

# Deep Gaussian Process Autoencoders for Novelty Detection

---

R. Domingues <sup>1</sup>, P. Michiardi <sup>1</sup>, J. Zouaoui <sup>2</sup>, M. Filippone <sup>1</sup>

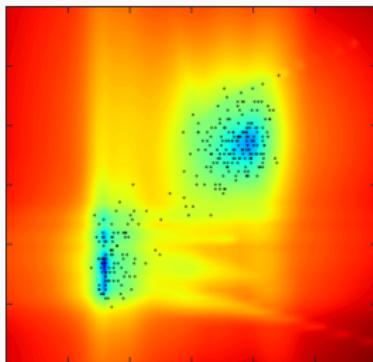
<sup>1</sup> EURECOM, Sophia Antipolis, France

<sup>2</sup> Amadeus, Sophia Antipolis, France

ECML-PKDD 2018 - Journal track  
September 10<sup>th</sup> 2018

# Novelty detection

- Recognition of anomalies in test data which differ significantly from the training set
- Estimate the distribution of nominal samples
- Similar to a one-class classification problem
- Usually addressed by unsupervised learning
- Training set may be contaminated by outliers



# Deep Gaussian Process autoencoder

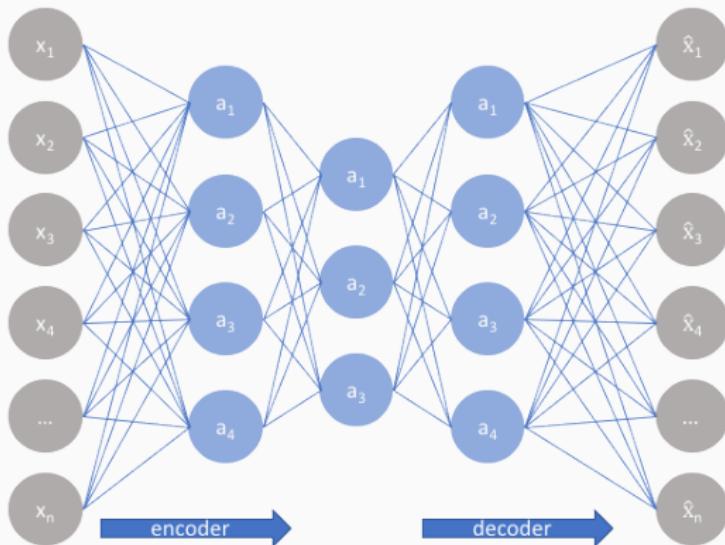
- Unsupervised and probabilistic
- Suitable for any type of data
- Training only requires tensor products
- Inference through stochastic variational inference
- Mini-batch learning

## DGP-AE Architecture

---

# Autoencoders

- Learn a compressed representation of the training data by minimizing the error between the input data and the reconstructed output

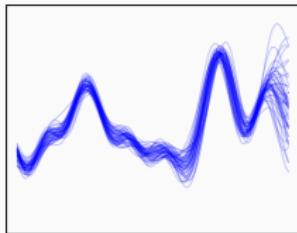


# Deep Gaussian Process autoencoders

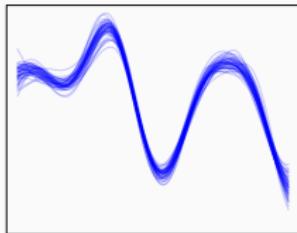
- Deep probabilistic models
- Composition of functions

$$f(\mathbf{x}) = \left( h^{(N_h-1)} \left( \theta^{(N_h-1)} \right) \circ \dots \circ h^{(0)} \left( \theta^{(0)} \right) \right) (\mathbf{x})$$

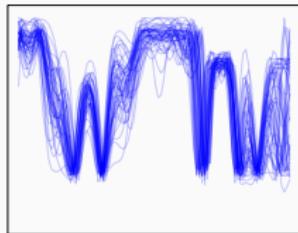
$h^{(0)}(\mathbf{x})$



$h^{(1)}(\mathbf{x})$



$h^{(1)}(h^{(0)}(\mathbf{x}))$



## Marginal likelihood

- Inference requires calculating the marginal likelihood:

$$\begin{aligned} p(\mathcal{X}|\boldsymbol{\theta}) &= \int p\left(\mathcal{X}|\mathcal{F}^{(N_L)}, \boldsymbol{\theta}^{(N_L)}\right) \times \\ &\quad p\left(\mathcal{F}^{(N_L)}|\mathcal{F}^{(N_L-1)}, \boldsymbol{\theta}^{(N_L-1)}\right) \times \dots \times \\ &\quad p\left(\mathcal{F}^{(1)}|\mathcal{F}^{(N_0)}, \boldsymbol{\theta}^{(0)}\right) d\mathcal{F}^{(N_L)} \dots d\mathcal{F}^{(1)} \end{aligned}$$

## Model inference

---

# DGPs with Random Features

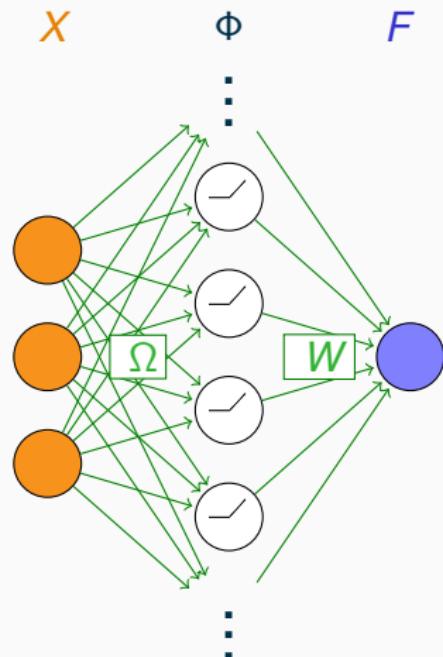
- GPs are single-layered Neural Nets with an infinite number of hidden units

- Weight-space view of a GP

$$F = \Phi W$$

- The priors over the weights are

$$p(W_{\cdot i}) = \mathcal{N}(\mathbf{0}, I)$$



# Random Feature Expansion of Kernels

- Low-rank approximation of GP covariance functions
- The **RBF kernel** can be approximated using **trigonometric functions**

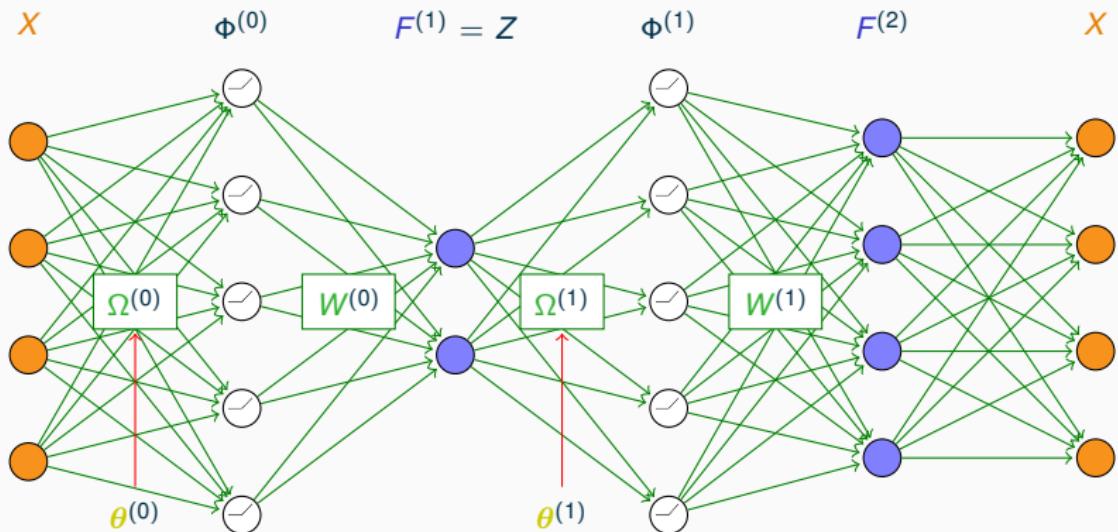
$$\Phi_{\text{RBF}} = \sqrt{\frac{\sigma^2}{N_{\text{RF}}}} [\cos(F\Omega), \sin(F\Omega)] \quad \text{with} \quad p(\Omega_j | \theta) = \mathcal{N}(\mathbf{0}, \Lambda^{-1})$$

- The first order **Arc-Cosine kernel** can be approximated using **Rectified Linear Units (ReLU)**

$$\Phi_{\text{ARC}} = \sqrt{\frac{2\sigma^2}{N_{\text{RF}}}} \max(0, F\Omega) \quad \text{with} \quad p(\Omega_j | \theta) = \mathcal{N}(\mathbf{0}, \Lambda^{-1})$$

- Approximated multivariate GPs are **Bayesian linear models**

# DGP-AEs with RFs (2 layers)



## DGP-AE with RFs - Stochastic Variational Inference

- Define  $\Psi = (\Omega^{(0)}, \dots, \Omega^{(L)}, W^{(0)}, \dots, W^{(L)})$
- Lower bound on the marginal likelihood:

$$\log [p(\mathcal{X}|\theta)] \geq \mathbb{E}_{q(\Psi)} (\log [p(\mathcal{X}|\Psi, \theta)]) - D_{\text{KL}} [q(\Psi) \| p(\Psi)]$$

where  $q(\Psi)$  approximates  $p(\Psi|\mathcal{X}, \theta)$

- $D_{\text{KL}}$  computable analytically if  $q$  and  $p$  are Gaussian
- We assume an approximate factorized Gaussian distribution  $q(\Psi)$

## DGPs with RFs - Stochastic variational inference

- Stochastic **unbiased** estimate of the expectation term
  - **Mini-batch**

$$\mathbb{E}_{q(\Psi)} (\log [p(\textcolor{brown}{X}|\Psi, \theta)]) \approx \frac{n}{m} \sum_{k \in \mathcal{I}_m} \mathbb{E}_{q(\Psi)} (\log [p(\textcolor{brown}{x}_k|\Psi, \theta)])$$

# DGPs with RFs - Stochastic variational inference

- Stochastic **unbiased** estimate of the expectation term
  - **Mini-batch**

$$\mathbb{E}_{q(\Psi)} (\log [p(\mathcal{X}|\Psi, \theta)]) \approx \frac{n}{m} \sum_{k \in \mathcal{I}_m} \mathbb{E}_{q(\Psi)} (\log [p(\mathbf{x}_k|\Psi, \theta)])$$

- **Monte Carlo sampling**

$$\mathbb{E}_{q(\Psi)} (\log [p(\mathbf{x}_k|\Psi, \theta)]) \approx \frac{1}{N_{MC}} \sum_{r=1}^{N_{MC}} \log [p(\mathbf{x}_k|\tilde{\Psi}_r, \theta)]$$

with  $\tilde{\Psi}_r \sim q(\Psi)$

- The derivative of the estimate yields a **stochastic gradient**

## DGPs with RFs - Stochastic variational inference

- Reparameterization trick

$$\left( \tilde{W}_r^{(l)} \right)_{ij} = s_{ij}^{(l)} \epsilon_{rij}^{(l)} + m_{ij}^{(l)}$$

with  $\epsilon_{rij}^{(l)} \sim \mathcal{N}(0, 1)$

# Predictions

- Predictive distribution

$$p(\mathbf{x}_* | \mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{x}_* | \boldsymbol{\psi}, \boldsymbol{\theta}) p(\boldsymbol{\psi} | \mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\psi}$$

- Approximation

$$\begin{aligned} p(\mathbf{x}_* | \mathbf{X}, \boldsymbol{\theta}) &\approx \int p(\mathbf{x}_* | \boldsymbol{\psi}, \boldsymbol{\theta}) q(\boldsymbol{\psi}) d\boldsymbol{\psi} \\ &\approx \frac{1}{N_{\text{MC}}} \sum_{r=1}^{N_{\text{MC}}} p(\mathbf{x}_* | \tilde{\boldsymbol{\psi}}_r, \boldsymbol{\theta}) \end{aligned}$$

# Likelihood functions

- Model inference for mixed-type features
  - Normal:  $p(\mathbf{x}_{[G]} | \mathbf{f}^{(N_L)}) = \mathcal{N}(\mathbf{x}_{[G]} | \mathbf{f}_{[G]}^{(N_L)}, \sigma_{[G]}^2)$
  - Softmax:  $p((\mathbf{x}_{[C]})_j | \mathbf{f}^{(N_L)}) = \frac{\exp[(\mathbf{f}_{[C]}^{(N_L)})_j]}{\sum_i \exp[(\mathbf{f}_{[C]}^{(N_L)})_i]}$
  - Combined likelihood:  $p(\mathbf{x} | \mathbf{f}^{(N_L)}) = \prod_k p(\mathbf{x}_{[k]} | \mathbf{f}^{(N_L)})$

## Evaluation

---

# Algorithms

- **Isolation Forest:** IFOREST (Liu et al. 2008)
- **Robust Kernel Density Estimation:** RKDE (Kim and Scott 2012)
- **Feedforward Autoencoders:** AE-1, AE-5
- **Variational Autoencoders:** VAE-1, VAE-2 (Kingma and Welling 2014)
- **Variational Auto-Encoded DGP:** VAE-DGP-2 (Dai et al. 2016)
- **Neural Autoregressive Distribution Estimator:** NADE-2 (Uria et al. 2016)

# Algorithms

- **Isolation Forest:** IFOREST (Liu et al. 2008)
- **Robust Kernel Density Estimation:** RKDE (Kim and Scott 2012)
- **Feedforward Autoencoders:** AE-1, AE-5
- **Variational Autoencoders:** VAE-1, VAE-2 (Kingma and Welling 2014)
- **Variational Auto-Encoded DGP:** VAE-DGP-2 (Dai et al. 2016)
- **Neural Autoregressive Distribution Estimator:** NADE-2 (Uria et al. 2016)

# Algorithms

- **Isolation Forest:** IFOREST (Liu et al. 2008)
- **Robust Kernel Density Estimation:** RKDE (Kim and Scott 2012)
- **Feedforward Autoencoders:** AE-1, AE-5
- **Variational Autoencoders:** VAE-1, VAE-2 (Kingma and Welling 2014)
- **Variational Auto-Encoded DGP:** VAE-DGP-2 (Dai et al. 2016)
- **Neural Autoregressive Distribution Estimator:** NADE-2 (Uria et al. 2016)

# Algorithms

- **Isolation Forest:** IFOREST (Liu et al. 2008)
- **Robust Kernel Density Estimation:** RKDE (Kim and Scott 2012)
- **Feedforward Autoencoders:** AE-1, AE-5
- **Variational Autoencoders:** VAE-1, VAE-2 (Kingma and Welling 2014)
- **Variational Auto-Encoded DGP:** VAE-DGP-2 (Dai et al. 2016)
- **Neural Autoregressive Distribution Estimator:** NADE-2 (Uria et al. 2016)

# Algorithms

- **Isolation Forest:** IFOREST (Liu et al. 2008)
- **Robust Kernel Density Estimation:** RKDE (Kim and Scott 2012)
- **Feedforward Autoencoders:** AE-1, AE-5
- **Variational Autoencoders:** VAE-1, VAE-2 (Kingma and Welling 2014)
- **Variational Auto-Encoded DGP:** VAE-DGP-2 (Dai et al. 2016)
- **Neural Autoregressive Distribution Estimator:** NADE-2 (Uria et al. 2016)

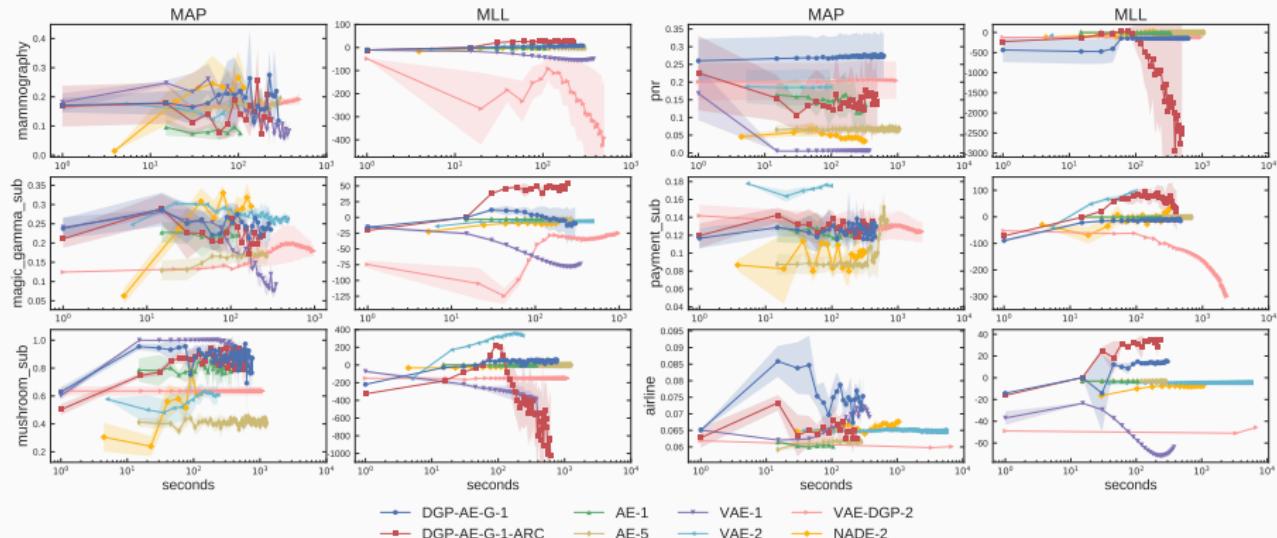
# Method comparison

- 11 datasets, mean area under the precision-recall curve (MAP)
- Some datasets contain over 3 millions samples and 100 features
- DGP-AE achieves the best results for novelty detection
- Softmax accurately models categorical variables

	DGP-AE	DGP-AE	DGP-AE	DGP-AE	VAE-DGP-2	AE-1	AE-5	VAE-1	VAE-2	NADE-2	RKDE	IFOREST
	G-1	G-2	GS-1	GS-2								
MAMMOGRAPHY	<b>0.222</b>	0.183	<b>0.222</b>	0.183	<b>0.221</b>	0.118	0.075	0.119	0.148	0.193	<b>0.231</b>	<b>0.244</b>
MAGIC-GAMMA-SUB	0.260	0.340	0.260	0.340	0.235	0.253	0.125	0.230	0.305	<b>0.398</b>	<b>0.402</b>	0.290
WINE-QUALITY	<b>0.224</b>	<b>0.203</b>	<b>0.224</b>	<b>0.203</b>	0.075	0.106	0.042	0.064	0.124	0.102	0.051	0.059
MUSHROOM-SUB	0.811	0.677	<b>0.940</b>	0.892	0.636	0.725	0.331	0.758	0.479	0.596	0.839	0.546
CAR	0.050	0.061	0.043	0.067	0.045	0.044	0.032	<b>0.071</b>	0.050	0.030	0.034	0.041
GERMAN-SUB	0.066	0.077	<b>0.106</b>	0.098	<b>0.113</b>	0.065	<b>0.103</b>	<b>0.104</b>	0.062	<b>0.118</b>	<b>0.109</b>	0.079
PNR	<b>0.190</b>	0.172	<b>0.190</b>	0.172	<b>0.201</b>	0.059	0.107	0.100	0.106	0.006	0.146	0.124
TRANSACTIONS	0.756	0.752	<b>0.810</b>	<b>0.835</b>	0.509	0.563	0.510	0.532	0.760	0.373	0.585	0.564
SHARED-ACCESS	0.692	0.738	0.692	0.738	0.668	0.546	<b>0.766</b>	0.471	0.527	0.239	<b>0.783</b>	0.746
PAYMENT-SUB	<b>0.173</b>	<b>0.173</b>	0.168	0.168	0.137	0.157	0.129	<b>0.175</b>	0.143	0.101	<b>0.180</b>	0.142
AIRLINE	<b>0.081</b>	<b>0.079</b>	<b>0.081</b>	<b>0.079</b>	0.060	0.063	0.059	0.068	0.074	0.064	-	0.069
AVERAGE	0.344	0.338	0.366	0.370	0.284	0.264	0.222	0.262	0.270	0.216	0.336	0.284

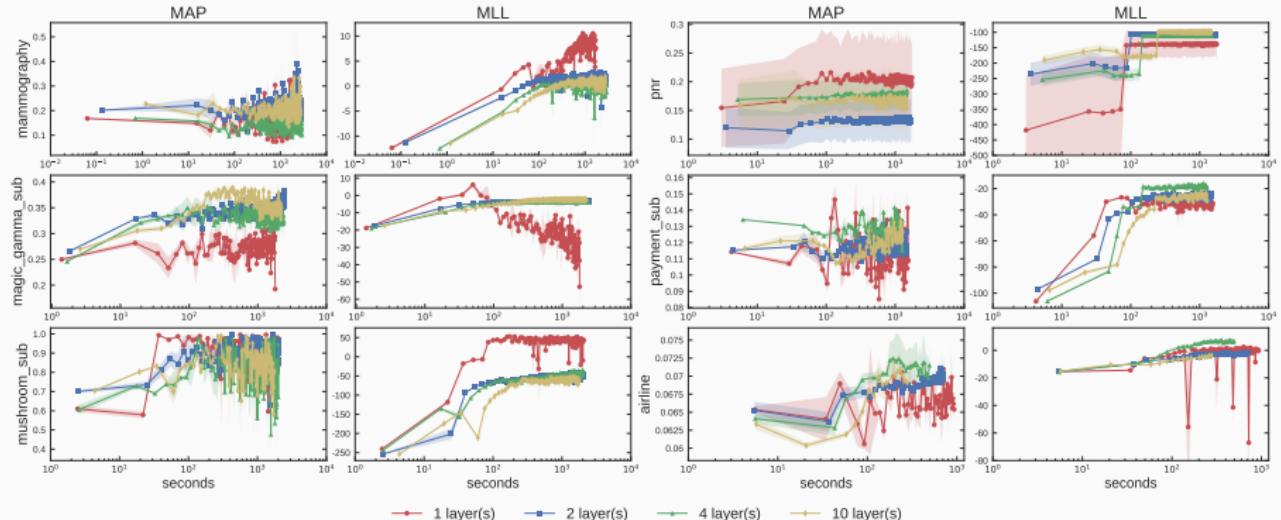
# Convergence monitoring - Networks

- MAP and mean log-likelihood (MLL). The higher the better



- DGP-AE shows the best likelihood
- MAP quickly stabilizes while the likelihood is continuously refined

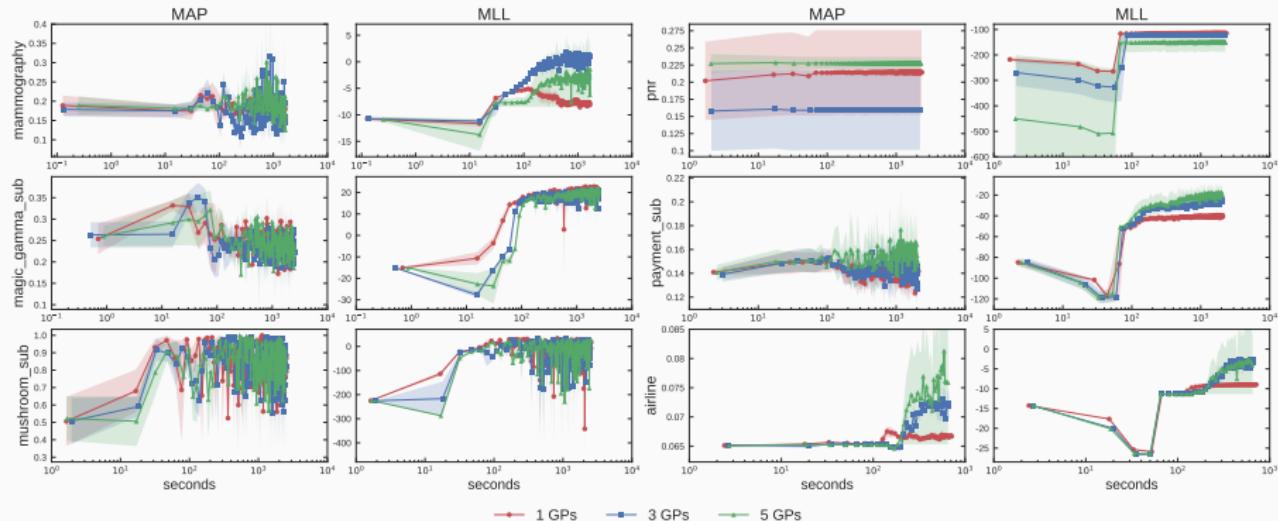
# Convergence monitoring - Depth



- Correlation between a higher test likelihood and a higher MAP
- Moderately deep networks capture the complexity of data without an important convergence overhead

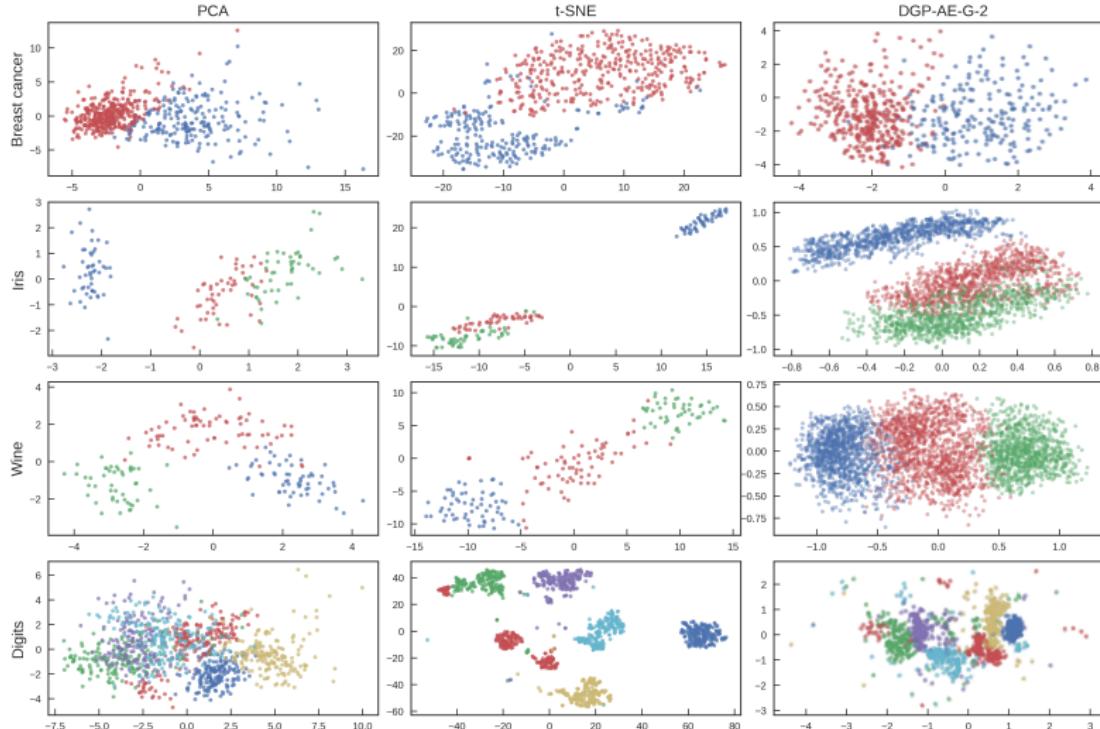
# Convergence monitoring - GPS

- Dimensionality reduction capabilities of a DGP-AE-G-2



- Increasing the number of GPs results in a slower convergence
- 5 GPs achieve good novelty detection performance despite a significant dimensionality reduction

# Latent representation



- Meaningful low-dimensional representations, comparable with state-of-the-art manifold learning methods

## Conclusions

---

# Conclusions

- **Contributions**

- Novel deep probabilistic model for novelty detection
- Competitive with state-of-the-art and DNN-based novelty detection methods
- Good dimensionality reduction abilities
- Tractable and scalable inference through random feature expansions and stochastic variational inference
- Suitable to model mixed-types features

# Conclusions

- **Contributions**
  - Novel deep probabilistic model for novelty detection
  - Competitive with state-of-the-art and DNN-based novelty detection methods
  - Good dimensionality reduction abilities
  - Tractable and scalable inference through random feature expansions and stochastic variational inference
  - Suitable to model mixed-types features
- **Future work**
  - Model discrete event sequences with structured DGP-AE
  - Generative DGP-AE

**Thank you!**