-dim code                  reconstruction

$$x(y) = g_\phi(y)$$

$$\hat{y}(x) = f_\theta(x)$$

cost function

$$\arg \min_{\theta, \phi} \sum_{n=1}^N \|y^{(n)} - \hat{y}^{(n)}\|^2 + \text{constraints}$$

PCA:

linear

$$x = \Phi y$$

$$\dim(\Phi) = K \times D$$

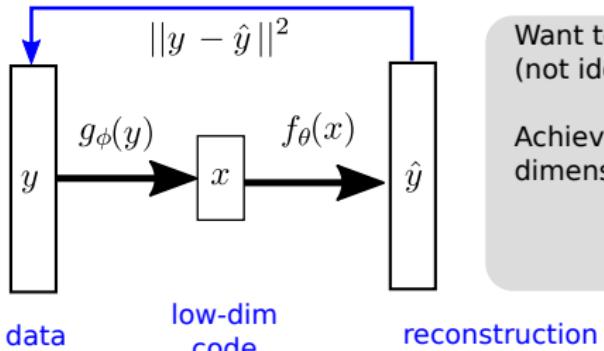
linear

$$\hat{y} = \Theta x$$

$$\dim(\Theta) = D \times K$$

$$\arg \min_{\Theta, \Phi} \sum_{n=1}^N \|y^{(n)} - \Theta \Phi y^{(n)}\|^2$$

# Dimensionality reduction using auto-encoders



Want to learn interesting embedding  
(not identity mapping)

Achieve via constraint:  
dimensionality, sparsity, function complexity

data                  low-dim code                  reconstruction

$$x(y) = g_\phi(y) \quad \hat{y}(x) = f_\theta(x)$$

PCA:

linear	linear
$x = \Phi y$	$\hat{y} = \Theta x$
$\dim(\Phi) = K \times D$	$\dim(\Theta) = D \times K$

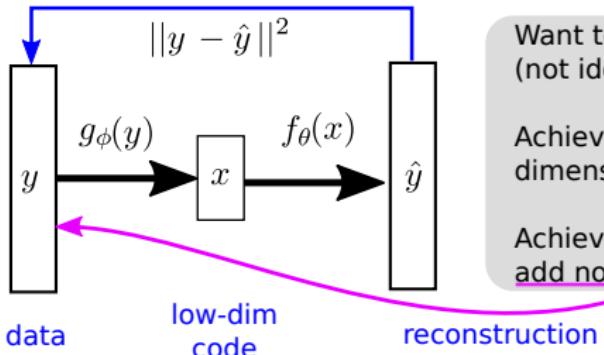
cost function

$$\arg \min_{\theta, \phi} \sum_{n=1}^N \|y^{(n)} - \hat{y}^{(n)}\|^2 + \text{constraints}$$

$$\arg \min_{\Theta, \Phi} \sum_{n=1}^N \|y^{(n)} - \Theta \Phi y^{(n)}\|^2$$

deep neural AA:  
deep neural network

# Dimensionality reduction using auto-encoders



Want to learn interesting embedding  
(not identity mapping)

Achieve via constraint:  
dimensionality, sparsity, function complexity

Achieve via data corruption:  
add noise, drop-out, transform

Vincent, Bengio et al  
JMLR, 2010

$$x(y) = g_\phi(y)$$

$$\hat{y}(x) = f_\theta(x)$$

cost function

$$\arg \min_{\theta, \phi} \sum_{n=1}^N \|y^{(n)} - \hat{y}^{(n)}\|^2 + \text{constraints}$$

PCA:

linear

$$x = \Phi y$$

$$\dim(\Phi) = K \times D$$

linear

$$\hat{y} = \Theta x$$

$$\dim(\Theta) = D \times K$$

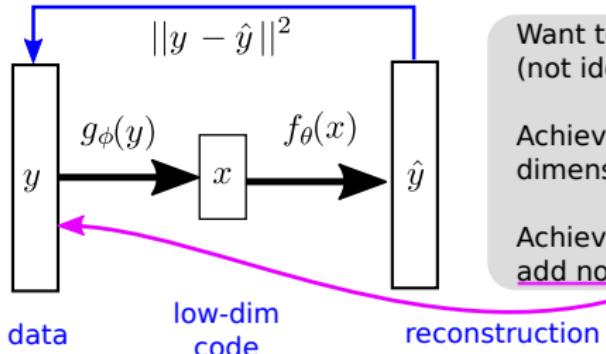
$$\arg \min_{\Theta, \Phi} \sum_{n=1}^N \|y^{(n)} - \Theta \Phi y^{(n)}\|^2$$

deep  
neural AA:

deep neural  
network

deep neural  
network

# Dimensionality reduction using auto-encoders



Want to learn interesting embedding  
(not identity mapping)

Achieve via constraint:  
dimensionality, sparsity, function complexity

Achieve via data corruption:  
add noise, drop-out, transform

Vincent, Bengio et al  
JMLR, 2010

$$\text{data} \quad \xrightarrow{\text{low-dim code}} \quad \text{reconstruction}$$
$$x(y) = g_\phi(y) \qquad \hat{y}(x) = f_\theta(x)$$

$$\text{cost function}$$
$$\arg \min_{\theta, \phi} \sum_{n=1}^N \|y^{(n)} - \hat{y}^{(n)}\|^2 + \text{constraints}$$

PCA:

linear	linear
$x = \Phi y$	$\hat{y} = \Theta x$
$\dim(\Phi) = K \times D$	$\dim(\Theta) = D \times K$

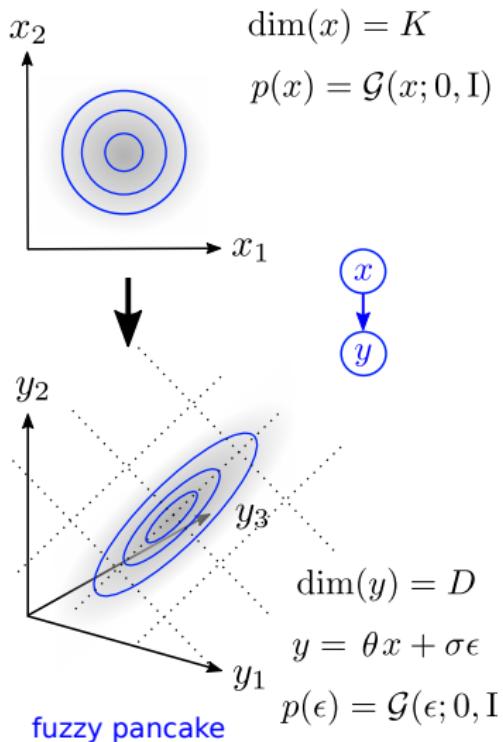
$$\arg \min_{\Theta, \Phi} \sum_{n=1}^N \|y^{(n)} - \Theta \Phi y^{(n)}\|^2$$

deep neural AA:  
Hinton et al Science, 2006

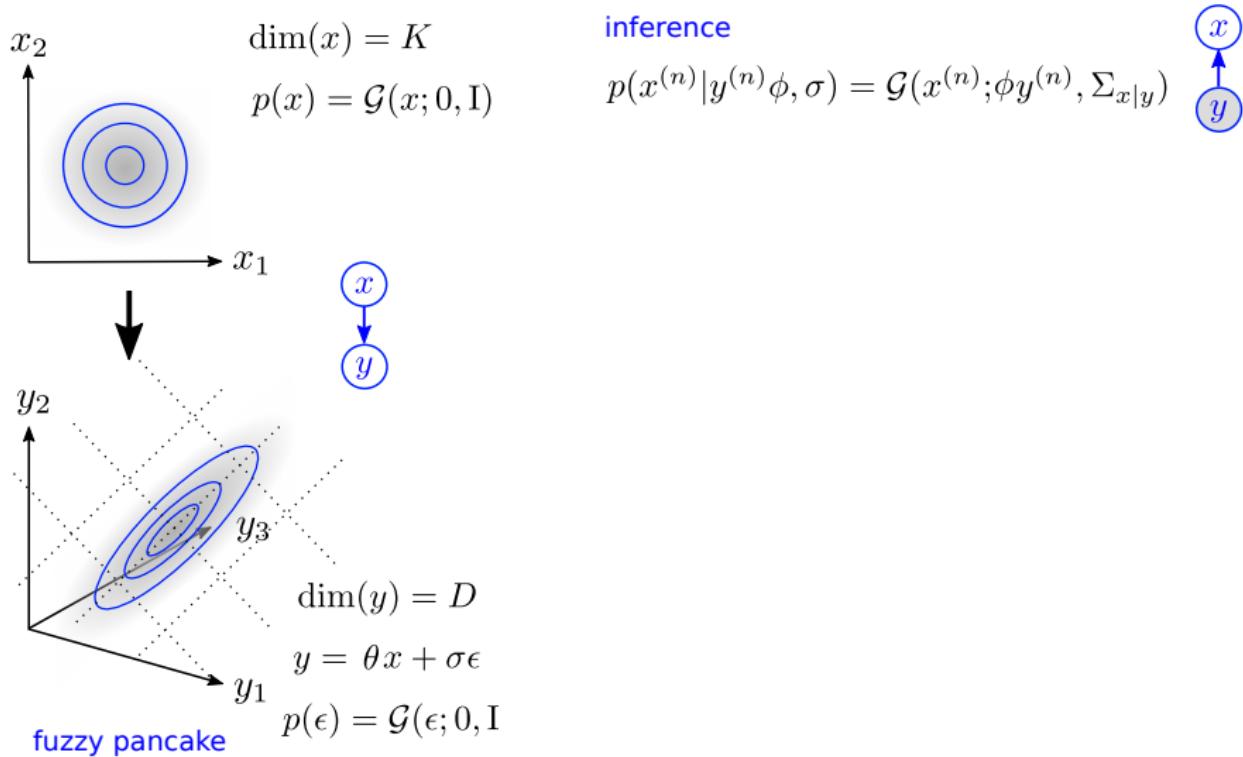
deep neural network

auto-encoders:  
1. learn functions for both mappings  
2. num. params. does not scale with N

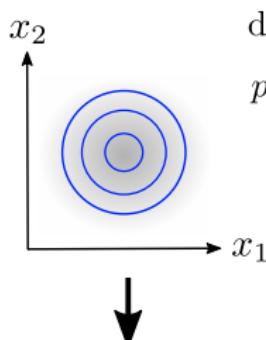
## Dimensionality reduction using probabilistic models: PCA



## Dimensionality reduction using probabilistic models: PCA



## Dimensionality reduction using probabilistic models: PCA



$$\dim(x) = K$$

$$p(x) = \mathcal{G}(x; 0, \mathbf{I})$$

**inference**

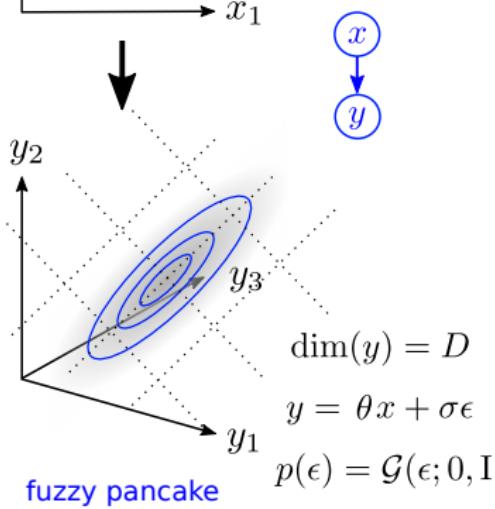
$$p(x^{(n)}|y^{(n)}\phi, \sigma) = \mathcal{G}(x^{(n)};\phi y^{(n)}, \Sigma_{x|y})$$



**maximum-likelihood learning**

$$\theta^{ML}, \sigma^{ML} = \arg \max_{\phi, \sigma} \sum_n \log p(y^{(n)}|\theta, \sigma)$$

$$p(y^{(n)}|\theta, \sigma) = \mathcal{G}(y^{(n)}; 0, \theta \theta^\top + \mathbf{I}\sigma^2)$$

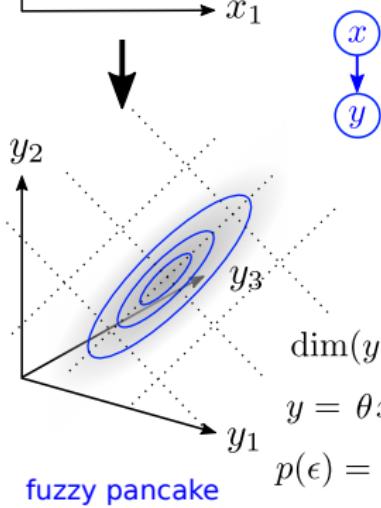
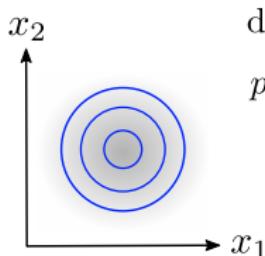


$$\dim(y) = D$$

$$y = \theta x + \sigma \epsilon$$

$$p(\epsilon) = \mathcal{G}(\epsilon; 0, \mathbf{I})$$

# Dimensionality reduction using probabilistic models: PCA



inference

$$p(x^{(n)}|y^{(n)}, \phi, \sigma) = \mathcal{G}(x^{(n)}; \phi y^{(n)}, \Sigma_{x|y})$$



maximum-likelihood learning

$$\theta^{ML}, \sigma^{ML} = \arg \max_{\phi, \sigma} \sum_n \log p(y^{(n)}|\theta, \sigma)$$

$$p(y^{(n)}|\theta, \sigma) = \mathcal{G}(y^{(n)}; 0, \theta \theta^\top + \mathbf{I}\sigma^2)$$

differentiate, set to zero, find:

$$\sigma^{ML} = \frac{1}{D-K} \sum_{k=K+1}^D \lambda_k \quad E_K = [e_1, \dots, e_K] \quad \Lambda_K = \text{diag}(\lambda_1, \dots, \lambda_K)$$

$$\theta^{ML} = E_K (\Lambda_K - \sigma^2 \mathbf{I})^{1/2} R$$

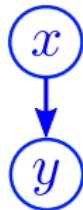
arbitrary rotation matrix

$$\hat{\Sigma}_y e_k = \lambda_k e_k \quad \hat{\Sigma}_y = \frac{1}{N} \sum_n y^{(n)} (y^{(n)})^\top$$

eigenvalues      eigenvectors

Tipping and Bishop 1999, Roweis 1997

## Dimensionality reduction using probabilistic models: a family of models



MODEL  
CLASS

full linear    factor analysis (FA)

$$p(x) = \mathcal{G}(x; 0, I)$$

diagonal

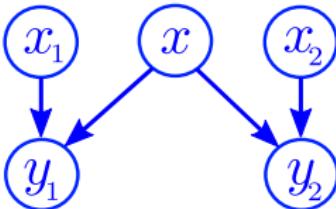
$$p(y|x) = \mathcal{G}(y; \theta x, D)$$

special                PCA  
linear

$$D = \sigma^2 I$$

train: eigenvalue problem

# Dimensionality reduction using probabilistic models: a family of models



audio



video

MODEL  
CLASS

full linear

factor analysis (FA)

$$\begin{aligned} p(x) &= \mathcal{G}(x; 0, I) & p(x), p(x_i) &= \mathcal{G}(x; 0, I) \\ p(y|x) &= \mathcal{G}(y; \theta x, D) & p(y_i|x) &= \mathcal{G}(y_i; \theta_i^{\text{sh}} x + \theta_i^{\text{pri}}, D) \end{aligned}$$

special  
linear

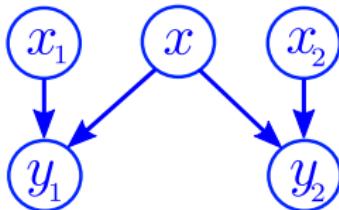
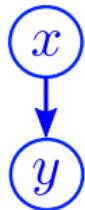
PCA

$$D = \sigma^2 I$$

train: eigenvalue problem

diagonal

# Dimensionality reduction using probabilistic models: a family of models



MODEL  
CLASS

full linear

factor analysis (FA)

$$p(x) = \mathcal{G}(x; 0, I) \quad p(x), p(x_i) = \mathcal{G}(x; 0, I)$$
$$p(y|x) = \mathcal{G}(y; \theta x, D) \quad p(y_i|x) = \mathcal{G}(y_i; \theta_i^{\text{sh}} x + \theta_i^{\text{pri}} x_i, D)$$

special  
linear

PCA

$$D = \sigma^2 I$$

train: eigenvalue problem

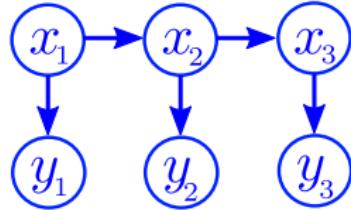
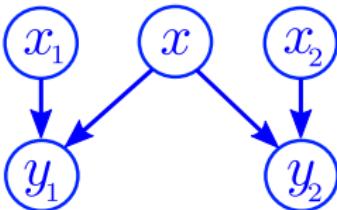
canonical correlation analysis

find linear projections  
of  $y_i$  that max correlation

$$D = \sigma^2 I$$

train: generalised eigenvalue  
problem

# Dimensionality reduction using probabilistic models: a family of models



MODEL  
CLASS

full linear

factor analysis (FA)

$$p(x) = \mathcal{G}(x; 0, I) \quad p(x), p(x_i) = \mathcal{G}(x; 0, I)$$
$$p(y|x) = \mathcal{G}(y; \theta x, D) \quad p(y_i|x) = \mathcal{G}(y_i; \theta_i^{\text{sh}} x + \theta_i^{\text{pri}} x_i, D)$$

special  
linear

PCA

$$D = \sigma^2 I$$

train: eigenvalue problem

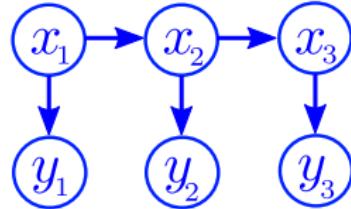
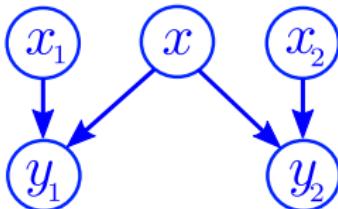
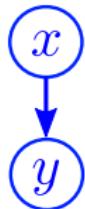
canonical correlation analysis

find linear projections  
of  $y_i$  that max correlation

$$D = \sigma^2 I$$

train: generalised eigenvalue  
problem

# Dimensionality reduction using probabilistic models: a family of models



MODEL  
CLASS

full linear

factor analysis (FA)

$$p(x) = \mathcal{G}(x; 0, I)$$

diagonal

$$p(y|x) = \mathcal{G}(y; \theta x, D)$$



audio



video

special  
linear

PCA

$$D = \sigma^2 I$$

train: eigenvalue problem

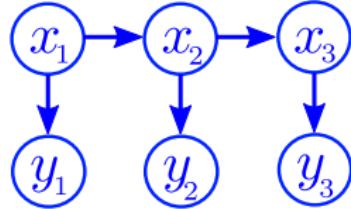
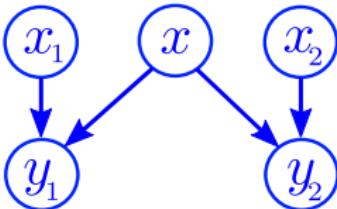
canonical correlation analysis

find linear projections  
of  $y_i$  that max correlation

$$D = \sigma^2 I$$

train: generalised eigenvalue  
problem

# Dimensionality reduction using probabilistic models: a family of models



MODEL  
CLASS

full linear    factor analysis (FA)

$$p(x) = \mathcal{G}(x; 0, I) \quad \text{diagonal}$$

$$p(y|x) = \mathcal{G}(y; \theta x, D)$$



audio



video

LGSSM

$$p(x_t|x_{t-1}) = \mathcal{G}(x_t; \Psi x_{t-1}, \Sigma)$$

$$p(y|x) = \mathcal{G}(y; \theta x, D)$$

special  
linear

PCA

$$D = \sigma^2 I$$

train: eigenvalue problem

find linear projections  
of  $y_i$  that max correlation

canonical correlation analysis

$$D = \sigma^2 I$$

train: generalised eigenvalue  
problem

find slowest  
projections of  
unit variance

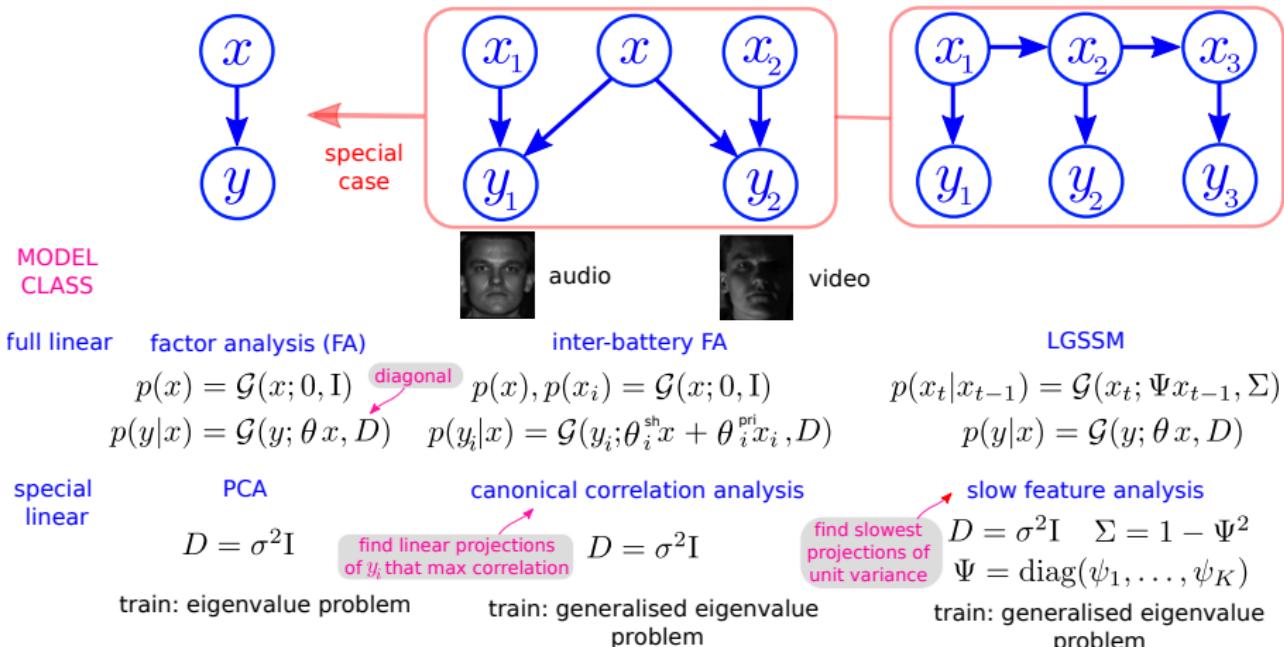
slow feature analysis

$$D = \sigma^2 I \quad \Sigma = 1 - \Psi^2$$

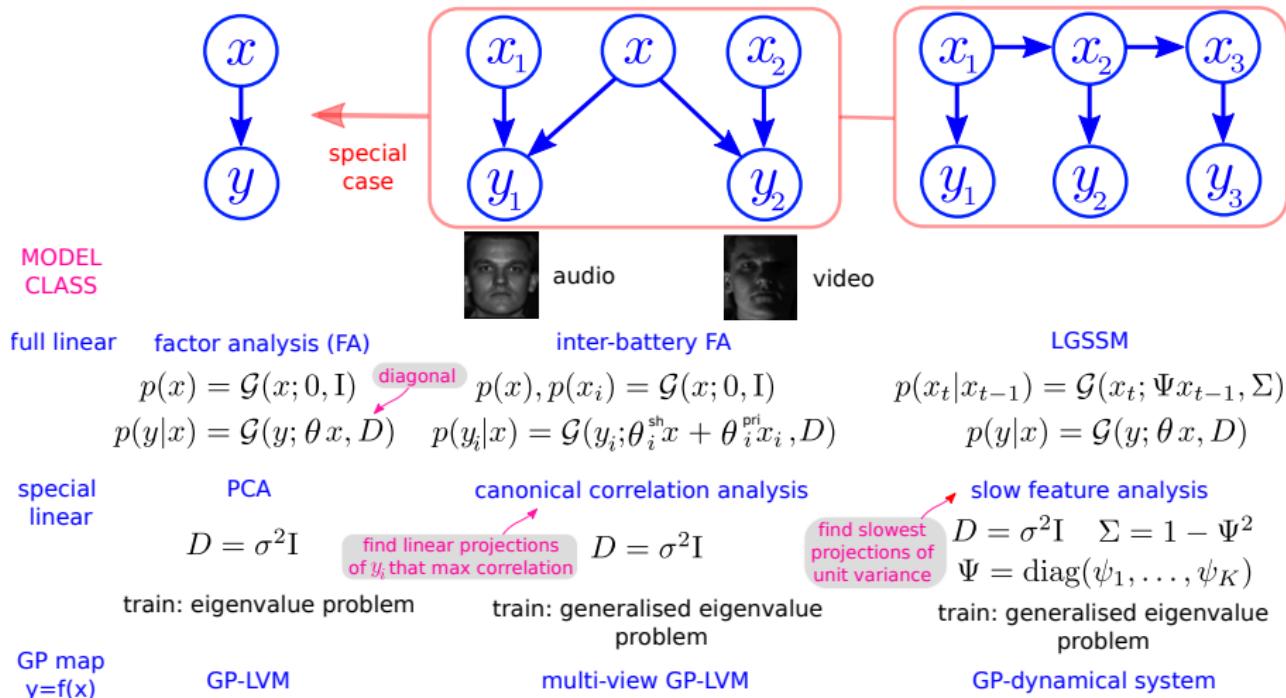
$$\Psi = \text{diag}(\psi_1, \dots, \psi_K)$$

train: generalised eigenvalue  
problem

# Dimensionality reduction using probabilistic models: a family of models



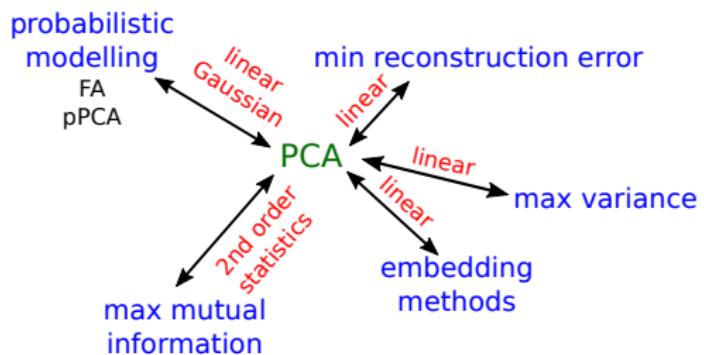
## Dimensionality reduction using probabilistic models: a family of models



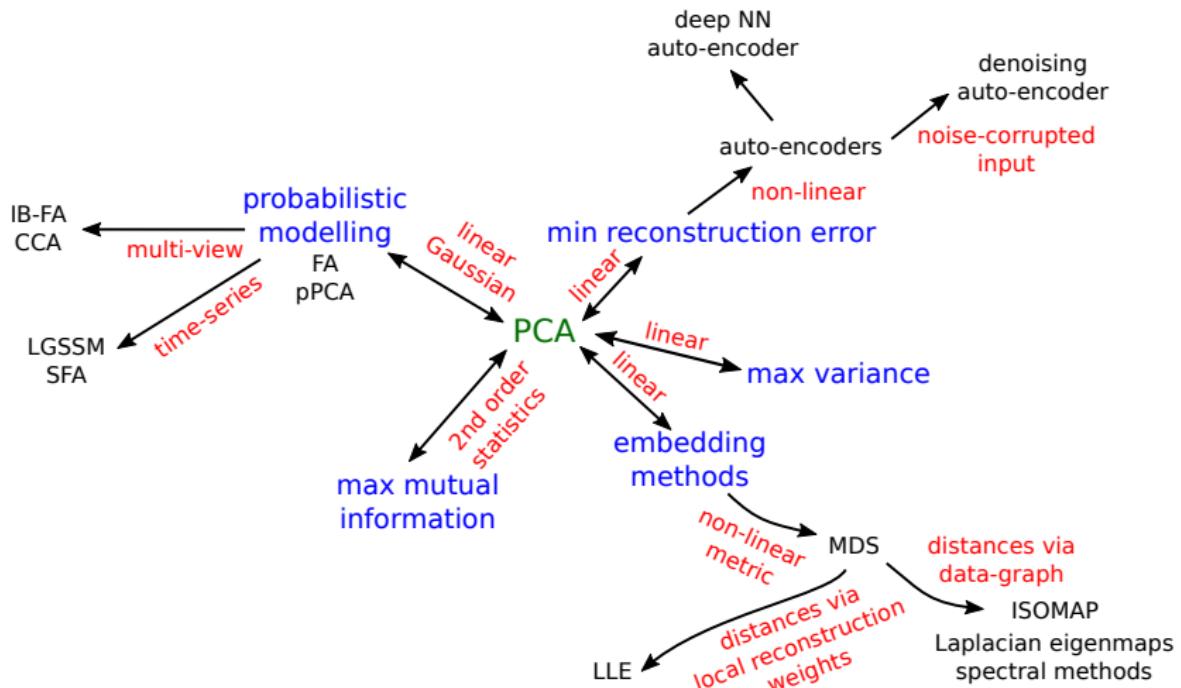
# Dimensionality reduction using probabilistic models: a family of models

MODEL CLASS			
full linear	factor analysis (FA)	inter-battery FA	LGSSM
	$p(x) = \mathcal{G}(x; 0, I)$ $p(y x) = \mathcal{G}(y; \theta x, D)$	$p(x), p(x_i) = \mathcal{G}(x; 0, I)$ $p(y_i x) = \mathcal{G}(y_i; \theta_i^{\text{sh}} x + \theta_i^{\text{pri}} x_i, D)$	$p(x_t x_{t-1}) = \mathcal{G}(x_t; \Psi x_{t-1}, \Sigma)$ $p(y x) = \mathcal{G}(y; \theta x, D)$
special linear	PCA $D = \sigma^2 I$	canonical correlation analysis $D = \sigma^2 I$ find linear projections of $y_i$ that max correlation	slow feature analysis $D = \sigma^2 I$ $\Sigma = 1 - \Psi^2$ $\Psi = \text{diag}(\psi_1, \dots, \psi_K)$ find slowest projections of unit variance
	train: eigenvalue problem	train: generalised eigenvalue problem	train: generalised eigenvalue problem
GP map $y=f(x)$	GP-LVM	multi-view GP-LVM	GP-dynamical system
other	many e.g. ICA	information bottleneck style and content	many e.g. GP-SSM

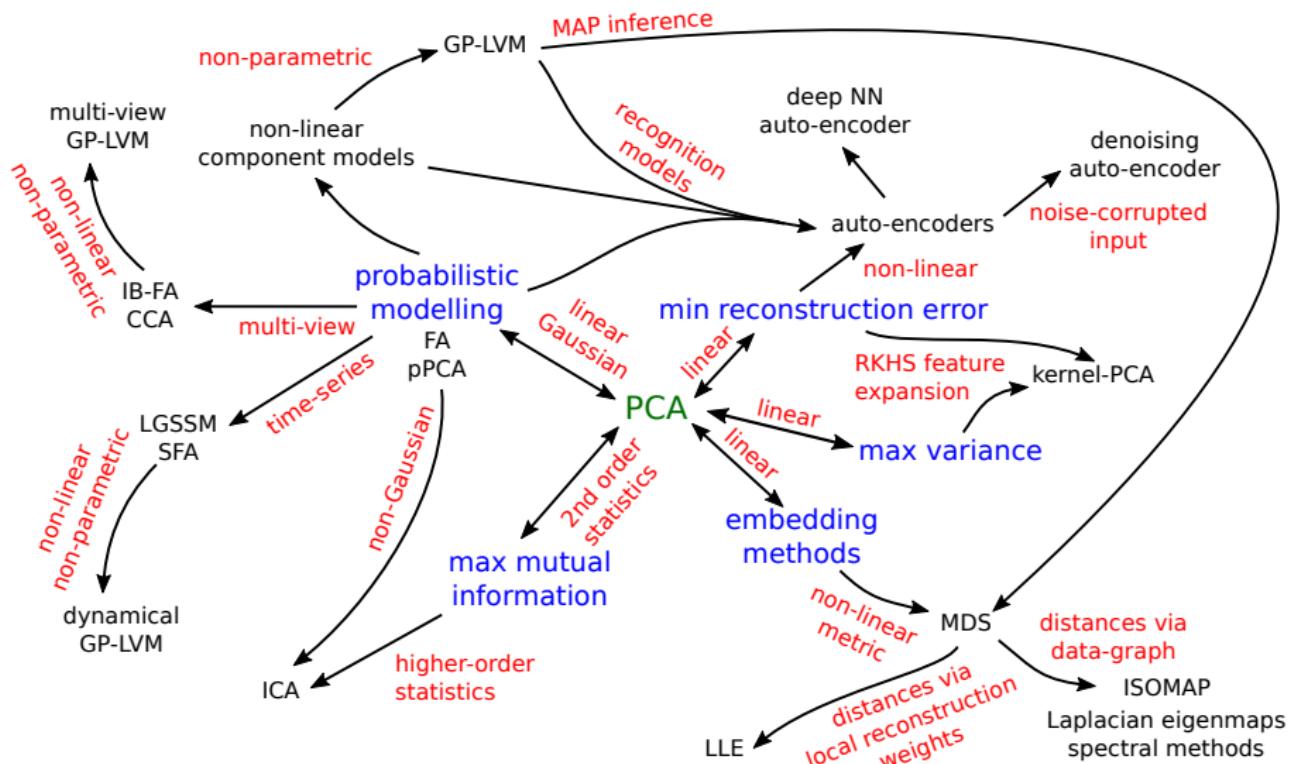
## Roadmap so far...



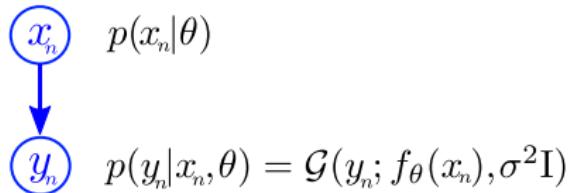
# Roadmap so far...



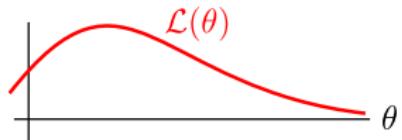
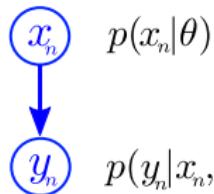
# Roadmap so far...



## Probabilistic inference as an auto-encoder



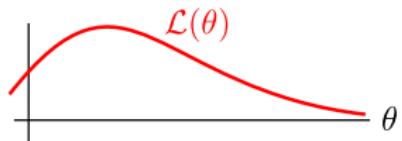
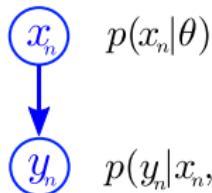
## Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

$$\mathcal{L}(\theta) = \sum_n \log p(y_n|\theta)$$

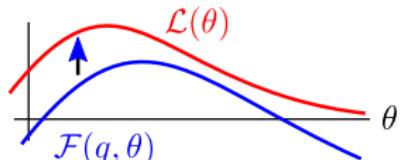
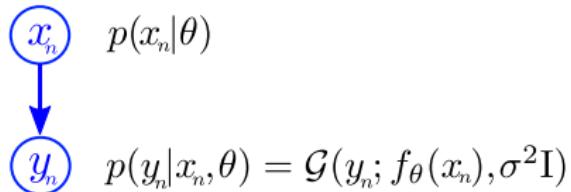
## Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

$$\mathcal{L}(\theta) = \sum_n \log p(y_n|\theta) = \sum_n \log \int p(y_n, x_n|\theta) dx_n = \sum_n \log \int \frac{q(x_n)}{q(x_n)} p(y_n, x_n|\theta) dx_n$$

## Probabilistic inference as an auto-encoder



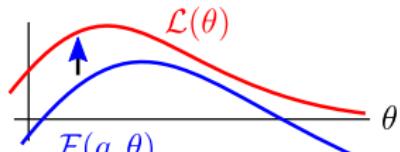
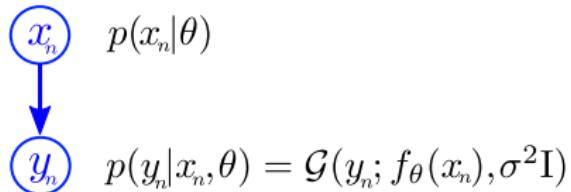
Goal: learn parameters via approximate maximum-likelihood

$$\mathcal{L}(\theta) = \sum_n \log p(y_n|\theta) = \sum_n \log \int p(y_n, x_n|\theta) dx_n = \sum_n \log \int \frac{q(x_n)}{q(x_n)} p(y_n, x_n|\theta) dx_n$$

Jensen's 

$$\mathcal{L}(\theta) \geq \sum_n \int q(x_n) \log \frac{1}{q(x_n)} p(y_n, x_n|\theta) dx_n = \mathcal{F}(q, \theta) \leftarrow \text{free-energy}$$

## Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

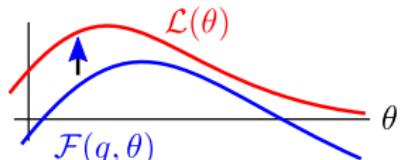
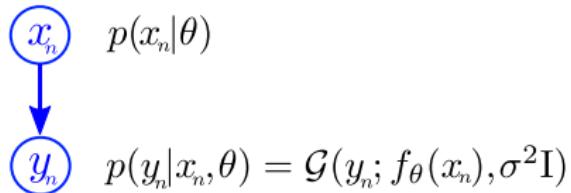
$$\mathcal{L}(\theta) = \sum_n \log p(y_n|\theta) = \sum_n \log \int p(y_n, x_n|\theta) dx_n = \sum_n \log \int \frac{q(x_n)}{q(x_n)} p(y_n, x_n|\theta) dx_n$$

Jensen's ↘

$$\mathcal{L}(\theta) \geq \sum_n \int q(x_n) \log \frac{1}{q(x_n)} p(y_n, x_n|\theta) dx_n = \mathcal{F}(q, \theta) \leftarrow \text{free-energy}$$

$$\mathcal{F}(q, \theta) = \sum_n \int q(x_n) \log p(y_n|x_n, \theta) dx_n + \sum_n \int q(x_n) \log \frac{p(x_n|\theta)}{q(x_n)} dx_n$$

## Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

$$\mathcal{L}(\theta) = \sum_n \log p(y_n|\theta) = \sum_n \log \int p(y_n, x_n|\theta) dx_n = \sum_n \log \int \frac{q(x_n)}{q(x_n)} p(y_n, x_n|\theta) dx_n$$

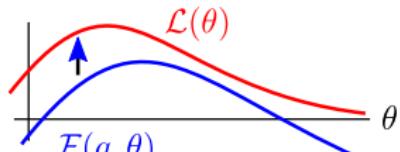
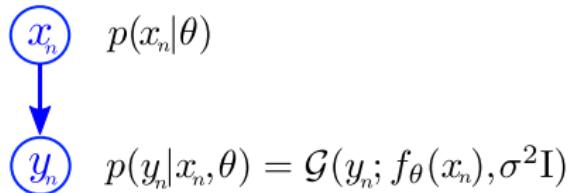
Jensen's ↘

$$\mathcal{L}(\theta) \geq \sum_n \int q(x_n) \log \frac{1}{q(x_n)} p(y_n, x_n|\theta) dx_n = \mathcal{F}(q, \theta) \leftarrow \text{free-energy}$$

$$\mathcal{F}(q, \theta) = \sum_n \int q(x_n) \log p(y_n|x_n, \theta) dx_n + \sum_n \int q(x_n) \log \frac{p(x_n|\theta)}{q(x_n)} dx_n$$

$$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \|y_n - f_\theta(x_n)\|^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) || p(x_n))$$

## Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

$$\mathcal{L}(\theta) = \sum_n \log p(y_n|\theta) = \sum_n \log \int p(y_n, x_n|\theta) dx_n = \sum_n \log \int \frac{q(x_n)}{q(x_n)} p(y_n, x_n|\theta) dx_n$$

Jensen's

$$\mathcal{L}(\theta) \geq \sum_n \int q(x_n) \log \frac{1}{q(x_n)} p(y_n, x_n|\theta) dx_n = \mathcal{F}(q, \theta) \leftarrow \text{free-energy}$$

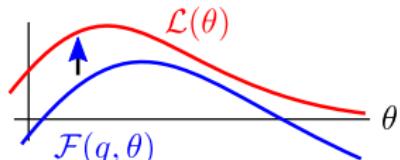
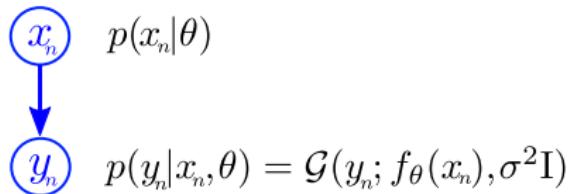
$$\mathcal{F}(q, \theta) = \sum_n \int q(x_n) \log p(y_n|x_n, \theta) dx_n + \sum_n \int q(x_n) \log \frac{p(x_n|\theta)}{q(x_n)} dx_n$$

reconstruction cost

soft-constraint

$$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \|y_n - f_\theta(x_n)\|^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) \| p(x_n))$$

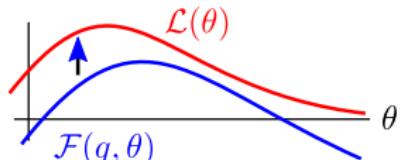
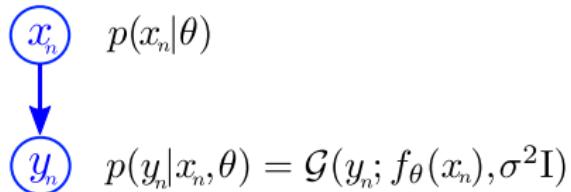
## Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood  
flavours of variational inference:

reconstruction cost	soft-constraint
$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \ y_n - f_\theta(x_n)\ ^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) \  p(x_n))$	

## Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

flavours of variational inference:

fixed family       $q(x_n) = \mathcal{G}(x_n; \mu_{q_n}, \Sigma_{q_n})$

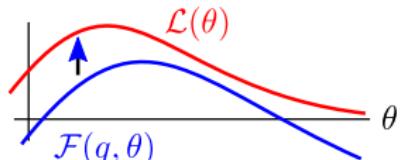
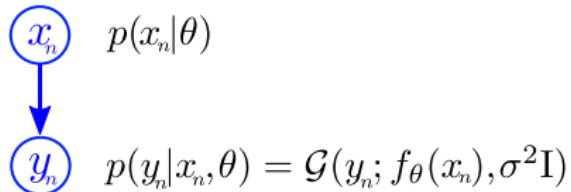
$$\arg \max_{\theta, \mu_{q_{1:N}}, \Sigma_{q_{1:N}}} \mathcal{F}(\mu_{q_{1:N}}, \Sigma_{q_{1:N}}, \theta)$$

reconstruction cost

$$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \|y_n - f_\theta(x_n)\|^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) || p(x_n))$$

soft-constraint

# Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

flavours of variational inference:

fixed family     $q(x_n) = \mathcal{G}(x_n; \mu_{q_n}, \Sigma_{q_n})$

$$\arg \max_{\theta, \mu_{q_{1:N}}, \Sigma_{q_{1:N}}} \mathcal{F}(\mu_{q_{1:N}}, \Sigma_{q_{1:N}}, \theta)$$

structured     $q(x_n) = \prod_{k=1}^K q_k(x_{n_k})$

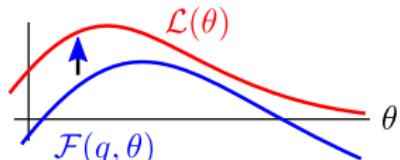
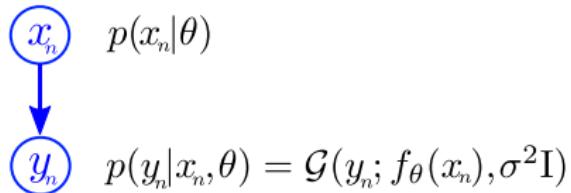
$$\arg \max_{\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta} \mathcal{F}(\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta)$$

reconstruction cost

$$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \|y_n - f_\theta(x_n)\|^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) || p(x_n))$$

soft-constraint

# Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

flavours of variational inference:

fixed family  $q(x_n) = \mathcal{G}(x_n; \mu_{q_n}, \Sigma_{q_n})$

$$\arg \max_{\theta, \mu_{q_{1:N}}, \Sigma_{q_{1:N}}} \mathcal{F}(\mu_{q_{1:N}}, \Sigma_{q_{1:N}}, \theta)$$

structured  $q(x_n) = \prod_{k=1}^K q_k(x_{n_k})$

$$\arg \max_{\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta} \mathcal{F}(\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta)$$

recognition model  $q_\phi(x_n) = \mathcal{G}(x_n; \mu_\phi(y_n), \Sigma_\phi(y_n))$

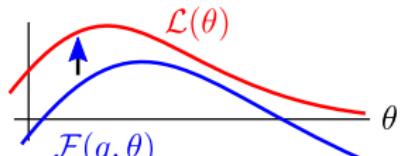
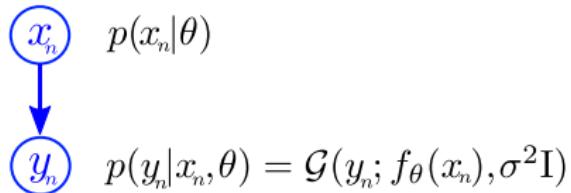
$$\arg \max_{\phi, \theta} \mathcal{F}(\phi, \theta) \quad \text{variational auto-encoder}$$

reconstruction cost

soft-constraint

$$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \|y_n - f_\theta(x_n)\|^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) || p(x_n))$$

# Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

flavours of variational inference:

Roweis and Ghahramani, 1999 (AEs and exact inference)  
Kingma and Welling, 2013 (AEs and approximate inference)

fixed family       $q(x_n) = \mathcal{G}(x_n; \mu_{q_n}, \Sigma_{q_n})$

$\arg \max_{\theta, \mu_{q_{1:N}}, \Sigma_{q_{1:N}}} \mathcal{F}(\mu_{q_{1:N}}, \Sigma_{q_{1:N}}, \theta)$

structured       $q(x_n) = \prod_{k=1}^K q_k(x_{n_k})$

$\arg \max_{\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta} \mathcal{F}(\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta)$

recognition  
model       $q_\phi(x_n) = \mathcal{G}(x_n; \mu_\phi(y_n), \Sigma_\phi(y_n))$

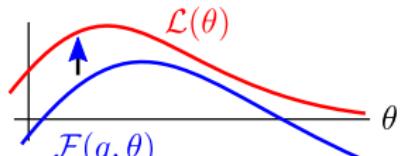
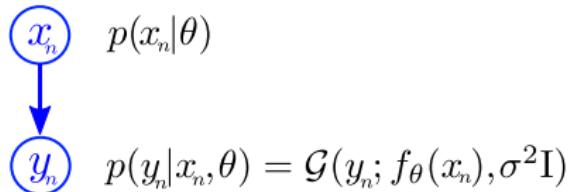
$\arg \max_{\phi, \theta} \mathcal{F}(\phi, \theta)$  **variational  
auto-encoder**

reconstruction cost

soft-constraint

$$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \|y_n - f_\theta(x_n)\|^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) || p(x_n))$$

# Probabilistic inference as an auto-encoder



Goal: learn parameters via approximate maximum-likelihood

flavours of variational inference:

Roweis and Ghahramani, 1999 (AEs and exact inference)  
Kingma and Welling, 2013 (AEs and approximate inference)

fixed family     $q(x_n) = \mathcal{G}(x_n; \mu_{q_n}, \Sigma_{q_n})$

$\arg \max_{\theta, \mu_{q_{1:N}}, \Sigma_{q_{1:N}}} \mathcal{F}(\mu_{q_{1:N}}, \Sigma_{q_{1:N}}, \theta)$

structured     $q(x_n) = \prod_{k=1}^K q_k(x_{n_k})$

$\arg \max_{\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta} \mathcal{F}(\{q_k(x_{n_k})\}_{k=1, n=1}^{K, N}, \theta)$

recognition  
model

$q_\phi(x_n) = \mathcal{G}(x_n; \mu_\phi(y_n), \Sigma_\phi(y_n))$   
GP LVM     $q_\phi(x) = \delta(x - g_\phi(y))$

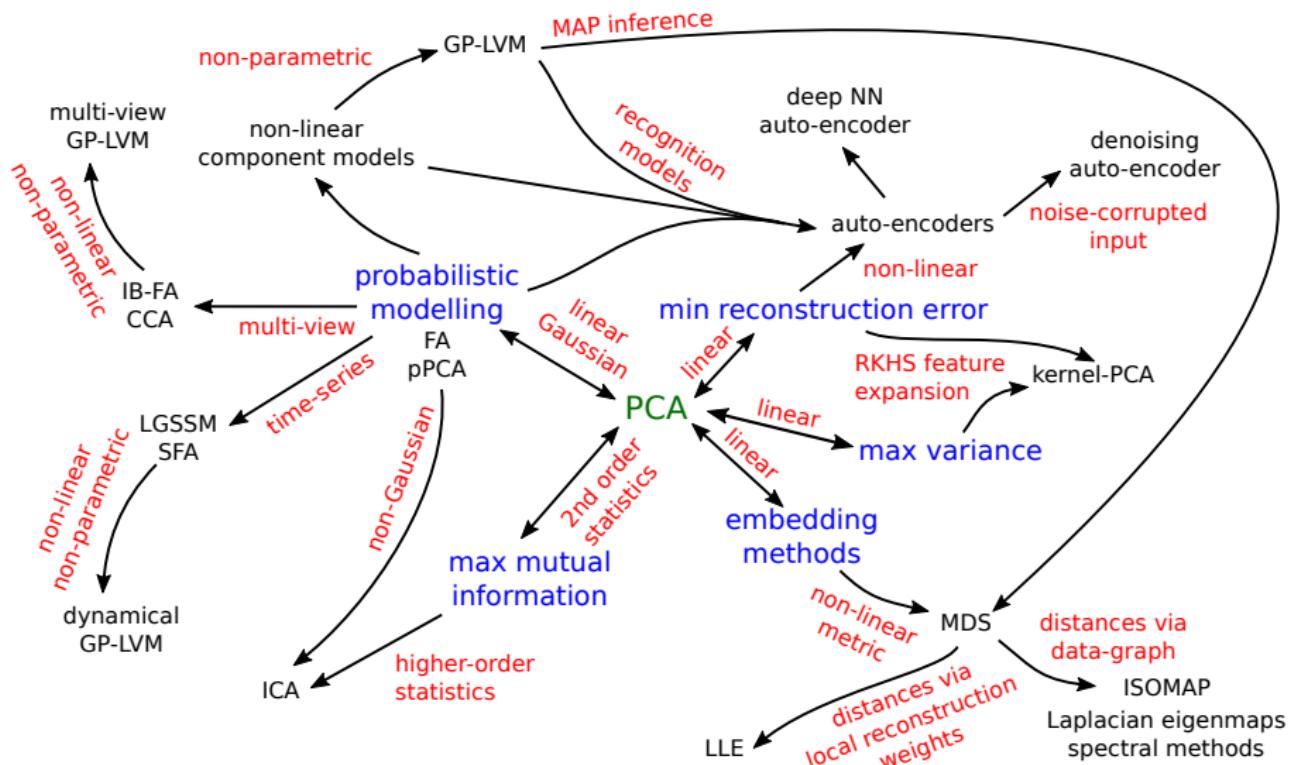
$\arg \max_{\phi, \theta} \mathcal{F}(\phi, \theta)$     **variational  
auto-encoder**

reconstruction cost

soft-constraint

$$\mathcal{F}(q, \theta) = -\sum_n \frac{1}{2\sigma^2} \langle \|y_n - f_\theta(x_n)\|^2 \rangle_{q(x_n)} - \frac{ND}{2} \log \sigma^2 - \sum_n \text{KL}(q(x_n) || p(x_n))$$

# Dimensionality reduction: conceptual space



## Additional Slides

## Gaussian Process Regression Model: Recap

Generative model (like non-linear regression)

$$y(x) = f(x) + \epsilon \sigma_y$$

$$p(\epsilon) = \mathcal{N}(0, 1)$$

place GP prior over the non-linear function

$$p(f(x)|\theta) = \mathcal{GP}(0, K(x, x'))$$

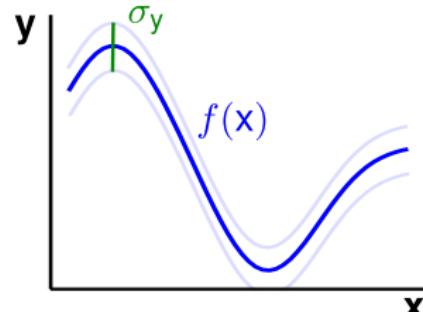
"multivariate Gaussian of infinite dimension"

(any finite subset of variables are multivariate Gaussian)

$$K(x, x') = \sigma^2 \exp\left(-\frac{1}{2l^2}(x - x')^2\right) \quad (\text{smoothly wiggling functions expected})$$

since the sum of two Gaussians is a Gaussian, the model induces a GP over  $y(x)$

$$p(y(x)|\theta) = \mathcal{GP}(0, K(x, x') + I\sigma_y^2)$$



## Gaussian Process Regression Model: Recap

Generative model (like non-linear regression)

$$y(x) = f(x) + \epsilon \sigma_y$$

$$p(\epsilon) = \mathcal{N}(0, 1)$$

place GP prior over the non-linear function

$$p(f(x)|\theta) = \mathcal{GP}(0, K(x, x'))$$

"multivariate Gaussian of infinite dimension"

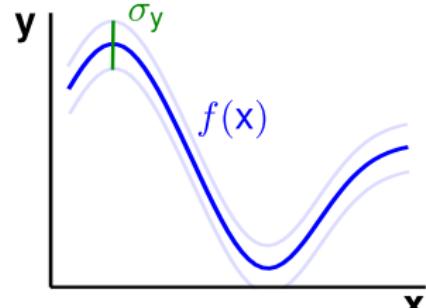
(any finite subset of variables are multivariate Gaussian)

$$K(x, x') = \sigma^2 \exp\left(-\frac{1}{2l^2}(x - x')^2\right) \quad (\text{smoothly wiggling functions expected})$$

since the sum of two Gaussians is a Gaussian, the model induces a GP over  $y(x)$

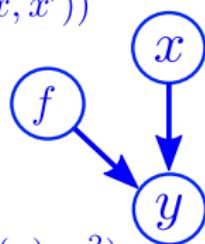
$$p(y(x)|\theta) = \mathcal{GP}(0, K(x, x') + I\sigma_y^2)$$

multi-output regression  $y_i(x) = f_i(x) + \epsilon_i \sigma_y$   $p(f_i(x)|\theta) = \mathcal{GP}(0, K_i(x, x'))$



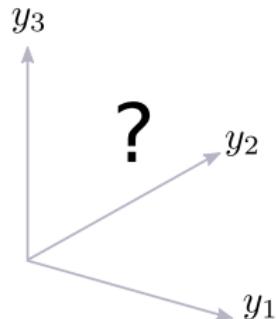
## Gaussian Process Latent Variable Model

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

$$p(x) = \mathcal{G}(x; 0, I)$$



**toy example:**

2 dimensional latents  $x$

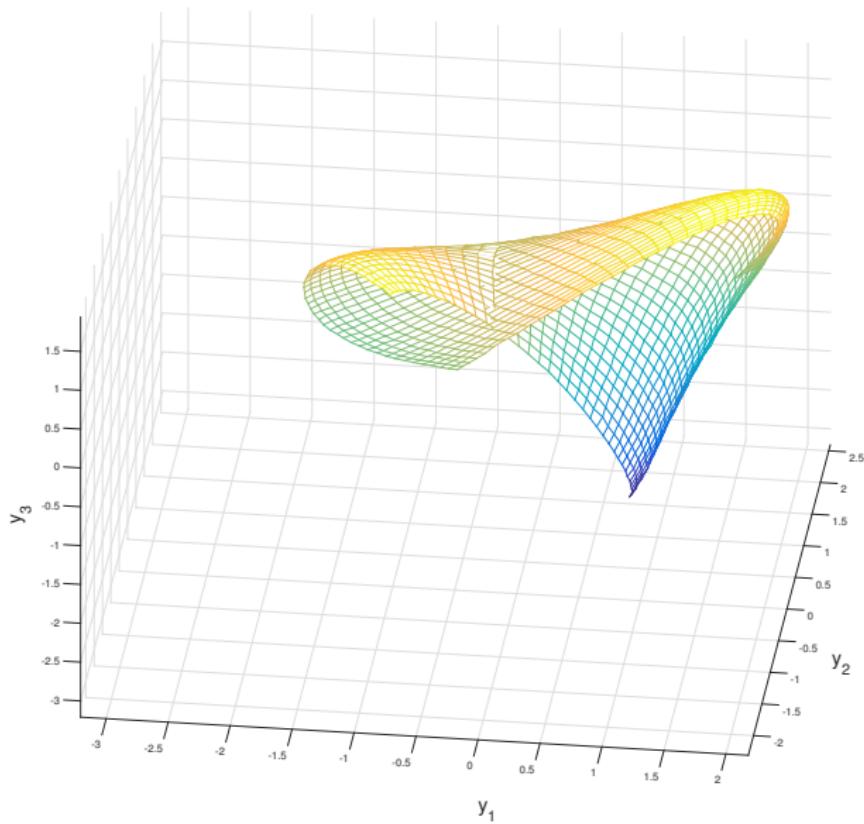
3 dimensional observed  $y$

sample  $f$

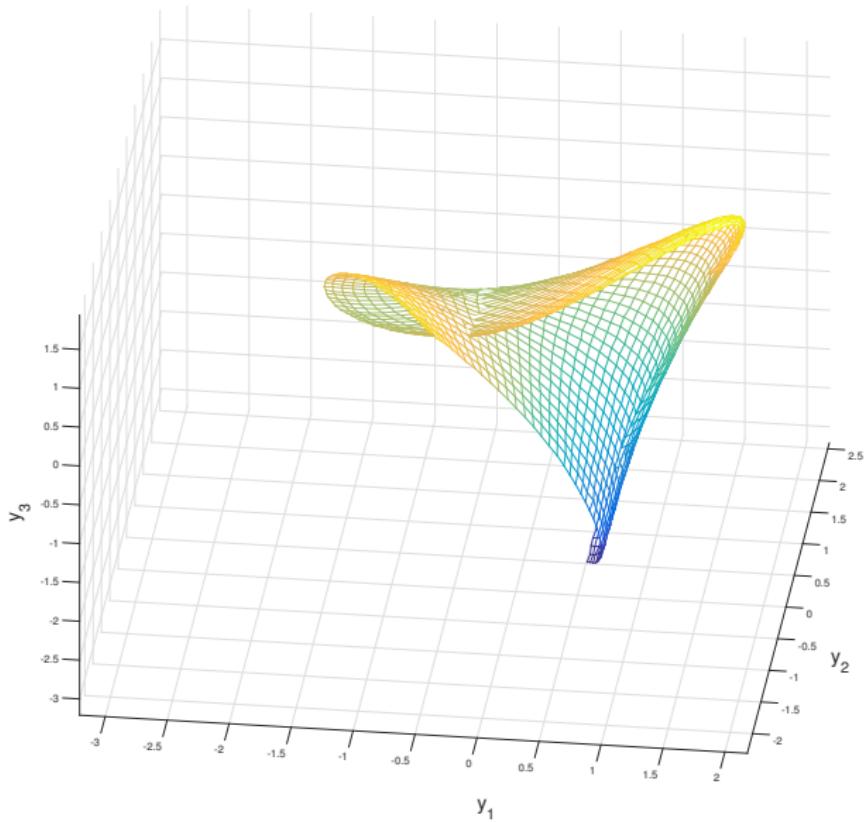
evaluate observed variables  $y$   
corresponding to a grid of  $x$

what does this look like?

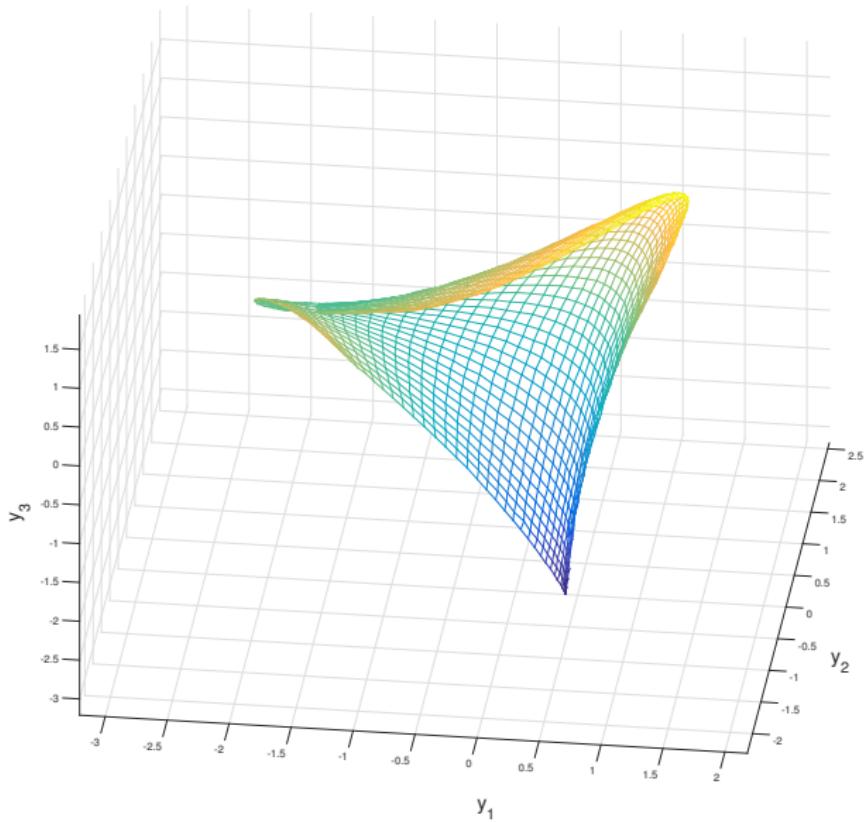
## Gaussian Process Latent Variable Model: manifold samples



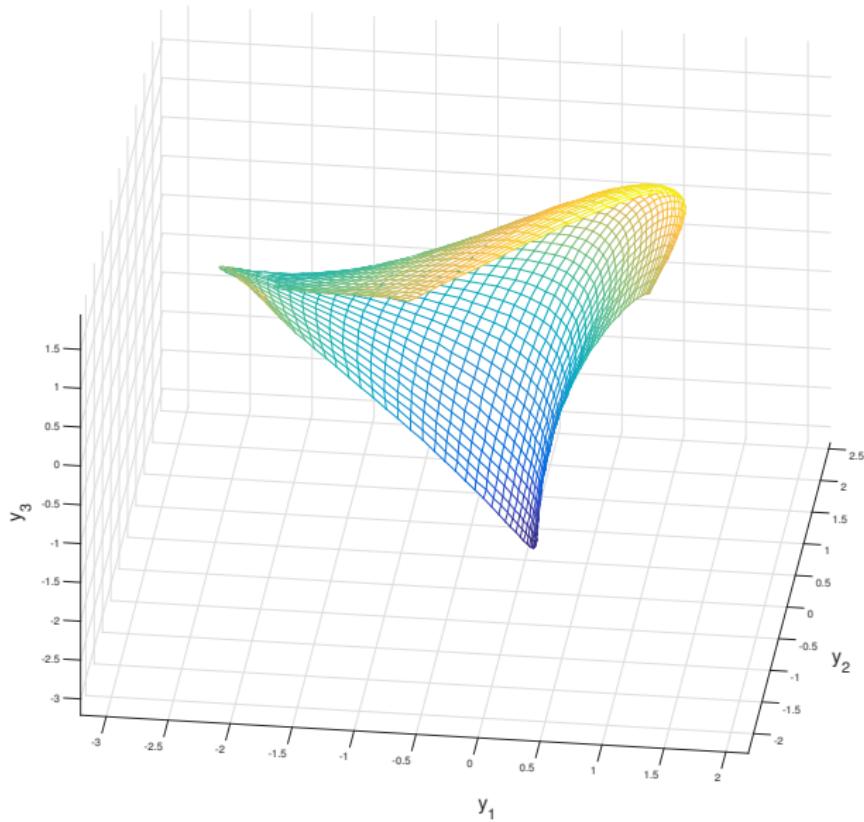
## Gaussian Process Latent Variable Model: manifold samples



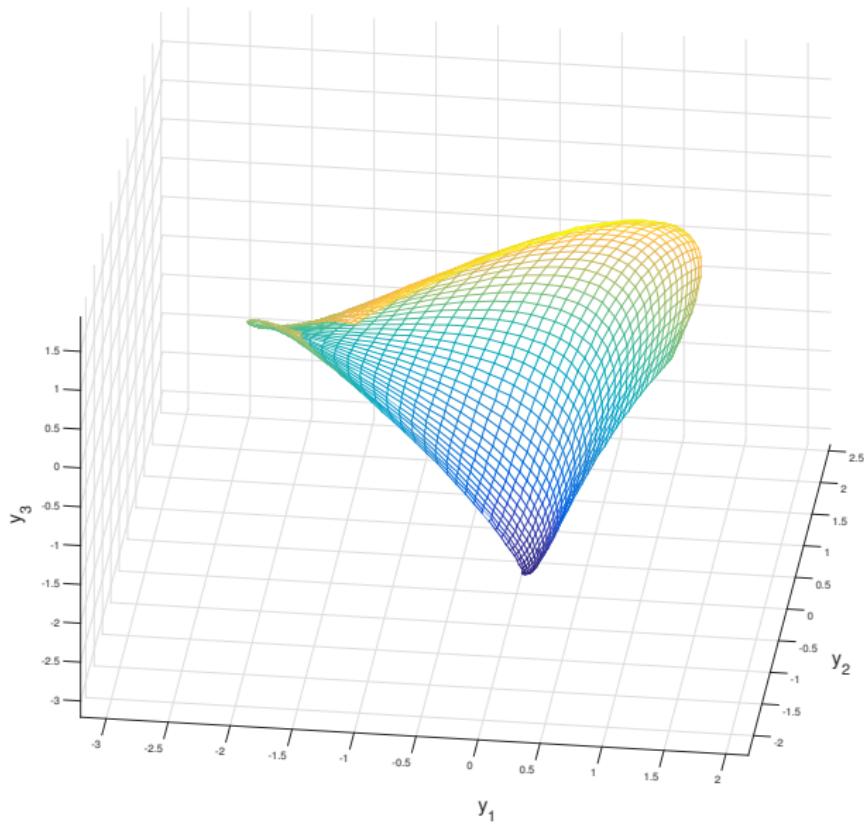
# Gaussian Process Latent Variable Model: manifold samples



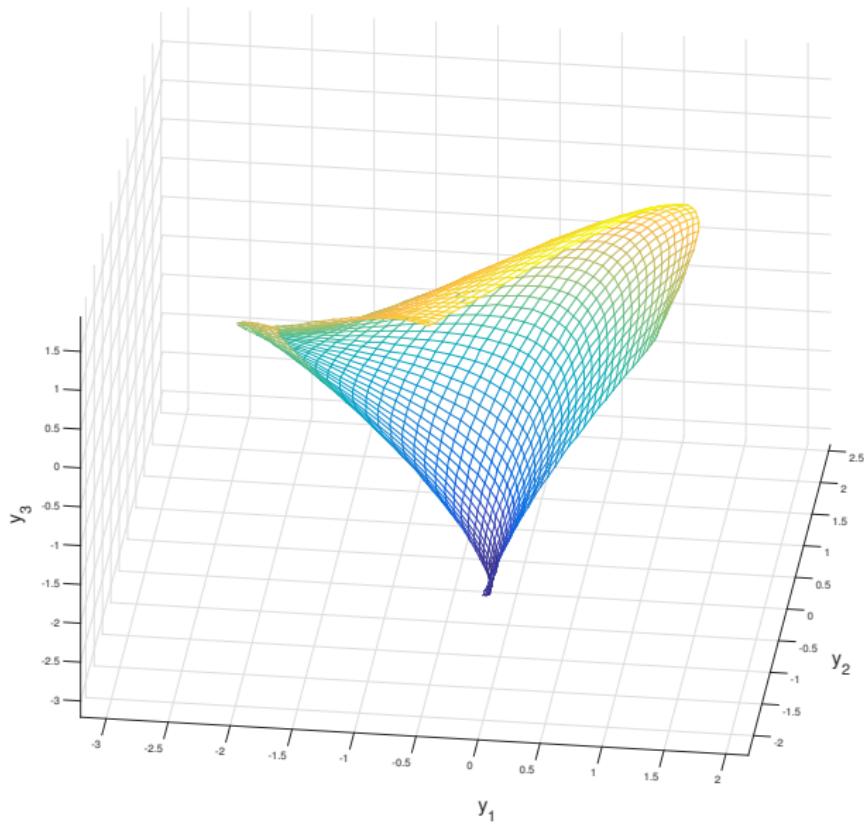
## Gaussian Process Latent Variable Model: manifold samples



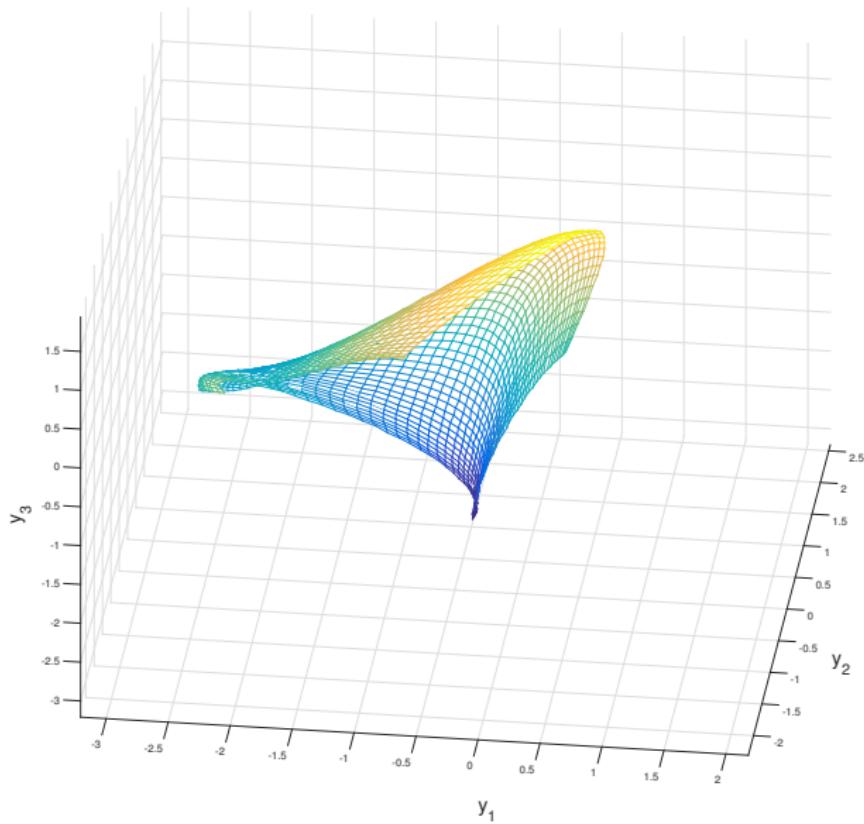
## Gaussian Process Latent Variable Model: manifold samples



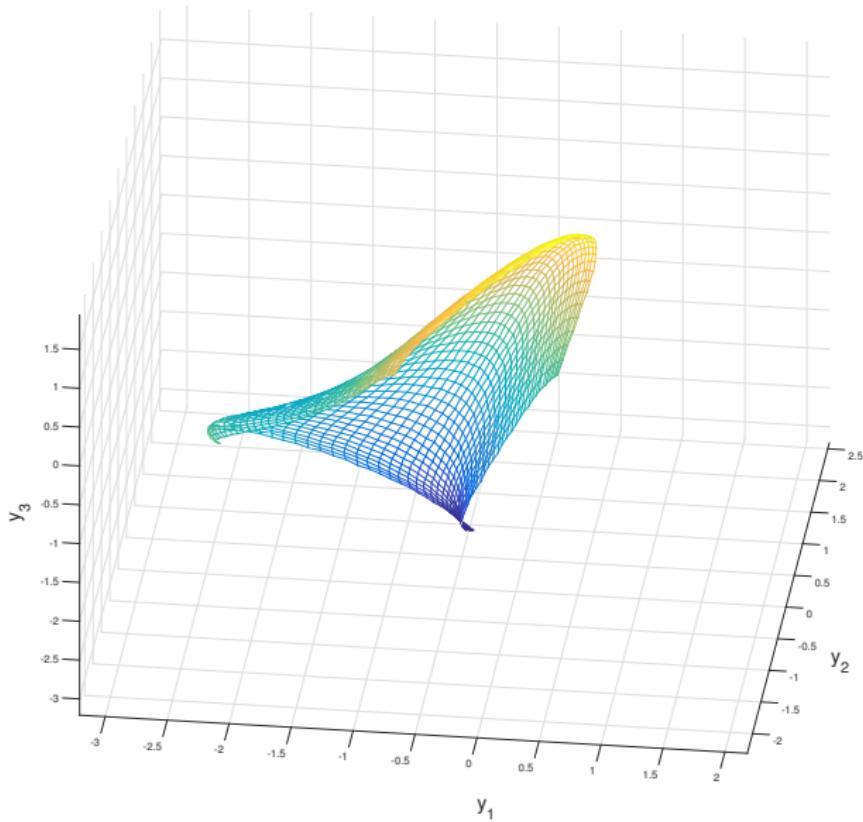
## Gaussian Process Latent Variable Model: manifold samples



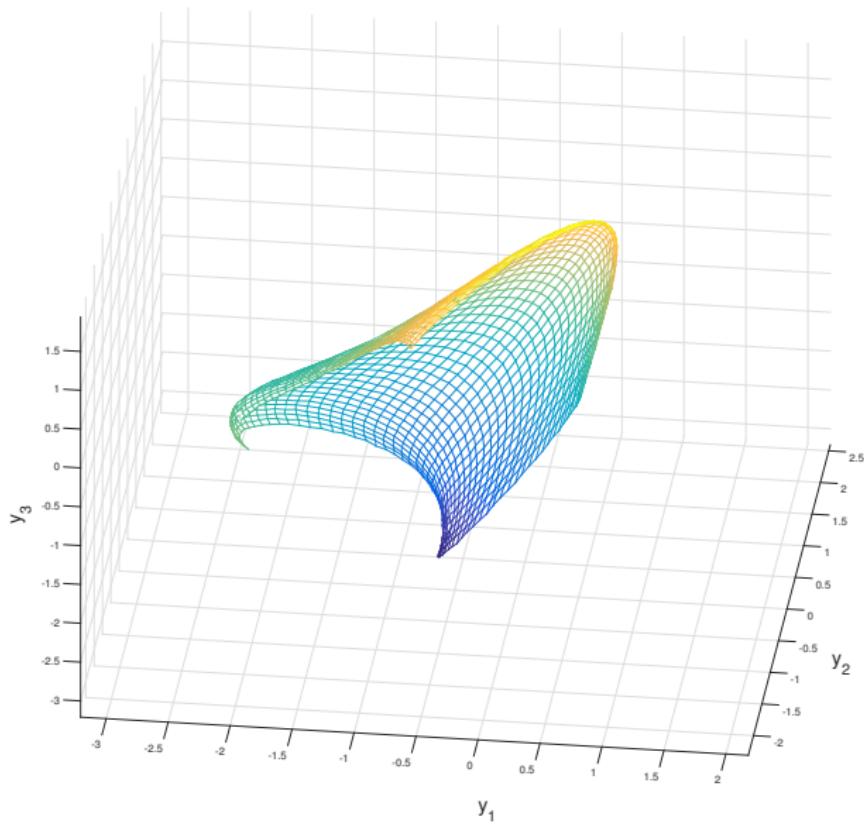
## Gaussian Process Latent Variable Model: manifold samples



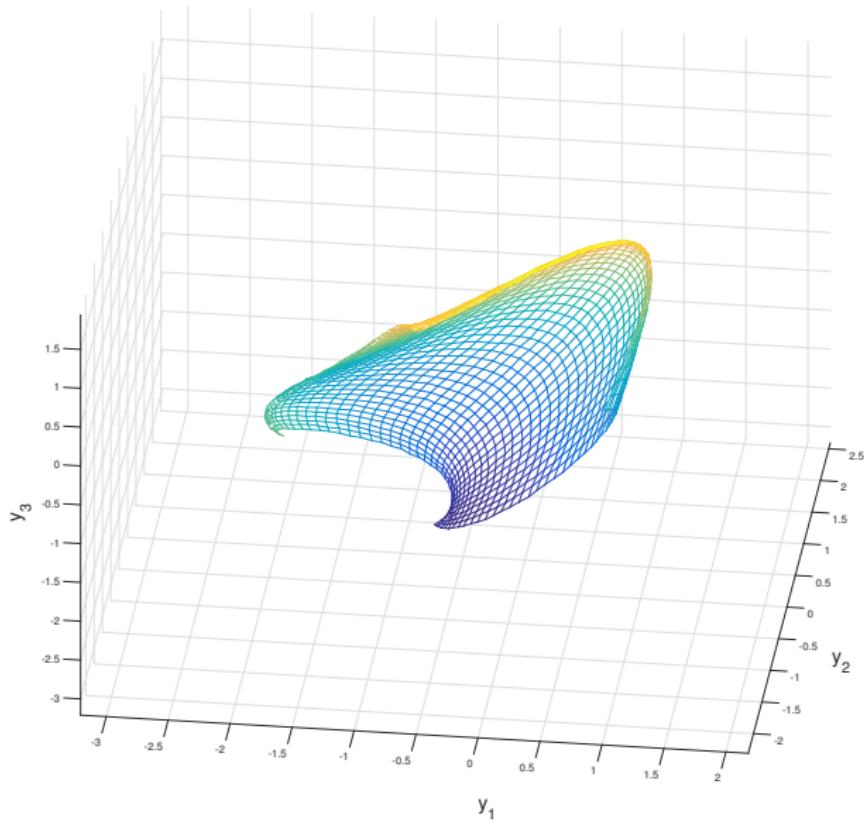
## Gaussian Process Latent Variable Model: manifold samples



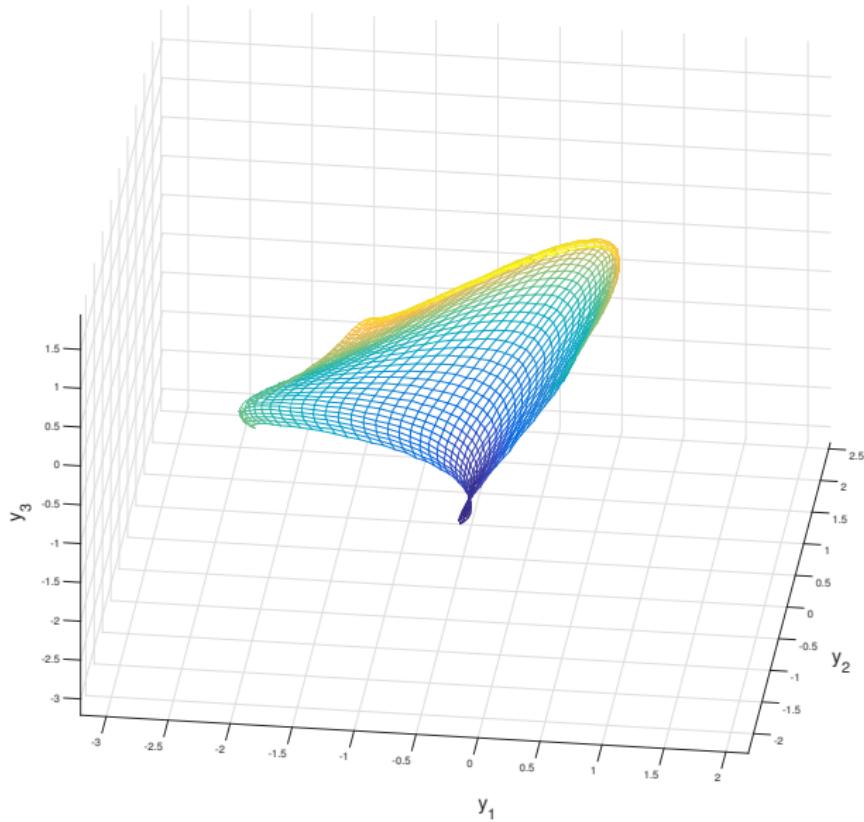
## Gaussian Process Latent Variable Model: manifold samples



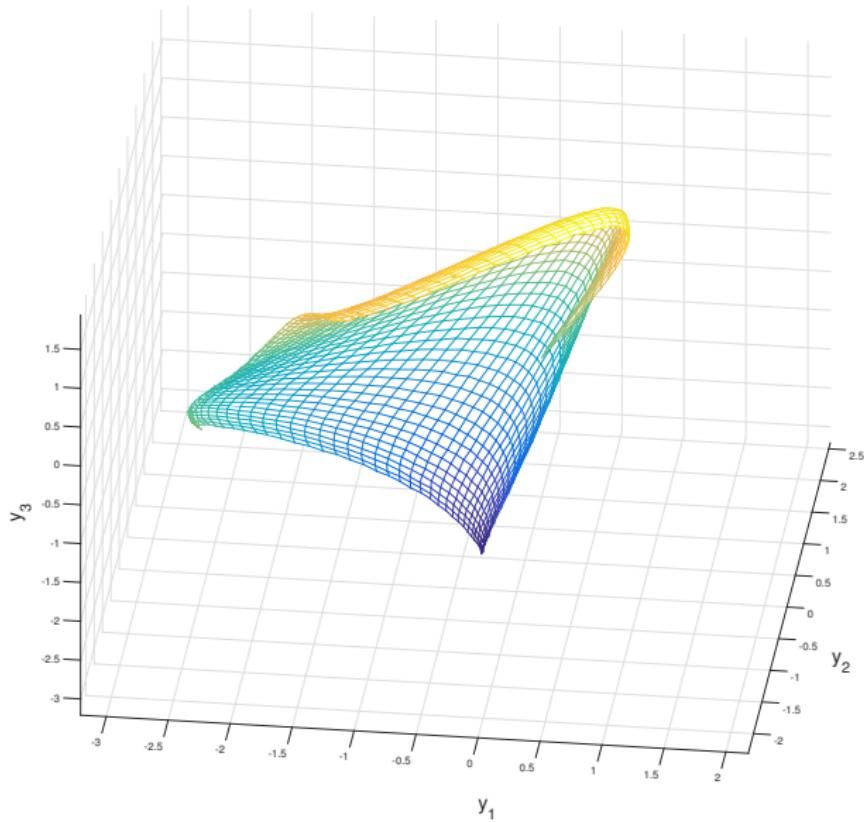
## Gaussian Process Latent Variable Model: manifold samples



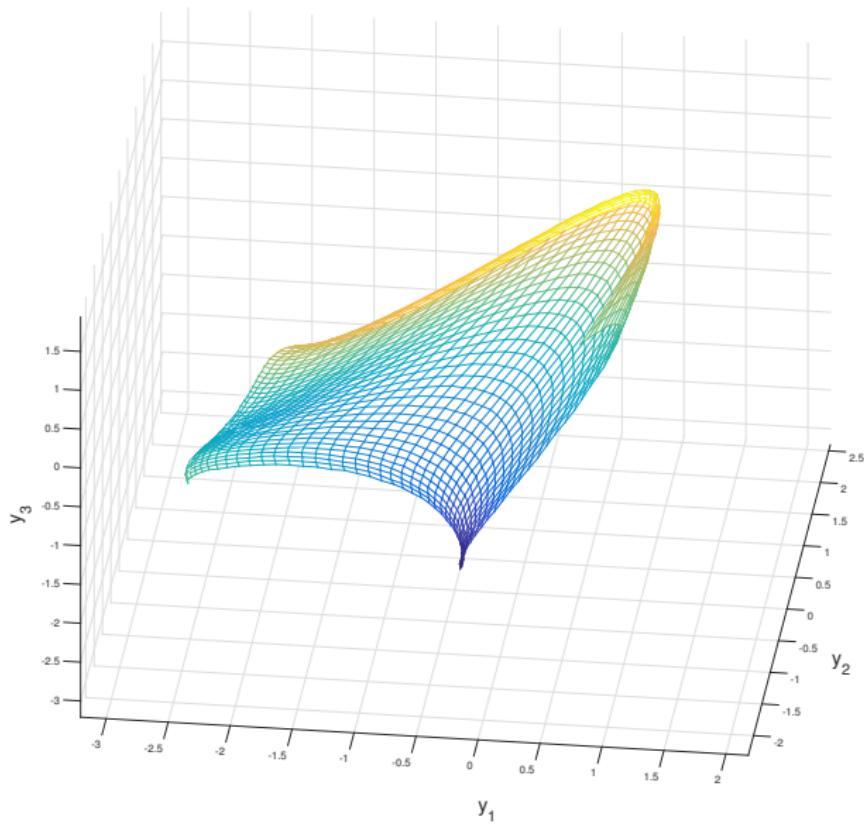
## Gaussian Process Latent Variable Model: manifold samples



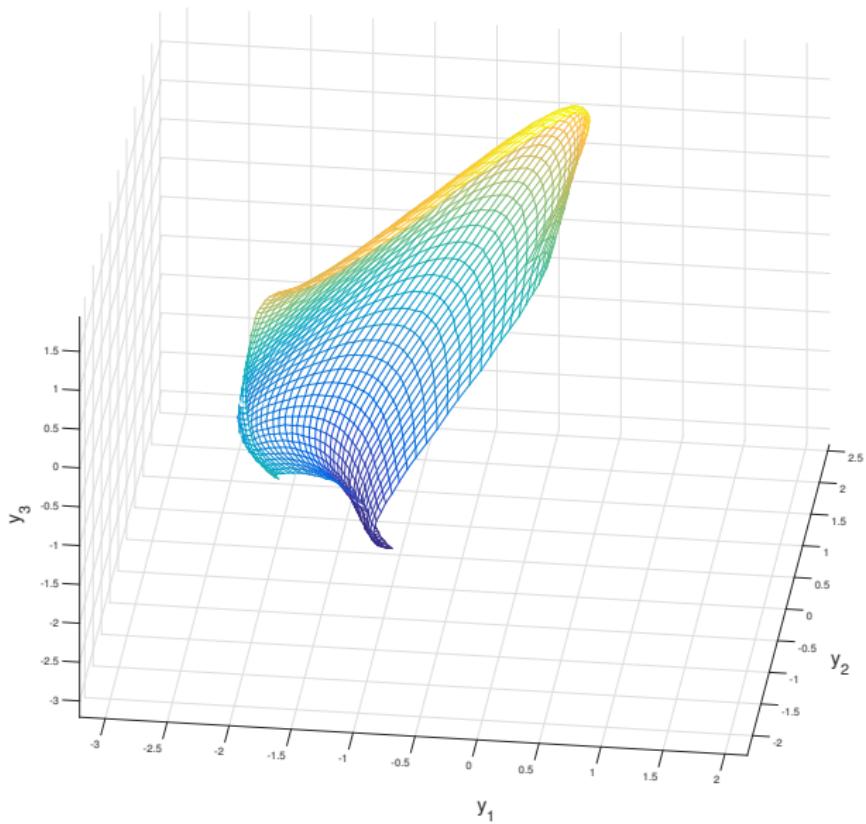
## Gaussian Process Latent Variable Model: manifold samples



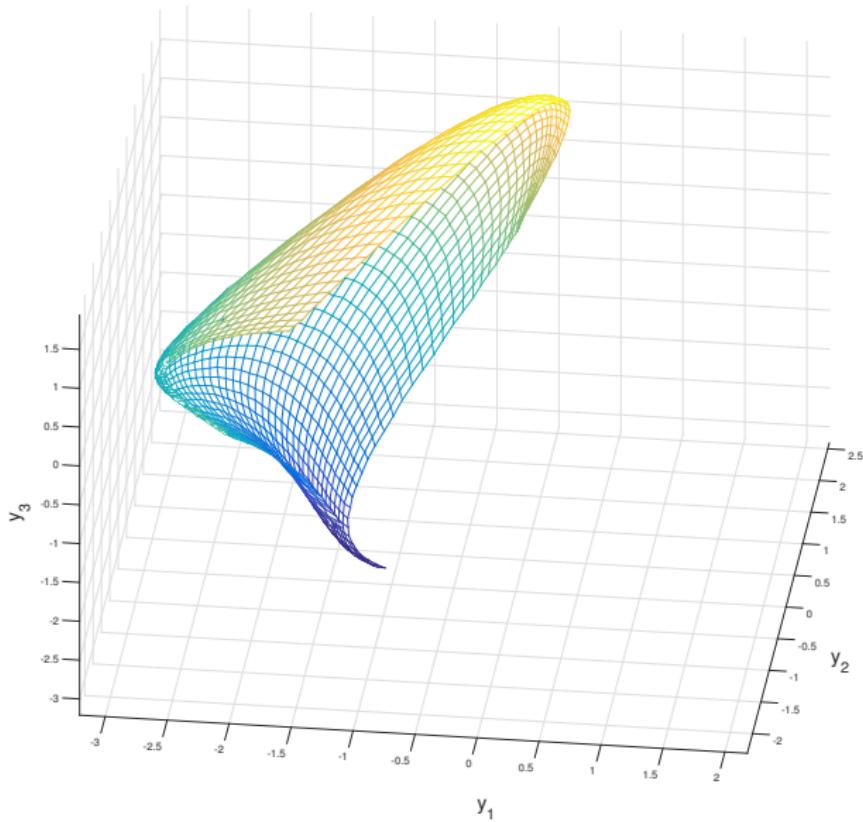
## Gaussian Process Latent Variable Model: manifold samples



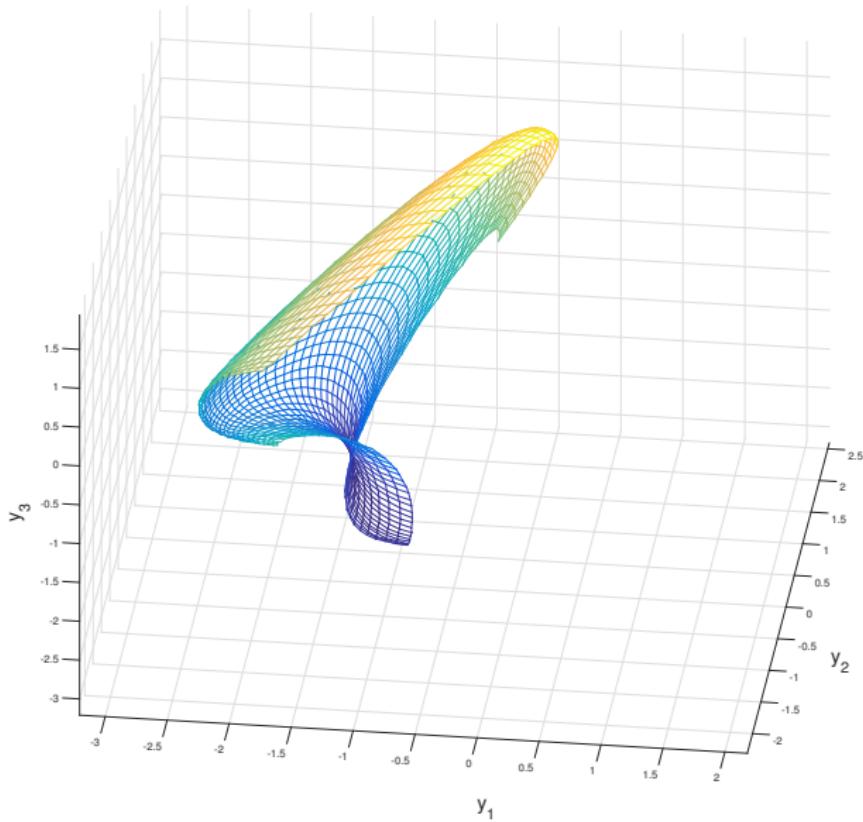
## Gaussian Process Latent Variable Model: manifold samples



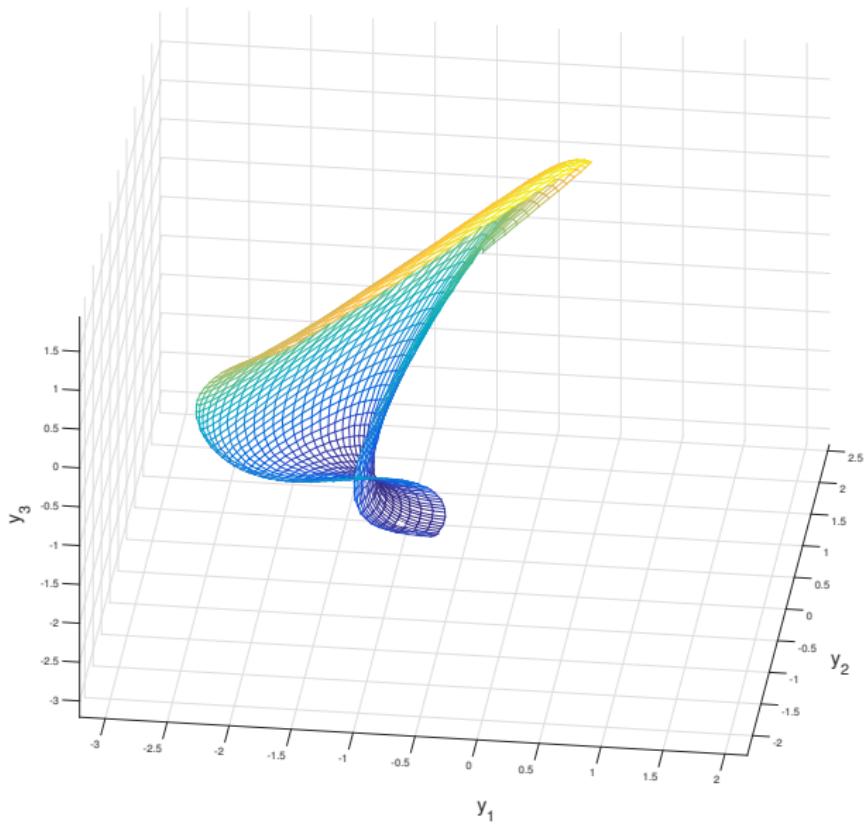
## Gaussian Process Latent Variable Model: manifold samples



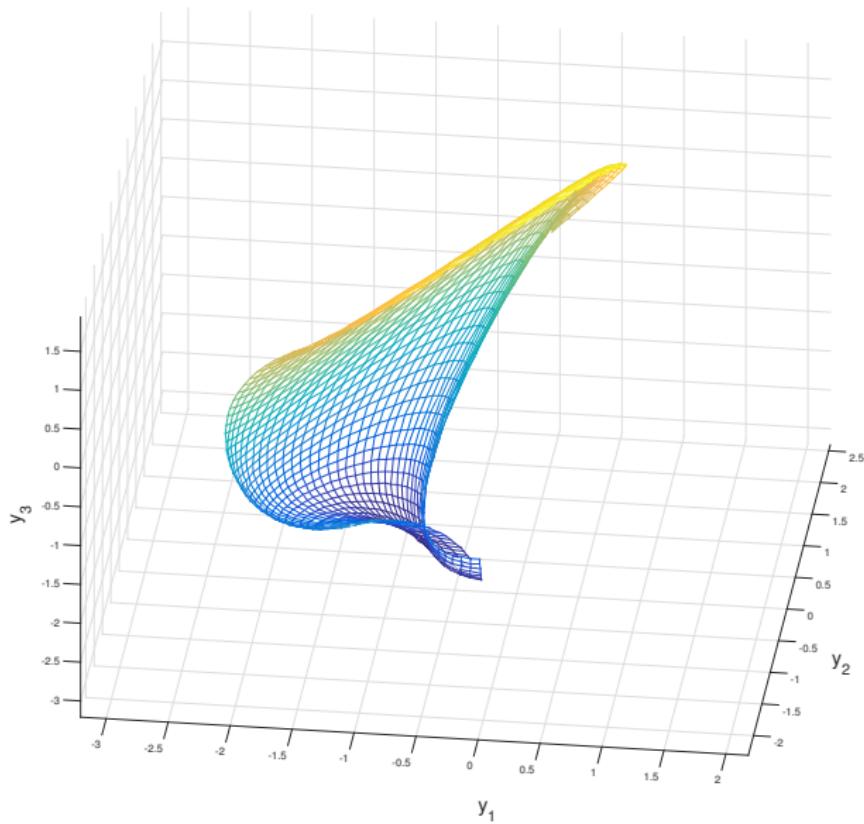
## Gaussian Process Latent Variable Model: manifold samples



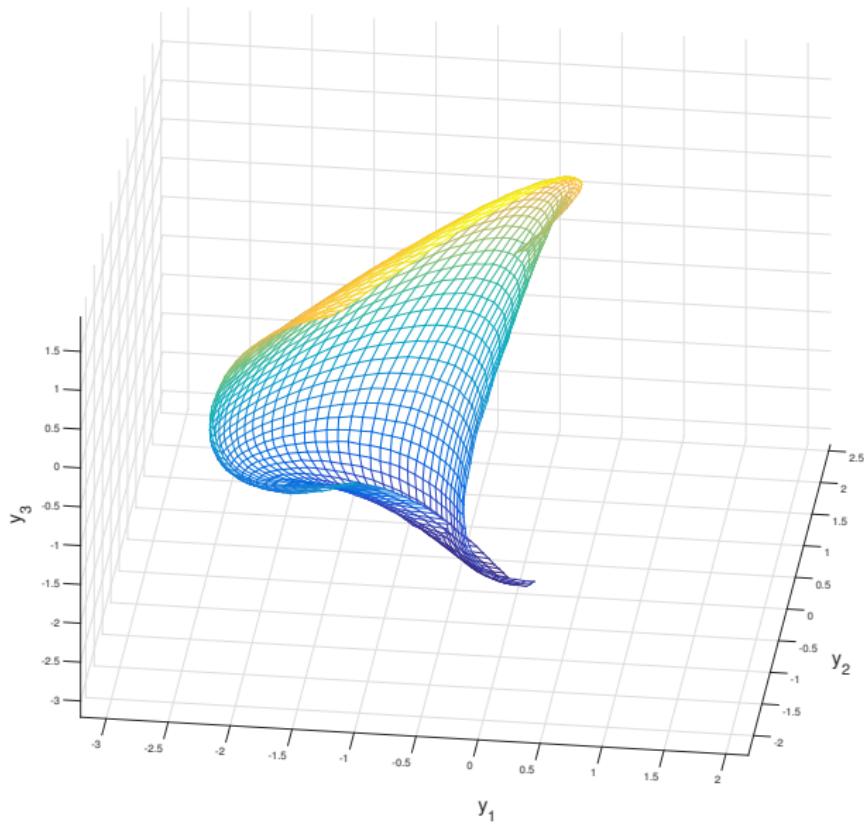
## Gaussian Process Latent Variable Model: manifold samples



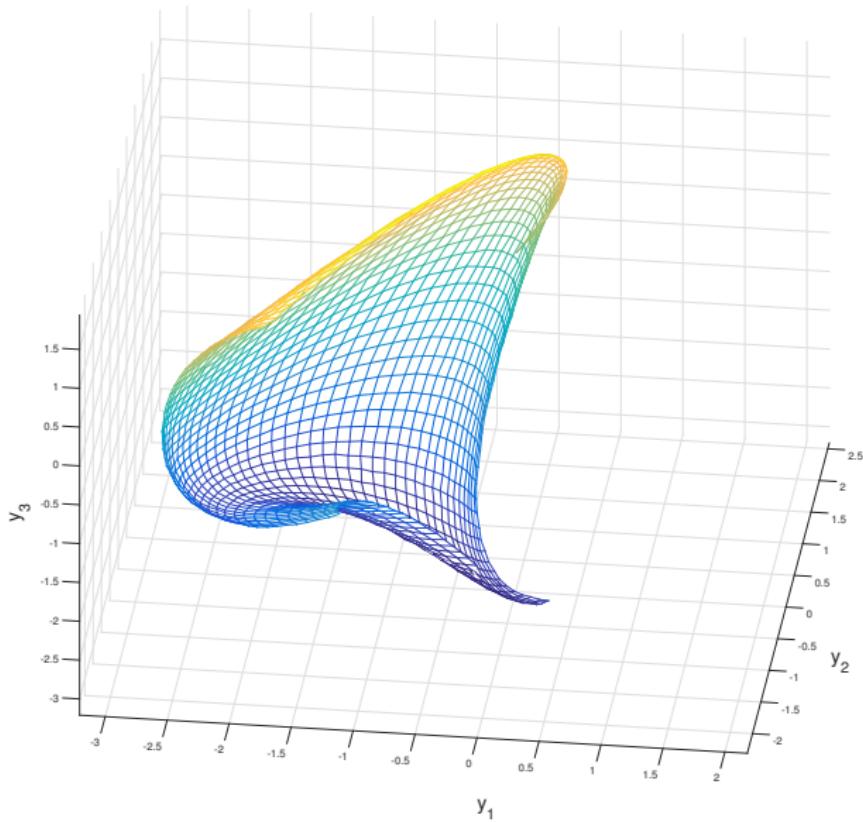
## Gaussian Process Latent Variable Model: manifold samples



## Gaussian Process Latent Variable Model: manifold samples

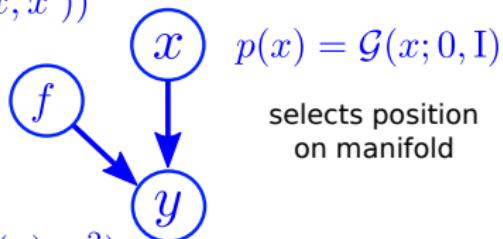


## Gaussian Process Latent Variable Model: manifold samples

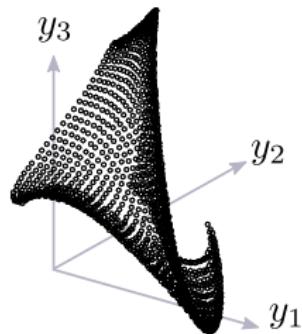


## Gaussian Process Latent Variable Model

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

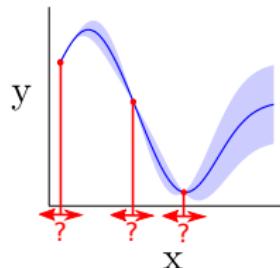


recover PCA using linear covariance  $C(x, x') = \sum_k x_k x'_k$

general case produces complex marginals

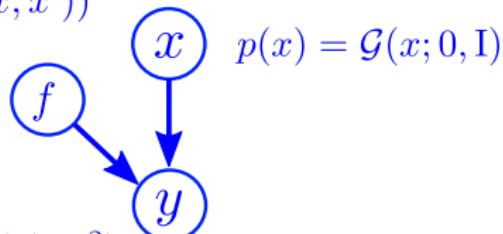
GP-regression with distribution over inputs  
(noisy/latent inputs)

Inference requires approximation...



## Gaussian Process Latent Variable Model: MAP Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(x) = \mathcal{G}(x; 0, I)$$

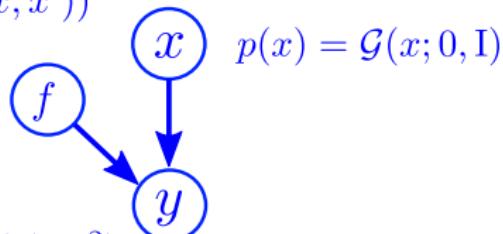
$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence, NIPS 2004

$$x_{\text{MAP}} = \arg \max_x p(x|y) = \arg \max_x \log p(x|y)$$

## Gaussian Process Latent Variable Model: MAP Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

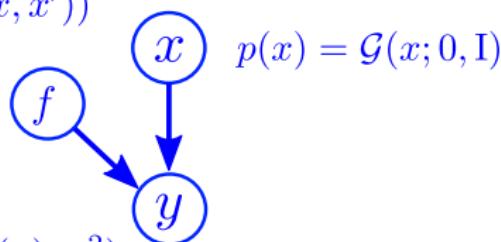
Lawrence, NIPS 2004

$$x_{\text{MAP}} = \arg \max_x p(x|y) = \arg \max_x \log p(x|y)$$

$$p(x|y) = \frac{1}{p(y)} p(y|x)p(x)$$

## Gaussian Process Latent Variable Model: MAP Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence, NIPS 2004

$$x_{\text{MAP}} = \arg \max_x p(x|y) = \arg \max_x \log p(x|y)$$

$$p(x|y) = \frac{1}{p(y)} p(y|x)p(x)$$

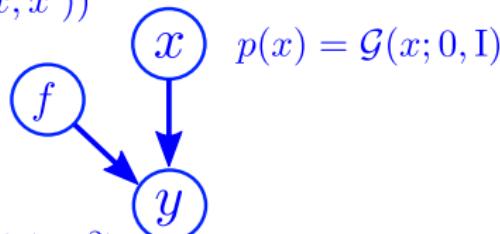
$$p(y|x) = \int p(y|x, f)p(f)df = \mathcal{G}(y_{1:N}; 0, \Sigma(x_{1:N}))$$

$y_{1:N} = [y_{1:N,1}; \dots; y_{1:N,D}]$

multi-output regression

## Gaussian Process Latent Variable Model: MAP Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence, NIPS 2004

$$x_{\text{MAP}} = \arg \max_x p(x|y) = \arg \max_x \log p(x|y)$$

$$p(x|y) = \frac{1}{p(y)} p(y|x)p(x)$$

$$p(y|x) = \int p(y|x, f)p(f)df = \mathcal{G}(y_{1:N}; 0, \Sigma(x_{1:N}))$$

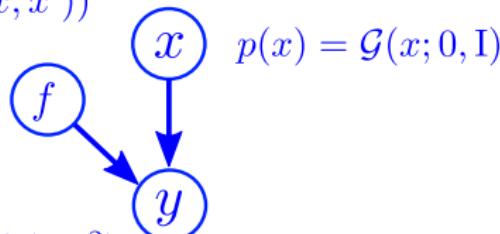
$y_{1:N} = [y_{1:N,1}; \dots; y_{1:N,D}]$

multi-output regression

$$x_{\text{MAP}} = \arg \max_x \log p(x) - \frac{1}{2} \log \det \Sigma(x_{1:N}) - \frac{1}{2} \text{trace}(\Sigma(x_{1:N})^{-1} y_{1:N} y_{1:N}^\top)$$

## Gaussian Process Latent Variable Model: MAP Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence, NIPS 2004

$$x_{\text{MAP}} = \arg \max_x p(x|y) = \arg \max_x \log p(x|y)$$

$$p(x|y) = \frac{1}{p(y)} p(y|x)p(x)$$

$$p(y|x) = \int p(y|x, f)p(f)df = \mathcal{G}(y_{1:N}; 0, \Sigma(x_{1:N}))$$

$y_{1:N} = [y_{1:N,1}; \dots; y_{1:N,D}]$

multi-output regression

$$x_{\text{MAP}} = \arg \max_x \log p(x) - \frac{1}{2} \log \det \Sigma(x_{1:N}) - \frac{1}{2} \text{trace}(\Sigma(x_{1:N})^{-1} y_{1:N} y_{1:N}^\top)$$

objective depends on:

$\|x_n\|$

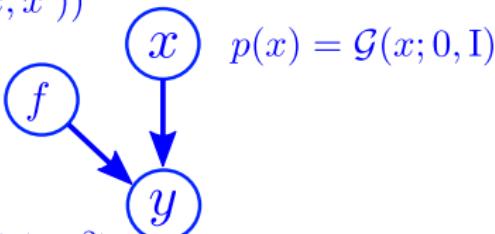
$\|x_n - x_m\|$

$\|y_n - y_m\|$

$\|y_n\|$

# Gaussian Process Latent Variable Model: MAP Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence, NIPS 2004

$$x_{\text{MAP}} = \arg \max_x p(x|y) = \arg \max_x \log p(x|y)$$

$$p(x|y) = \frac{1}{p(y)} p(y|x)p(x)$$

$$p(y|x) = \int p(y|x, f)p(f)df = \mathcal{G}(y_{1:N}; 0, \Sigma(x_{1:N}))$$

$$x_{\text{MAP}} = \arg \max_x \log p(x) - \frac{1}{2} \log \det \Sigma(x_{1:N}) - \frac{1}{2} \text{trace}(\Sigma(x_{1:N})^{-1} y_{1:N} y_{1:N}^\top)$$

objective depends on:  $\|x_n\|$   $\|x_n - x_m\|$   $\|y_n - y_m\|$   $\|y_n\|$

MAP inference in GPLVM  
 -- optimise positions of data in latent space  
 -- objective depends on distances between points in new/old space  
 -- no explicit mappings

ISOMAP-like

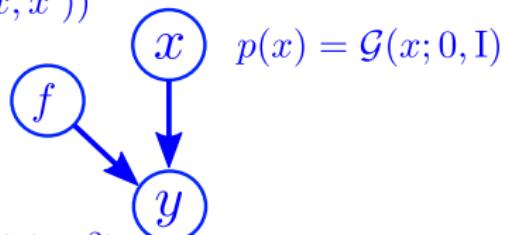
But  
 -- implicit soft decoding function:  
 $p(f|x_{\text{MAP}}, y)$   
 -- cost depends on mags.

$$y_{1:N} = [y_{1:N,1}; \dots; y_{1:N,D}]$$

multi-output regression

## Gaussian Process Latent Variable Model: Back-constrained Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(x) = \mathcal{G}(x; 0, I)$$

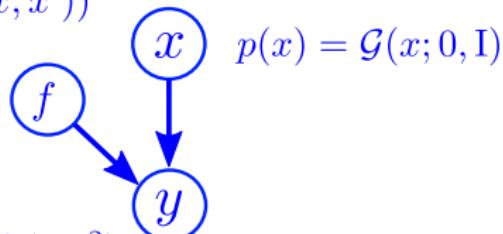
$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence et al, ICML 2006

$$x_{\text{MAP}} = \arg \max_x p(x|y, \theta) = \arg \max_x \log p(x, y|\theta)$$

## Gaussian Process Latent Variable Model: Back-constrained Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence et al, ICML 2006

$$x_{\text{MAP}} = \arg \max_x p(x|y, \theta) = \arg \max_x \log p(x, y|\theta)$$

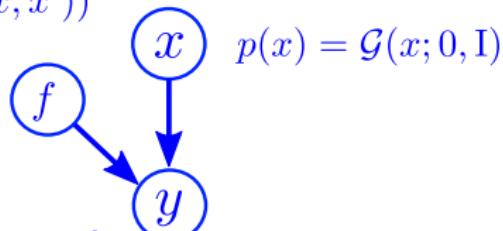
after optimisation:

$$x_{\text{MAP}} = x_{\text{MAP}}(y_{1:N}) \text{ approximate: } x_{\text{MAP}} \approx g_\phi(y_{1:N})$$

recognition model  
or back-constraint

# Gaussian Process Latent Variable Model: Back-constrained Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence et al, ICML 2006

$$x_{\text{MAP}} = \arg \max_x p(x|y, \theta) = \arg \max_x \log p(x, y|\theta)$$

after optimisation:

$$x_{\text{MAP}} = x_{\text{MAP}}(y_{1:N}) \text{ approximate: } x_{\text{MAP}} \approx g_\phi(y_{1:N})$$

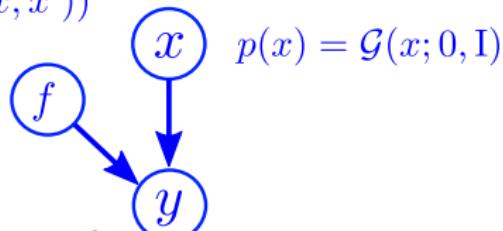
learn using same objective as before:

$$\phi = \arg \max_\phi p(x = g_\phi(y_{1:N})|y, \theta) = \arg \max_\phi \log p(x = g_\phi(y_{1:N}), y|\theta)$$

recognition model  
or back-constraint

# Gaussian Process Latent Variable Model: Back-constrained Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence et al, ICML 2006

$$x_{\text{MAP}} = \arg \max_x p(x|y, \theta) = \arg \max_x \log p(x, y|\theta)$$

after optimisation:

$$x_{\text{MAP}} = x_{\text{MAP}}(y_{1:N}) \text{ approximate: } x_{\text{MAP}} \approx g_\phi(y_{1:N})$$

learn using same objective as before:

$$\phi = \arg \max_\phi p(x = g_\phi(y_{1:N})|y, \theta) = \arg \max_\phi \log p(x = g_\phi(y_{1:N}), y|\theta)$$

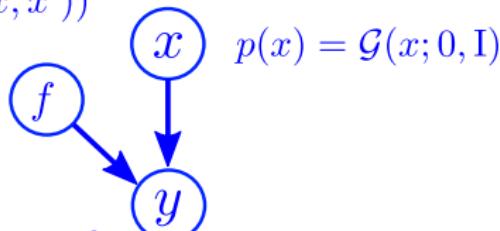
learn hyper-parameters using zero-temperature EM

$$\theta = \arg \max_\theta \log p(y|\theta) \approx \arg \max_\theta \log p(x = g_\phi(y_{1:N}), y|\theta)$$

recognition model  
or back-constraint

# Gaussian Process Latent Variable Model: Back-constrained Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$



$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence et al, ICML 2006

$$x_{\text{MAP}} = \arg \max_x p(x|y, \theta) = \arg \max_x \log p(x, y|\theta)$$

after optimisation:

$$x_{\text{MAP}} = x_{\text{MAP}}(y_{1:N}) \text{ approximate: } x_{\text{MAP}} \approx g_\phi(y_{1:N})$$

learn using same objective as before:

$$\phi = \arg \max_\phi p(x = g_\phi(y_{1:N})|y, \theta) = \arg \max_\phi \log p(x = g_\phi(y_{1:N}), y|\theta)$$

learn hyper-parameters using zero-temperature EM

$$\theta = \arg \max_\theta \log p(y|\theta) \approx \arg \max_\theta \log p(x = g_\phi(y_{1:N}), y|\theta)$$

recognition model  
or back-constraint

optimise same  
function

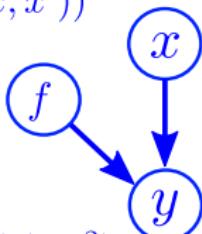
Gaussian Process Latent Variable Model: Back-constrained Inference

$$p(f_d) = \mathcal{GP}(f; 0, C(x, x'))$$

$$p(x) = \mathcal{G}(x; 0, \mathbf{I})$$

$$p(y_d|x, f_d) = \mathcal{G}(y; f_d(x), \sigma^2)$$

Lawrence et al. | ICML 2006



$$x_{\text{MAP}} = \arg \max_x p(x|y, \theta) = \arg \max_x \log p(x, y|\theta)$$

after optimisation:

$x_{\text{MAP}} = x_{\text{MAP}}(y_{1:N})$  approximate:  $x_{\text{MAP}} \approx g_\phi(y_{1:N})$

learn using same objective as before:

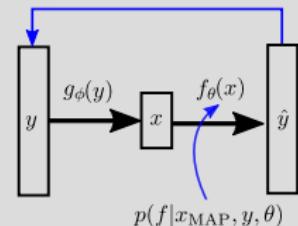
$$\phi = \arg \max_{\phi} p(x = g_{\phi}(y_{1:N}) | y, \theta) = \arg \max_{\phi} \log p(x = g_{\phi}(y_{1:N}), y | \theta)$$

learn hyper-parameters using zero-temperature EM

$$\theta = \arg \max_{\theta} \log p(y|\theta) \approx \arg \max_{\theta} \log p(x = g_{\phi}(y_{1:N}), y|\theta) \quad \text{optimise same function}$$

Back-constrained GPLVM  
Accelerate inference/  
learning via recognition  
model

$$\arg \max_{\theta, \phi} \log p(x = g_\phi(y_{1:N}), y | \theta)$$



auto-encoder-like

recognition model  
or back-constraint

$J(\theta)$  optimise same function

# Gaussian Process Latent Variable Model: Application

